

```
In [57]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings("ignore")

from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score, confusion_matrix, roc_curve, roc_auc_score
from sklearn.decomposition import PCA

from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
from sklearn.svm import SVC
```

```
In [58]: df=pd.read_csv((r'C:\Users\mmi\Desktop\cancer.csv'))
```

```
In [59]: df.head(5)
df.drop('id',axis=1,inplace=True)
df.drop('Unnamed: 32',axis=1,inplace=True)
df
```

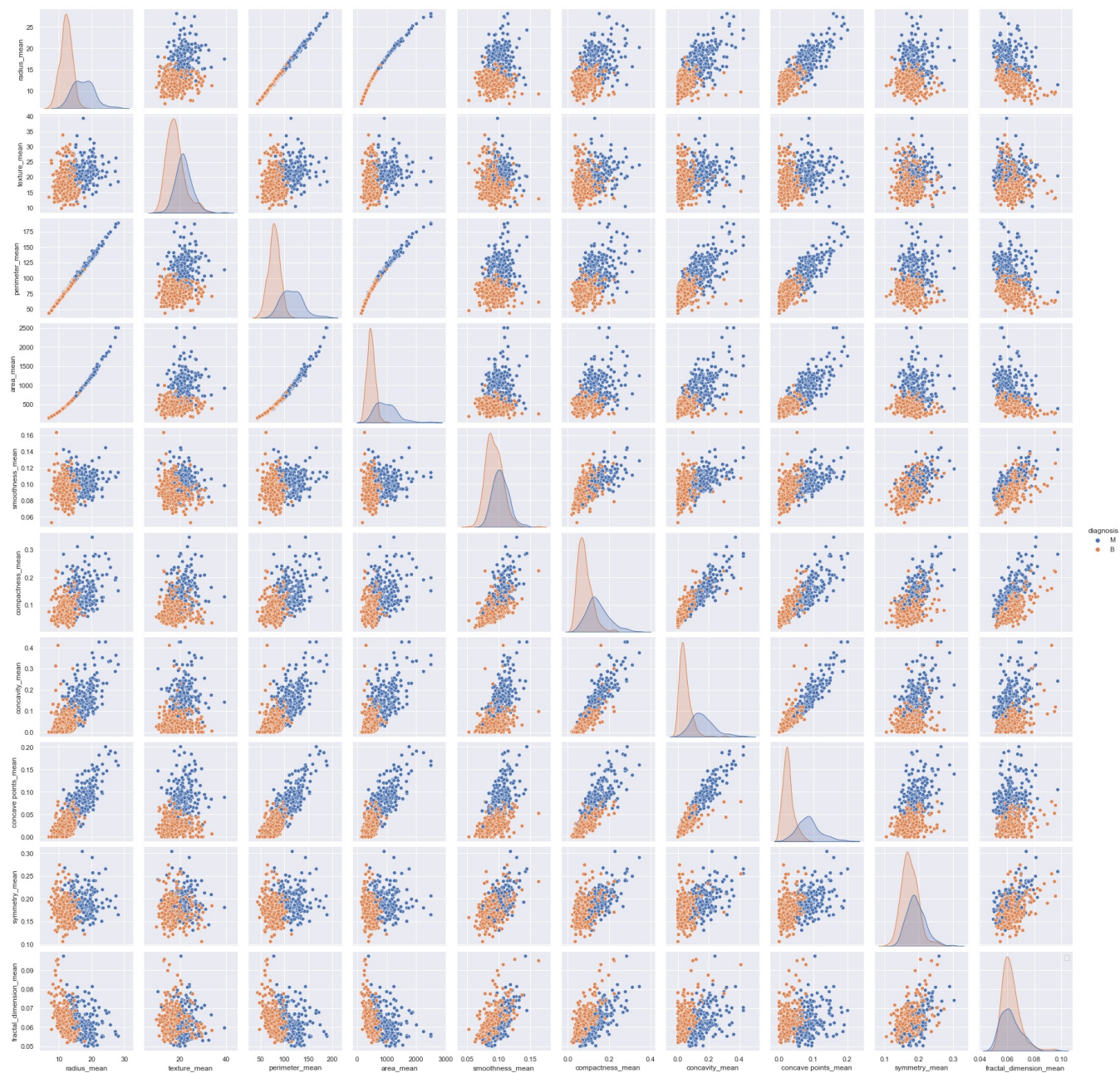
```
Out[59]:
```

	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	conca points_me
0	M	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.30010	0.147
1	M	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.08690	0.070
2	M	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.19740	0.127
3	M	11.42	20.38	77.58	386.1	0.14250	0.28390	0.24140	0.105
4	M	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.19800	0.104
...
564	M	21.56	22.39	142.00	1479.0	0.11100	0.11590	0.24390	0.138
565	M	20.13	28.25	131.20	1261.0	0.09780	0.10340	0.14400	0.097
566	M	16.60	28.08	108.30	858.1	0.08455	0.10230	0.09251	0.053
567	M	20.60	29.33	140.10	1265.0	0.11780	0.27700	0.35140	0.152
568	B	7.76	24.54	47.92	181.0	0.05263	0.04362	0.00000	0.000

569 rows × 31 columns

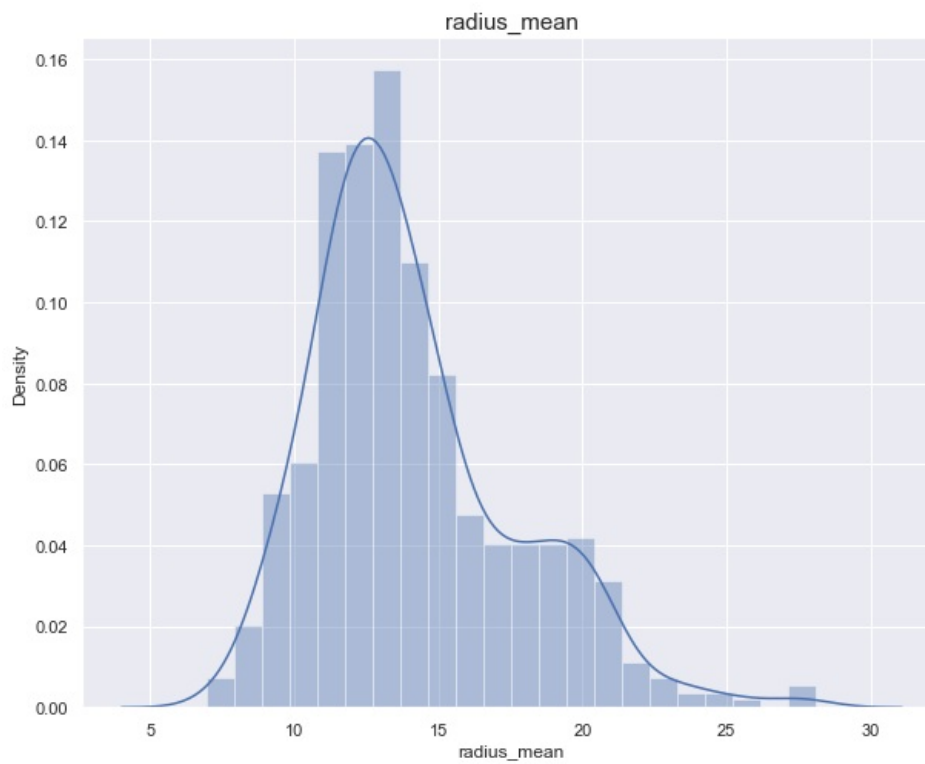
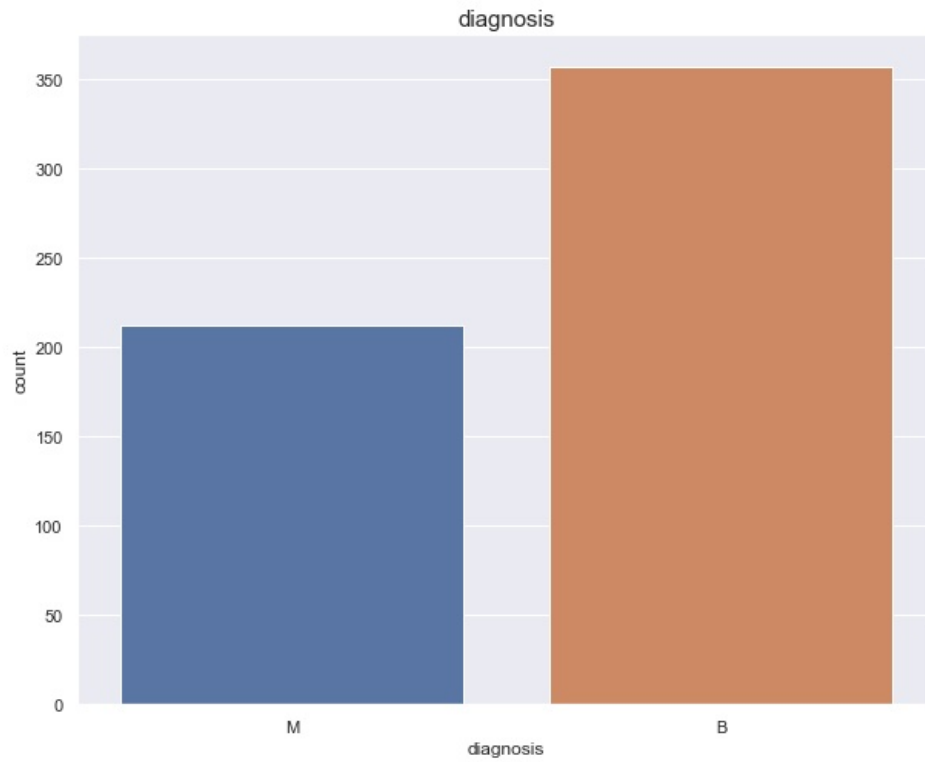
```
In [60]: sns.set()
cols_to_pairplot = df.columns[:11]
sns.pairplot(df[cols_to_pairplot], hue="diagnosis")
plt.legend()
plt.show()
```

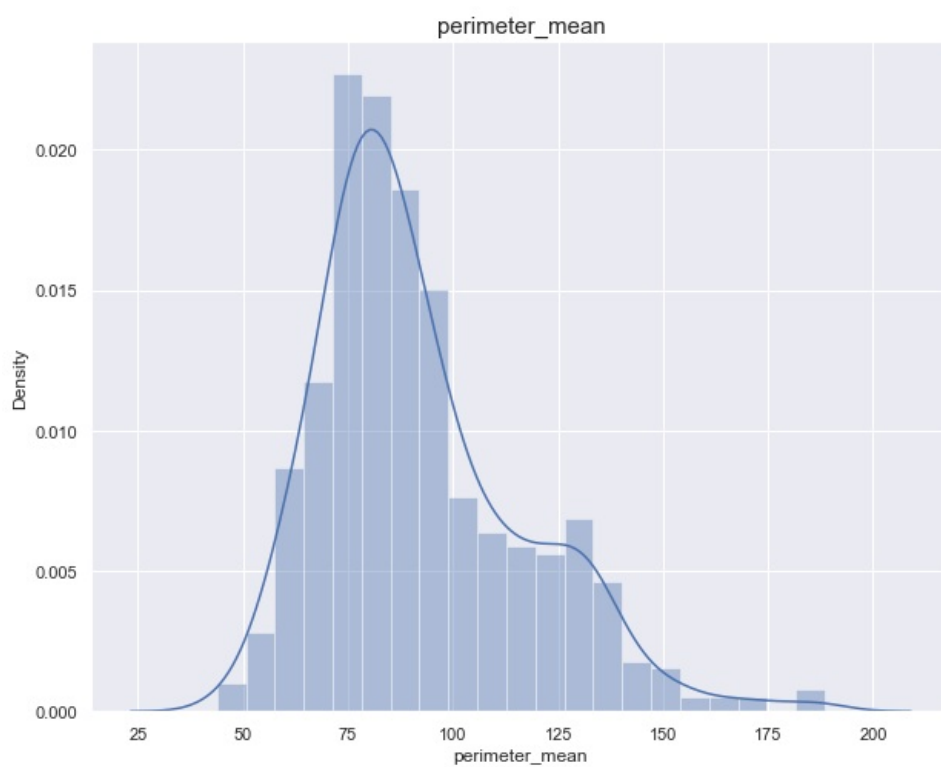
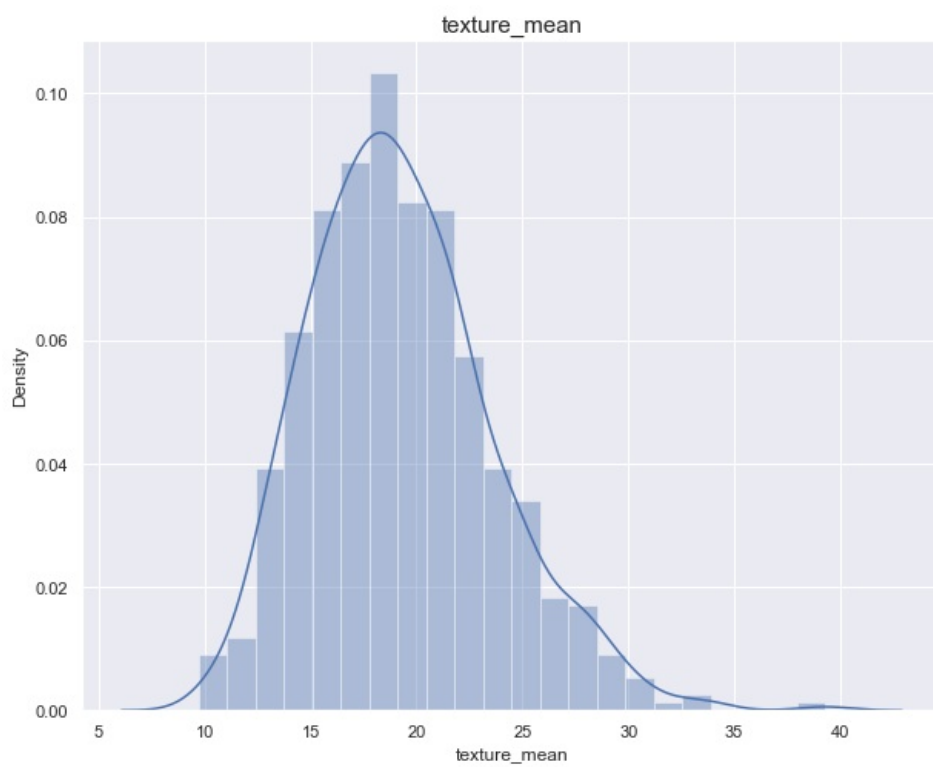
No artists with labels found to put in legend. Note that artists whose label start with an underscore are ignored when legend() is called with no argument.

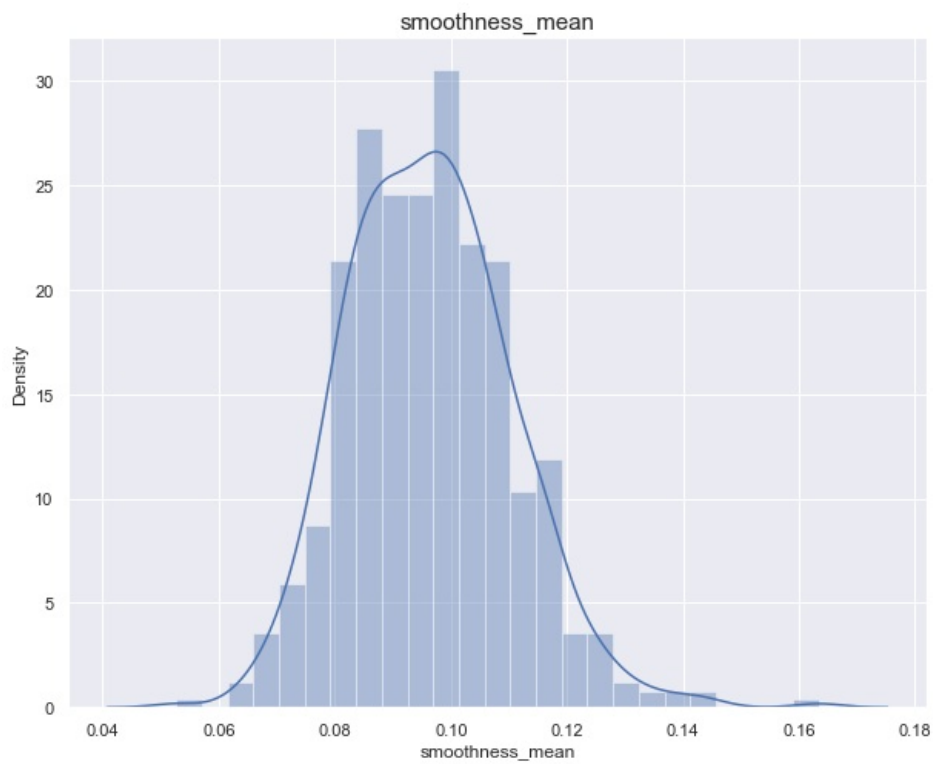
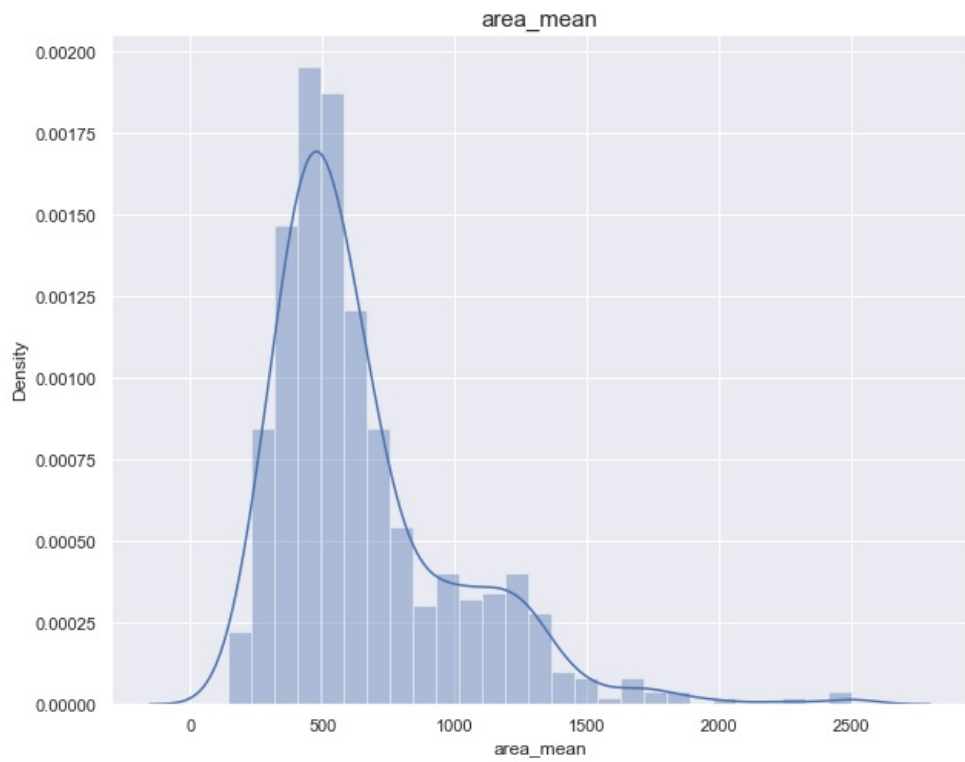


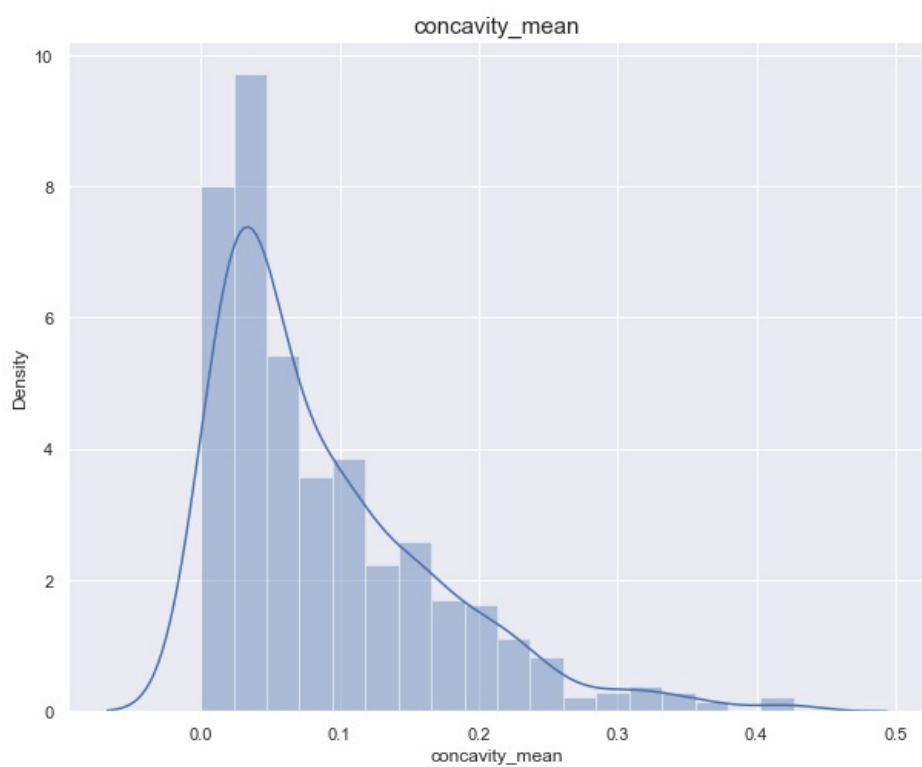
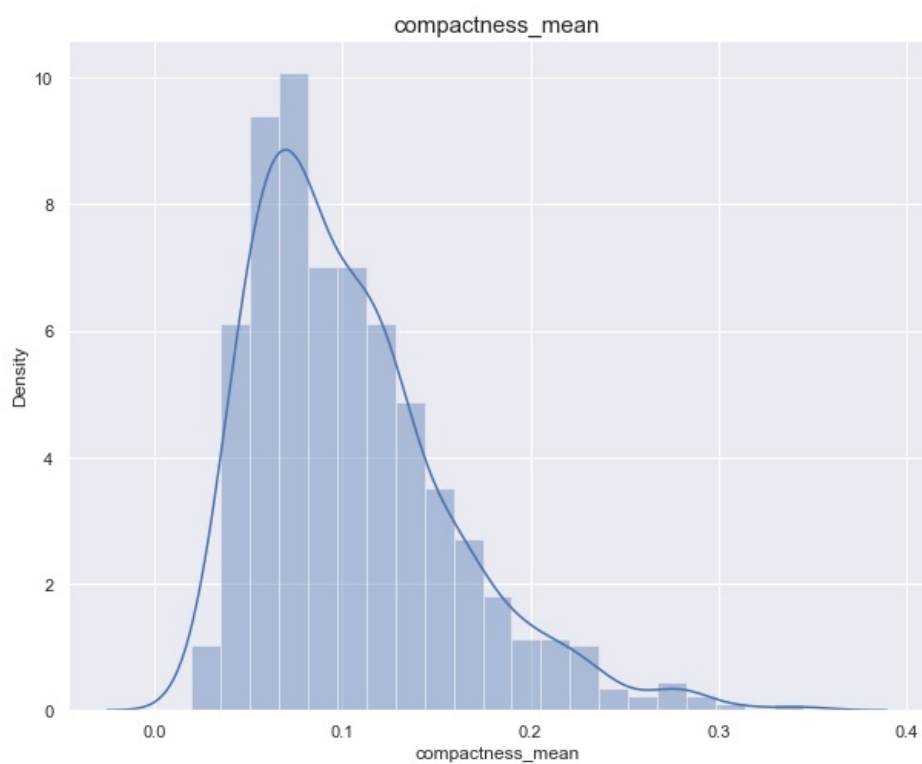
```
In [61]: plt.figure(figsize=(10,8))
sns.countplot(df["diagnosis"])
plt.title("diagnosis", size=15)
plt.show()
```

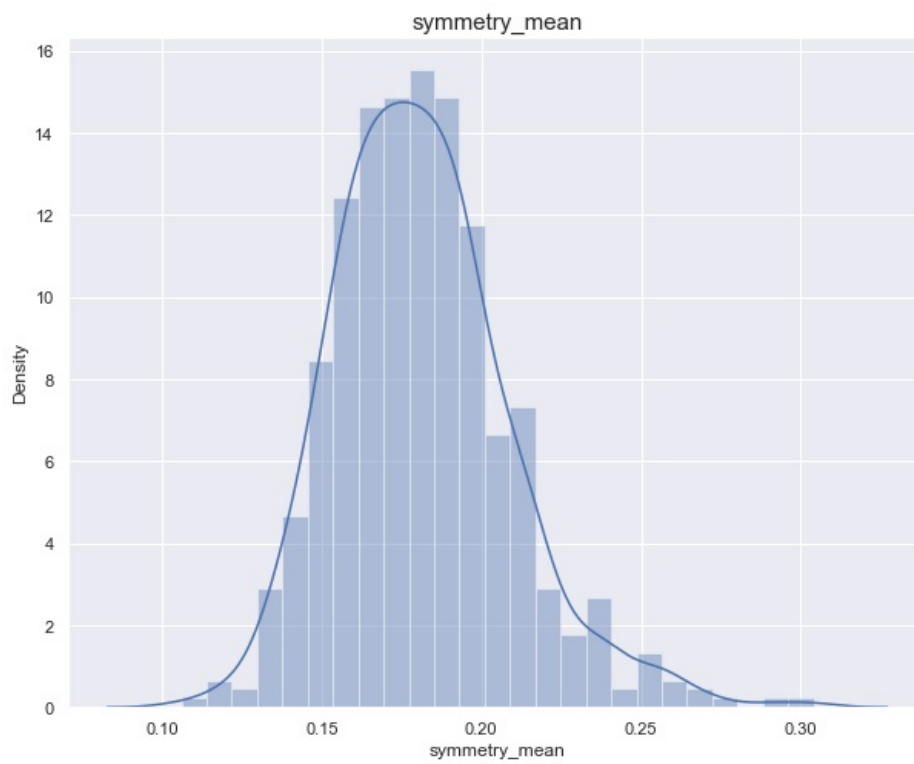
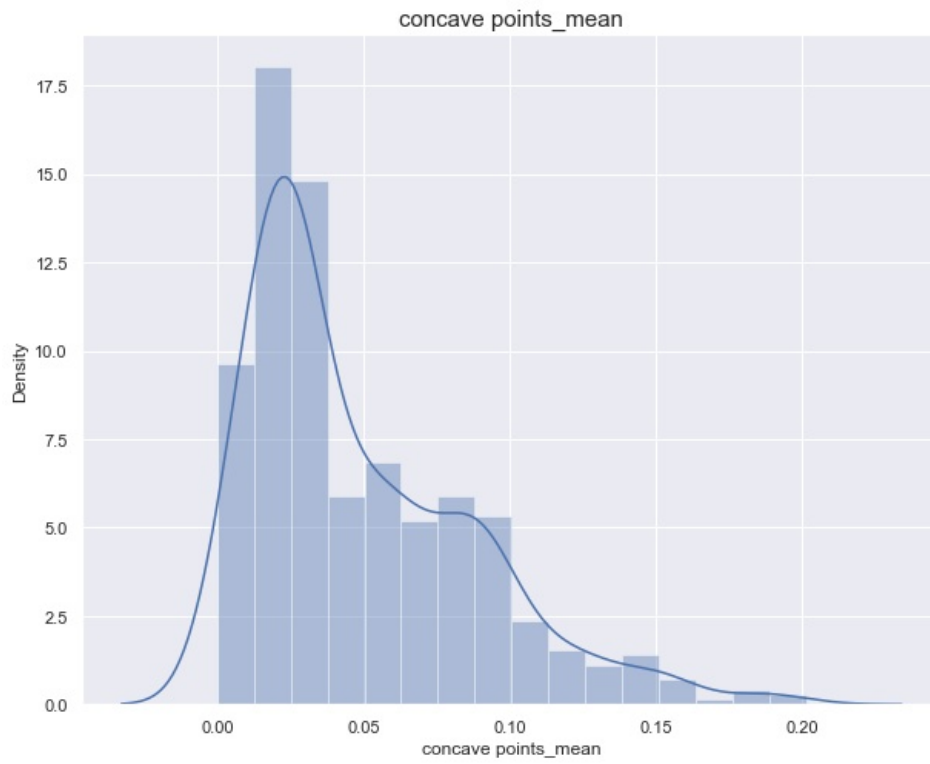
```
for col in df.drop("diagnosis", axis=1).columns:
    plt.figure(figsize=(10,8))
    sns.distplot(df[col])
    plt.title(f"{col}", size=15)
    plt.show()
```

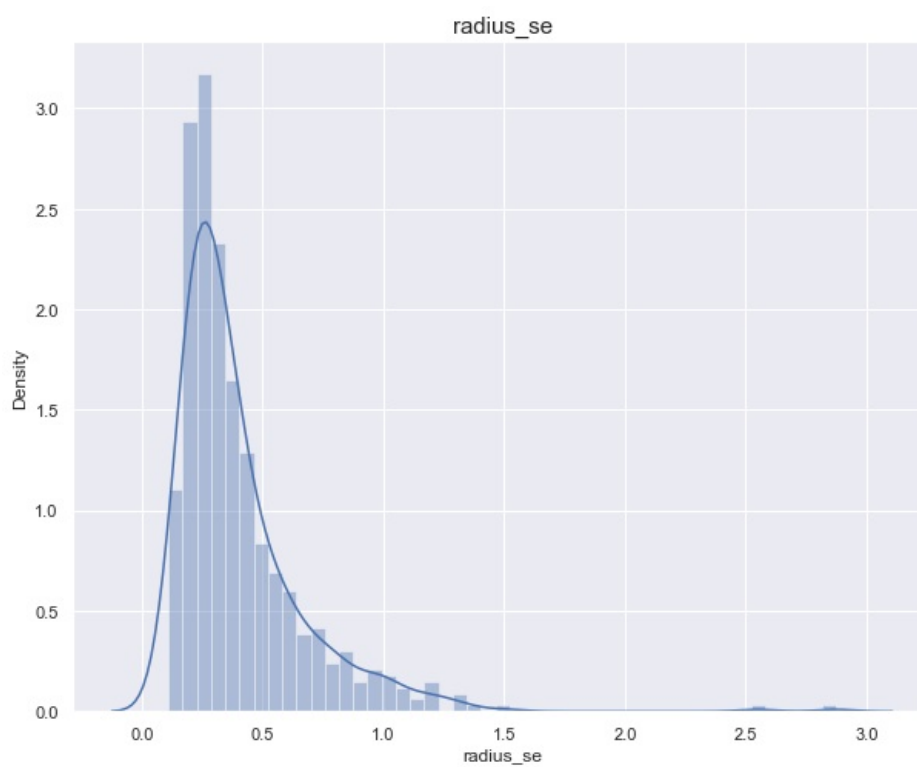
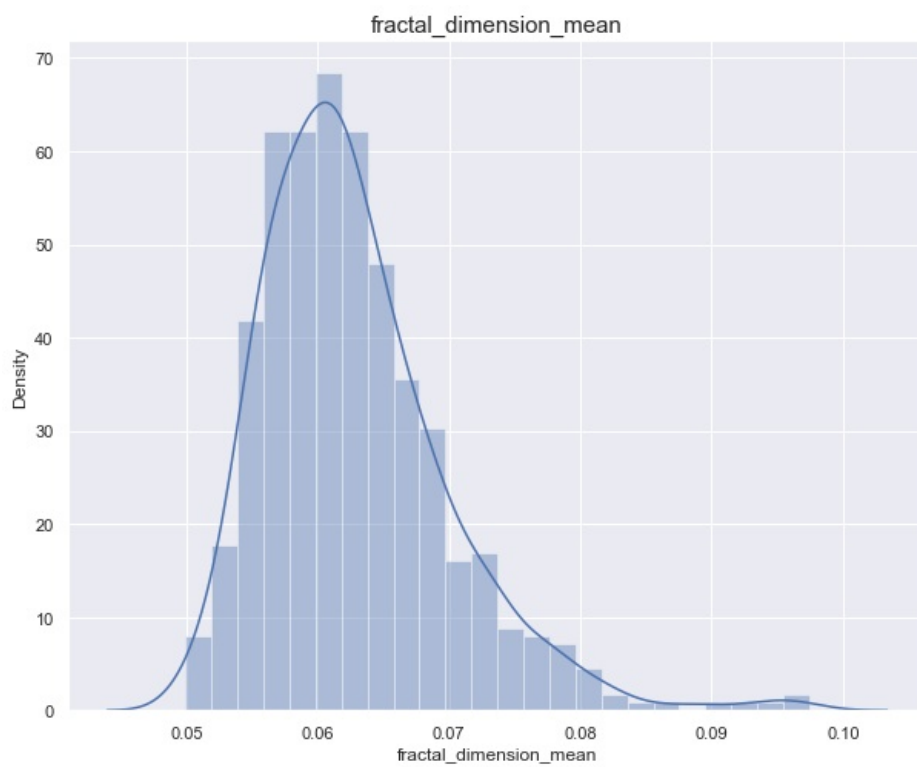


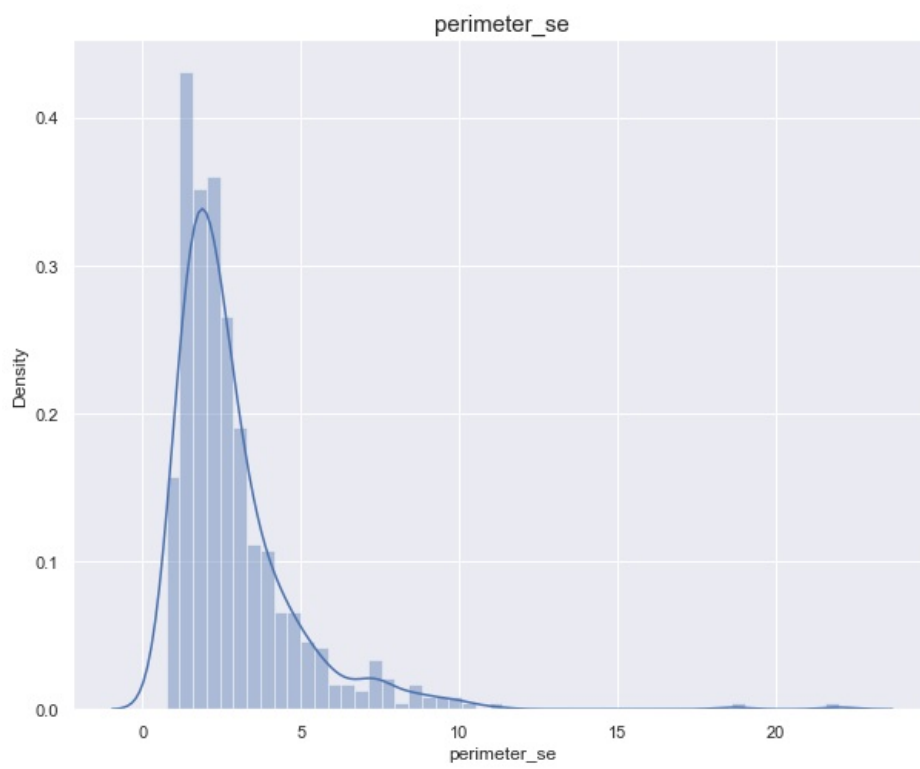
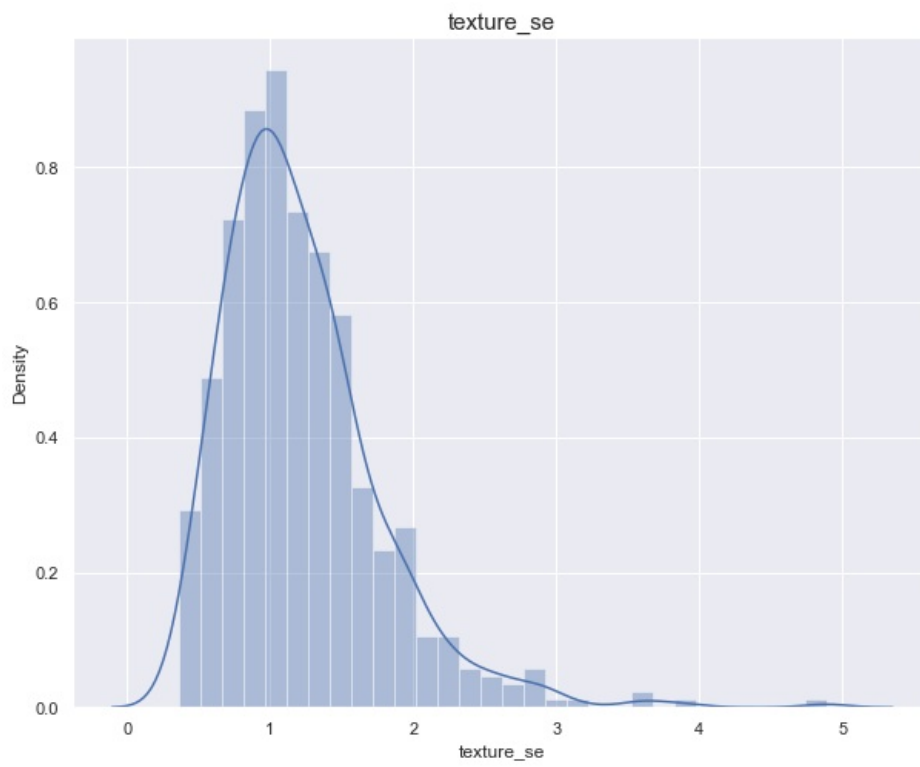


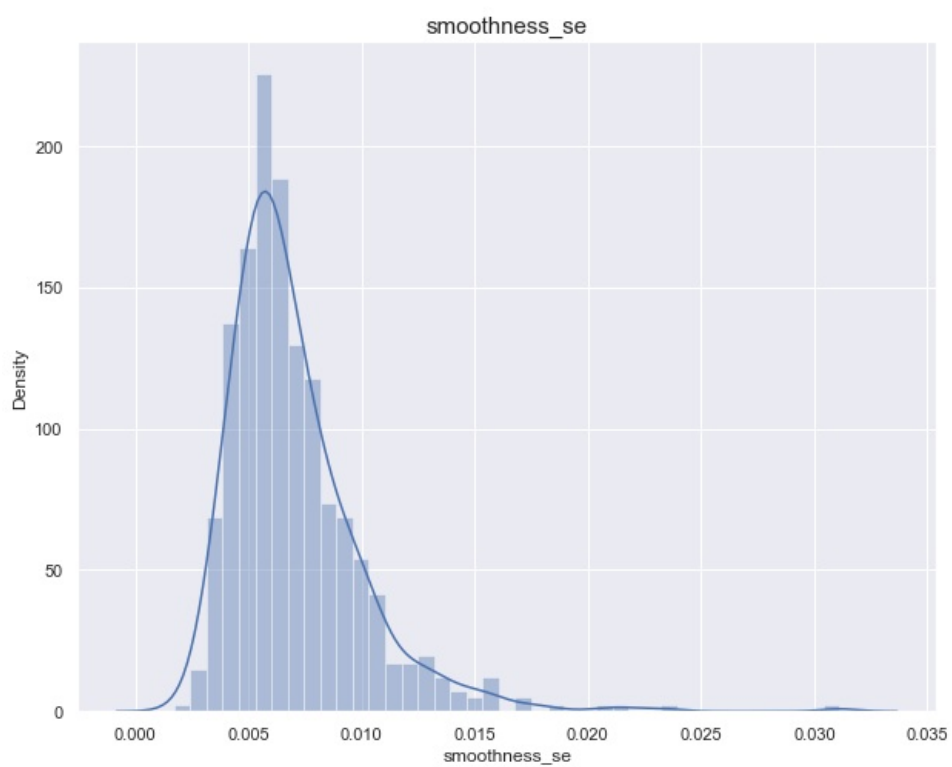
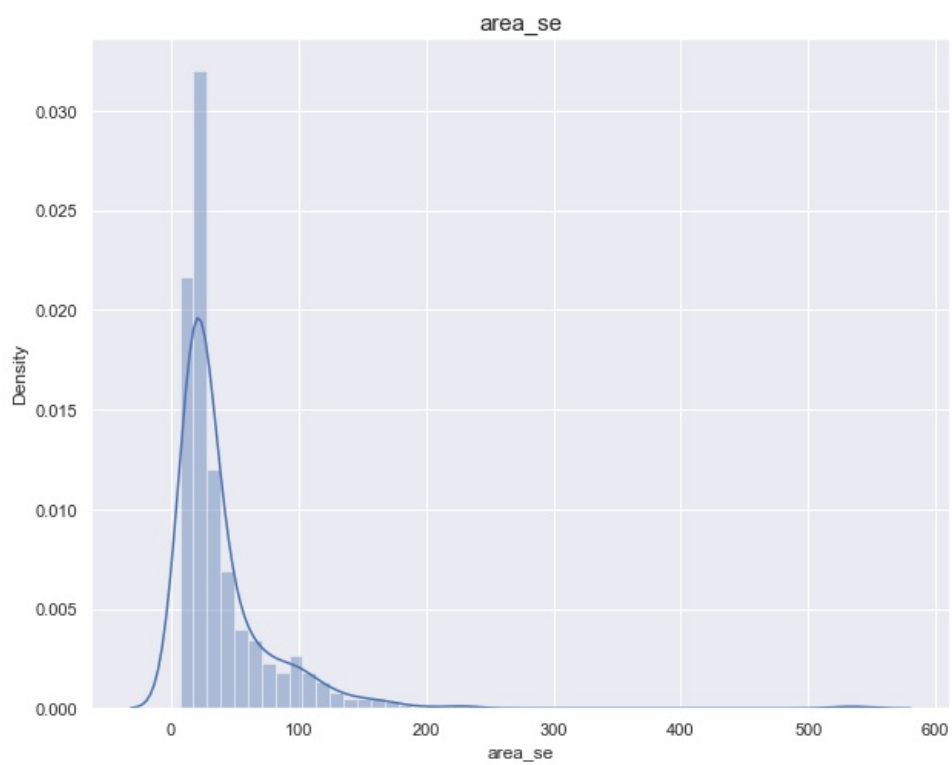


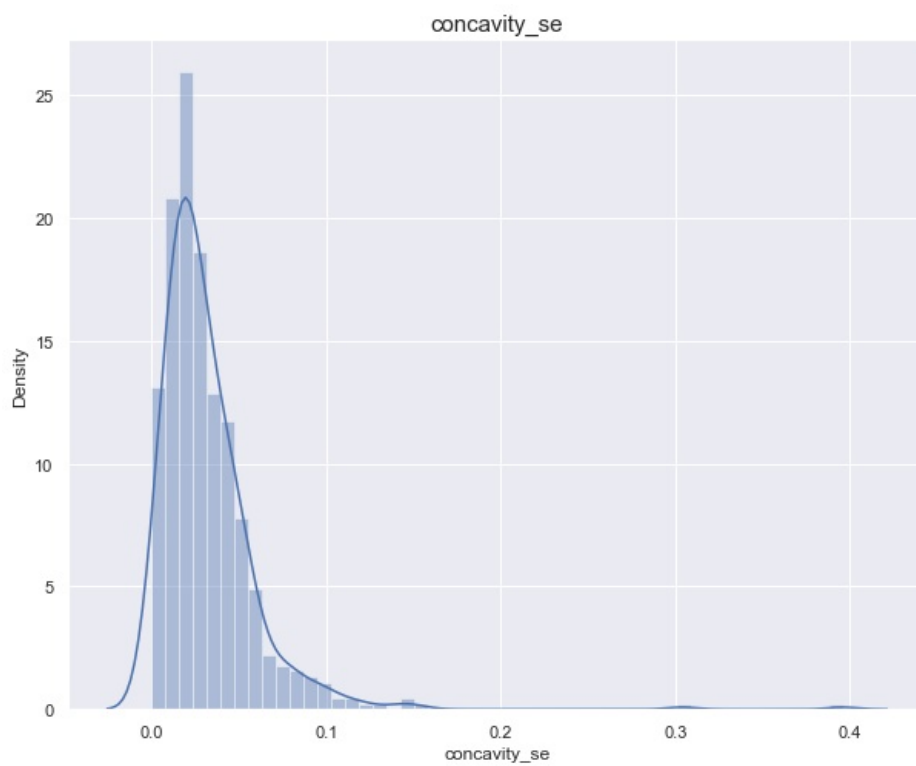
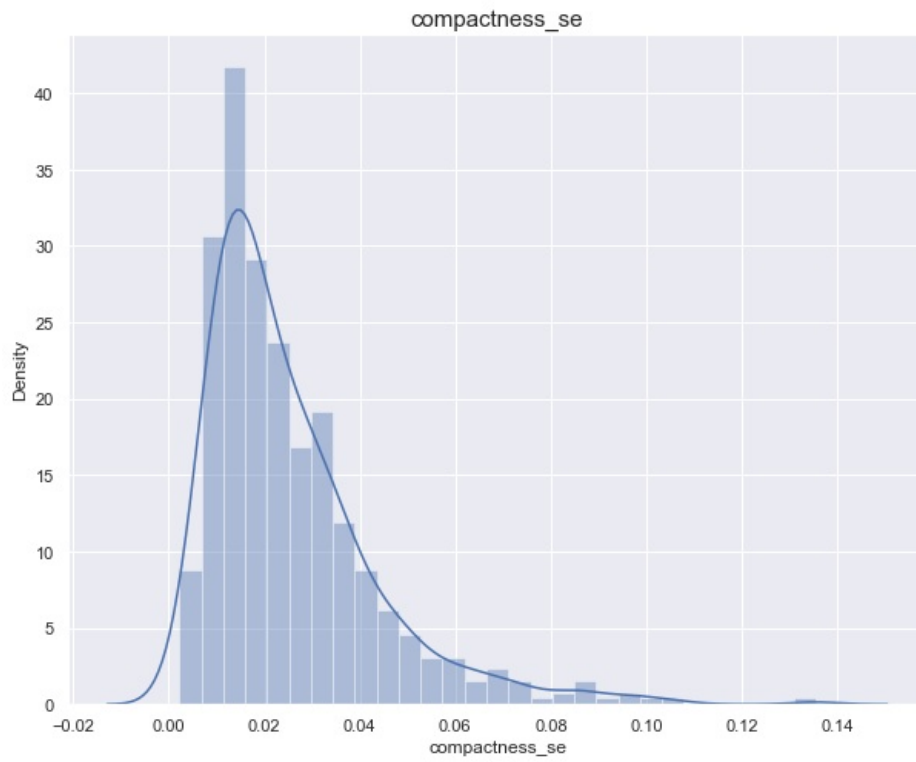


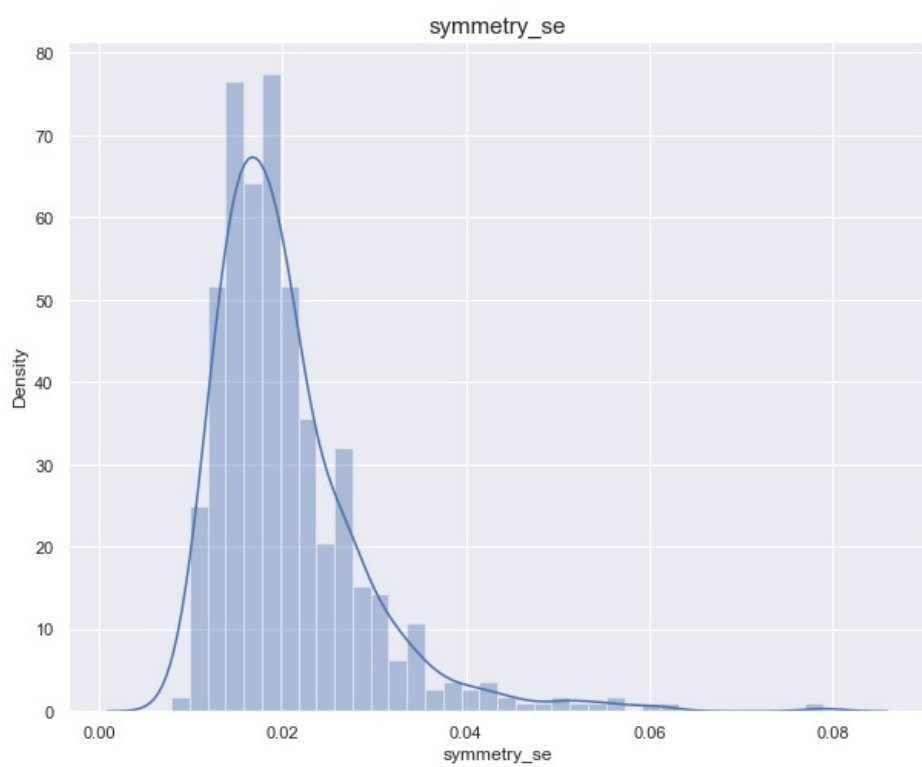
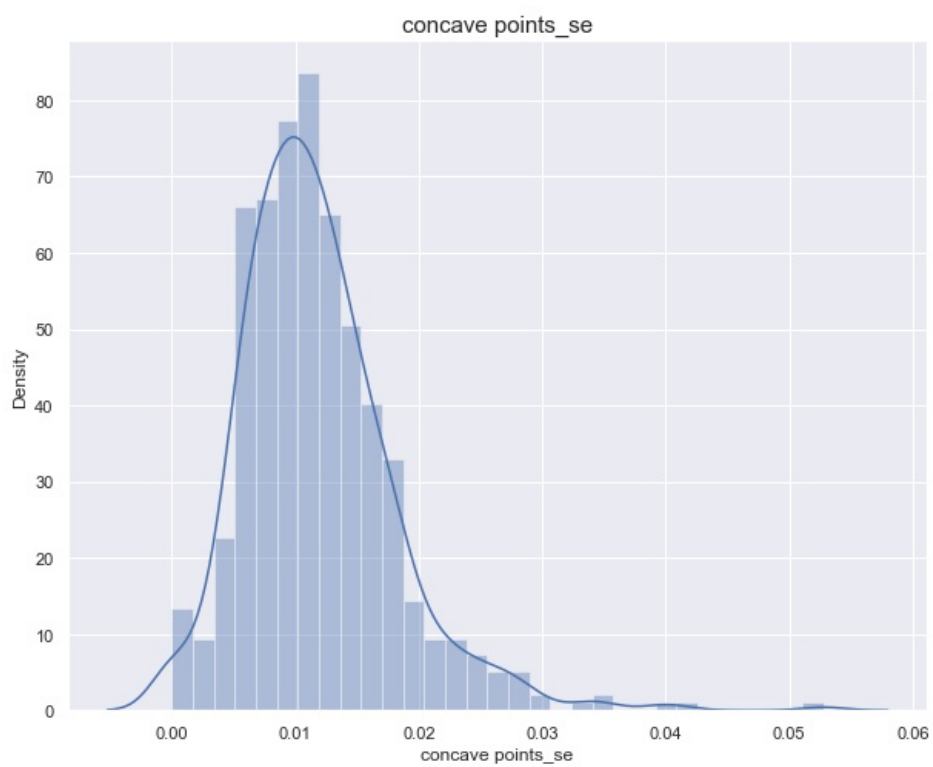


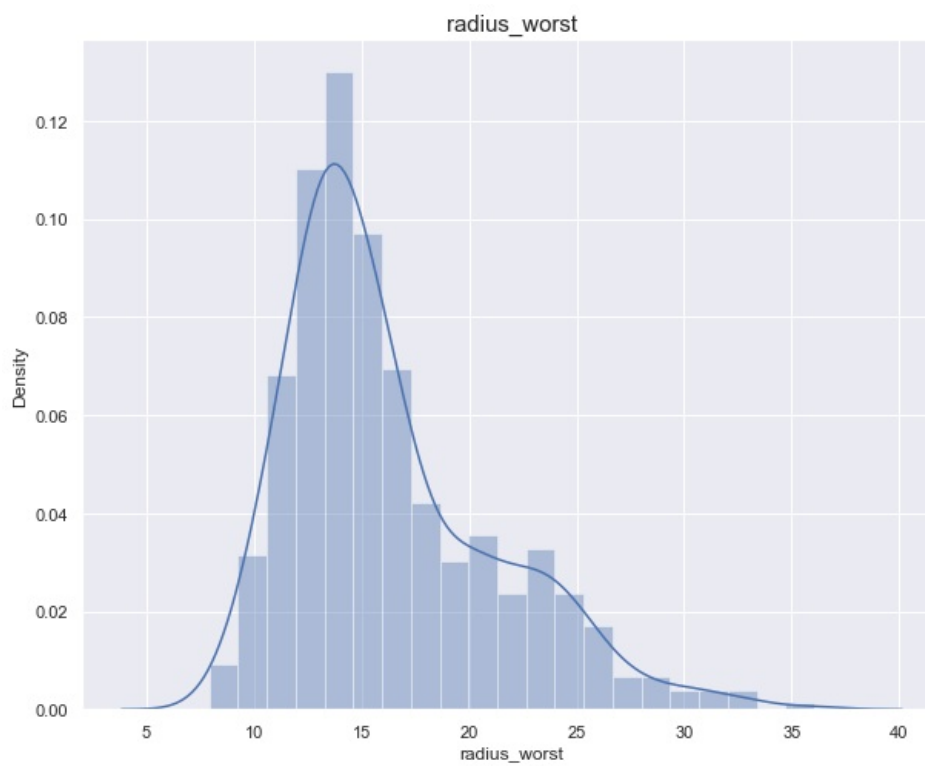
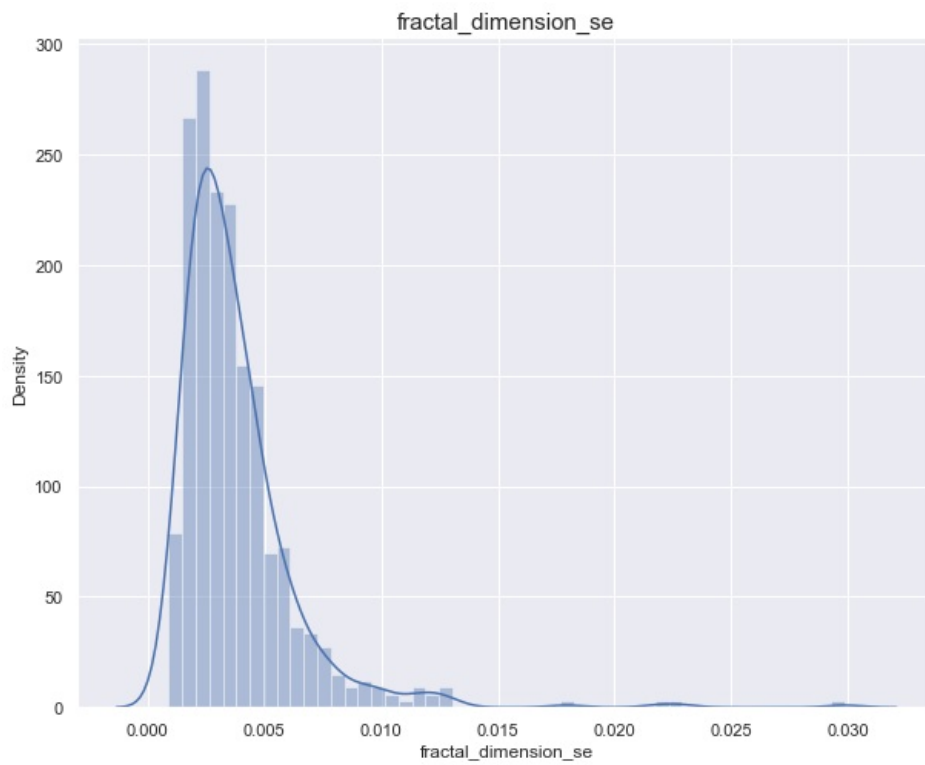


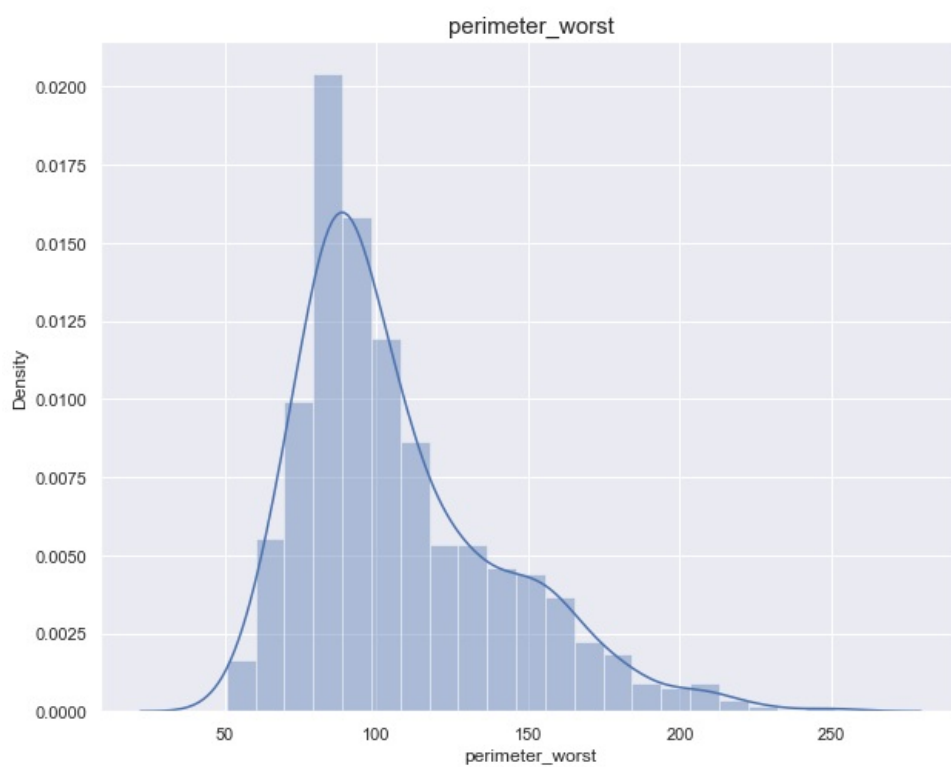
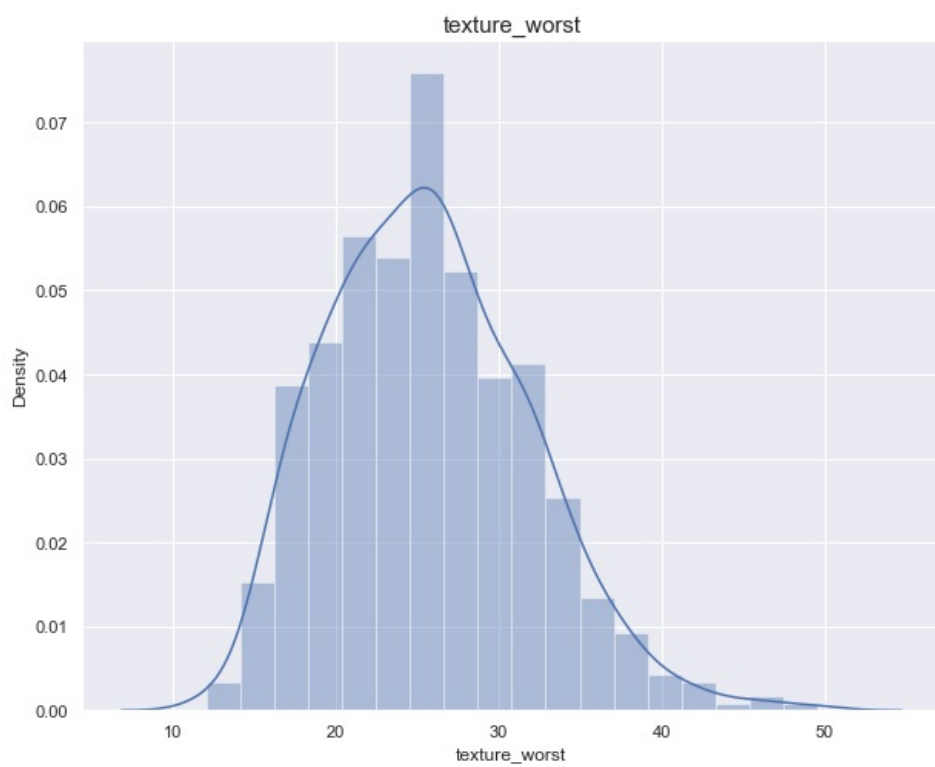


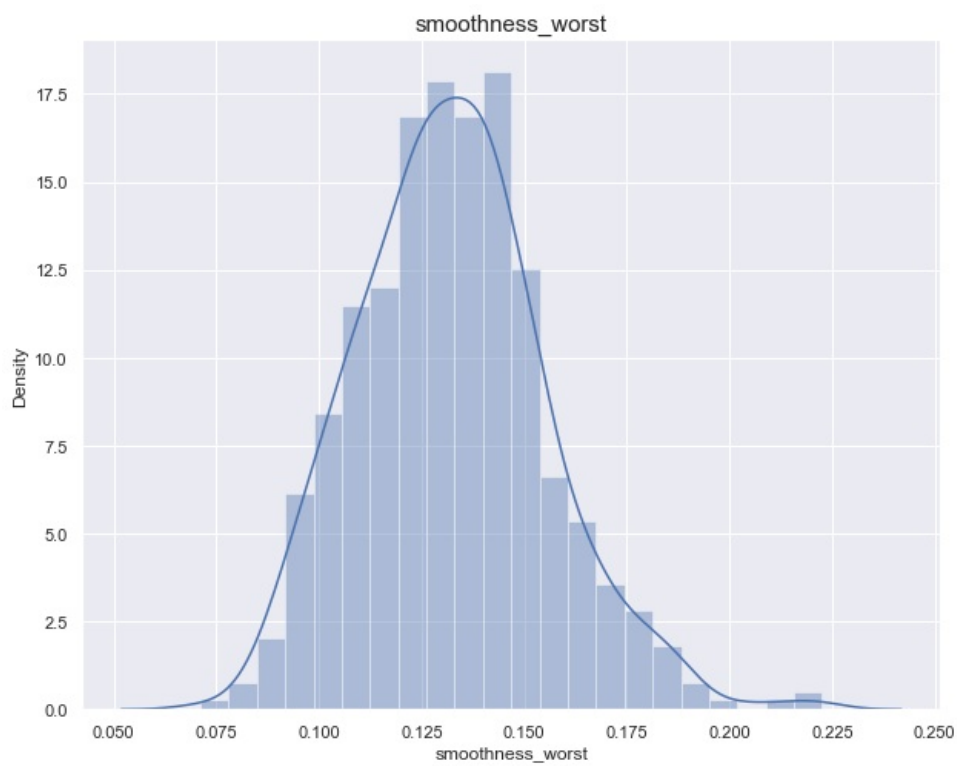
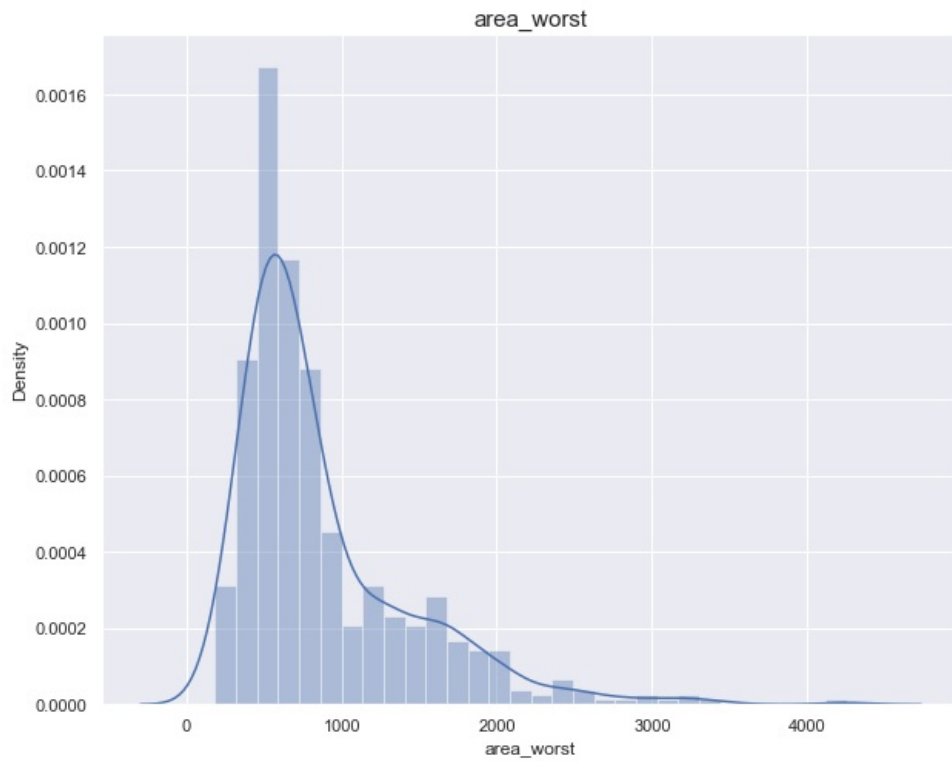


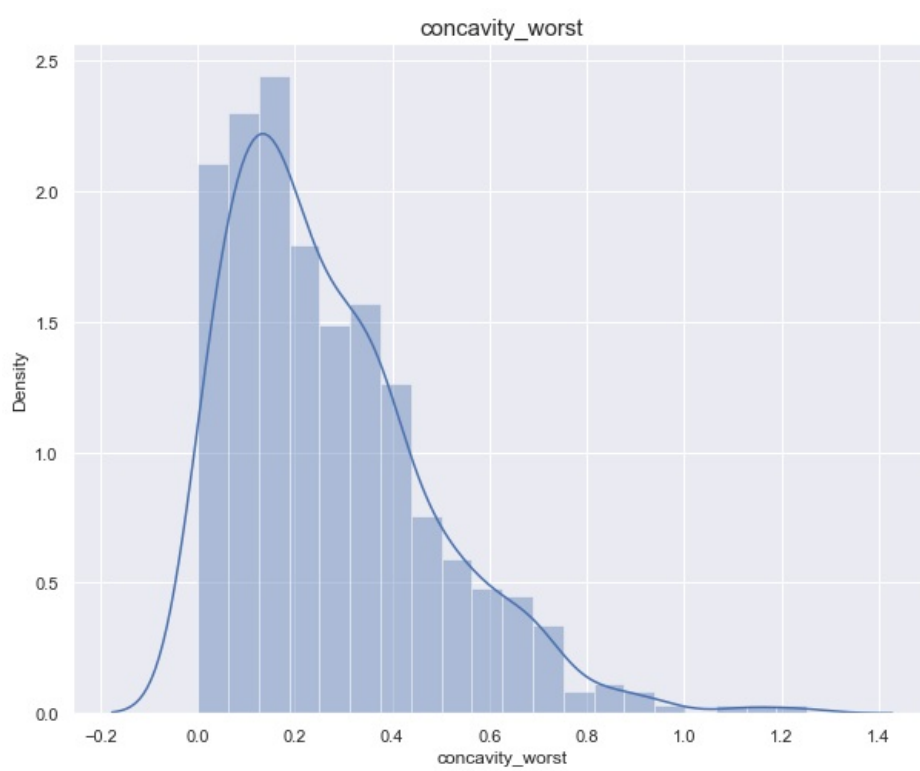
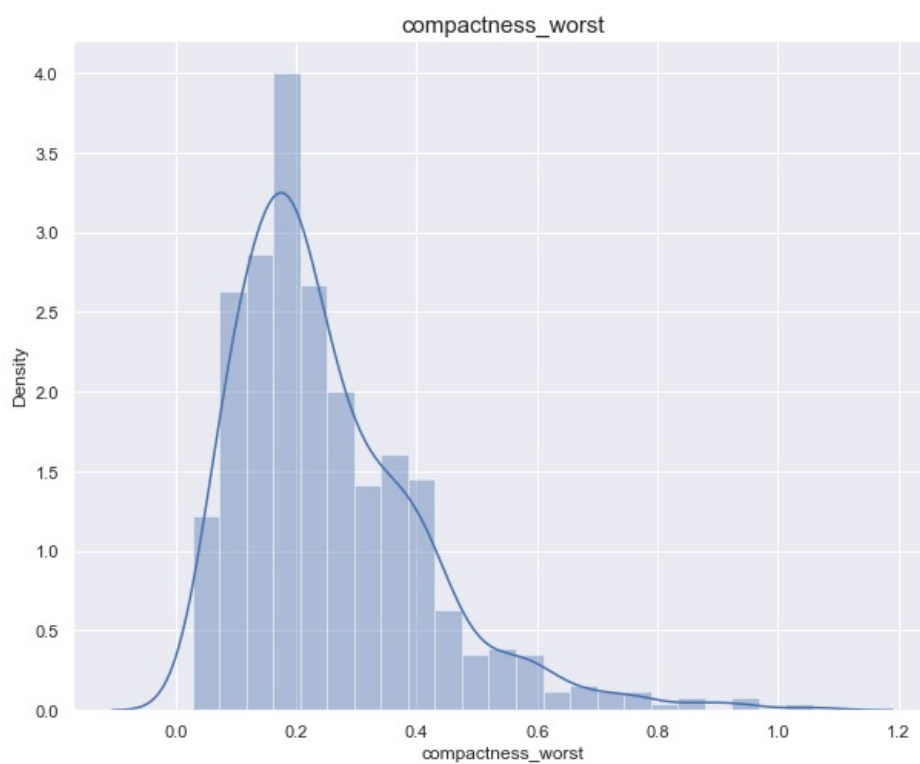


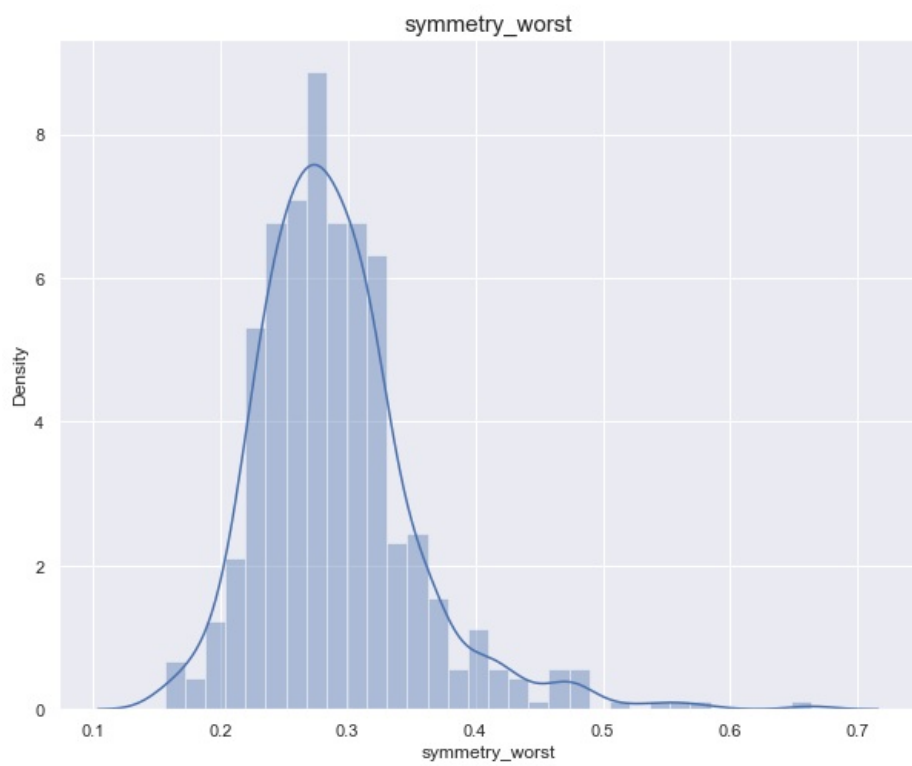
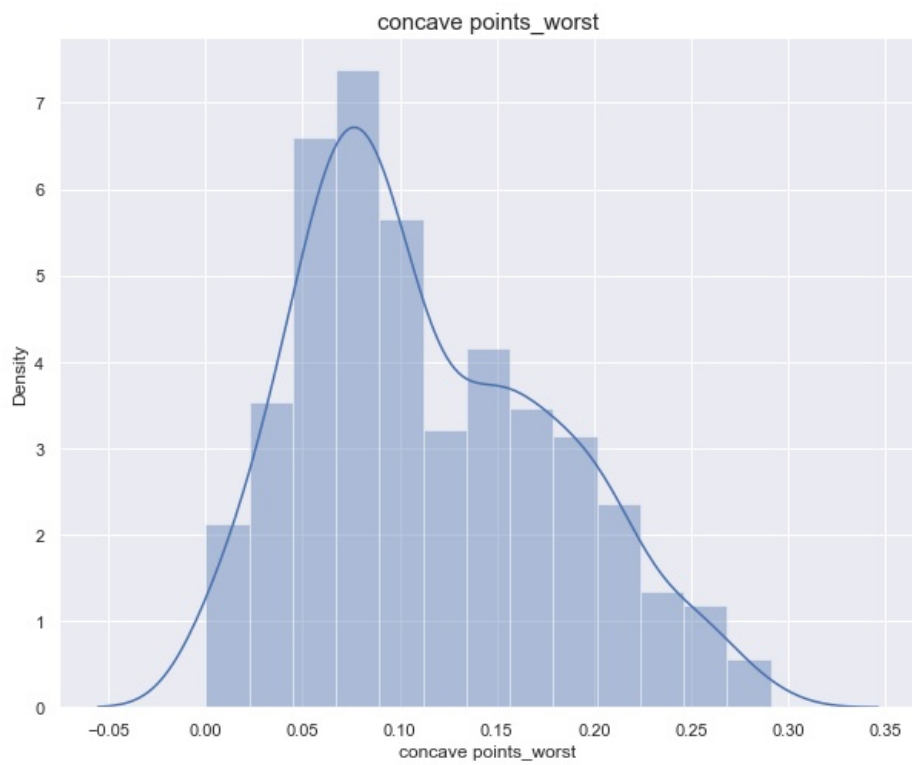


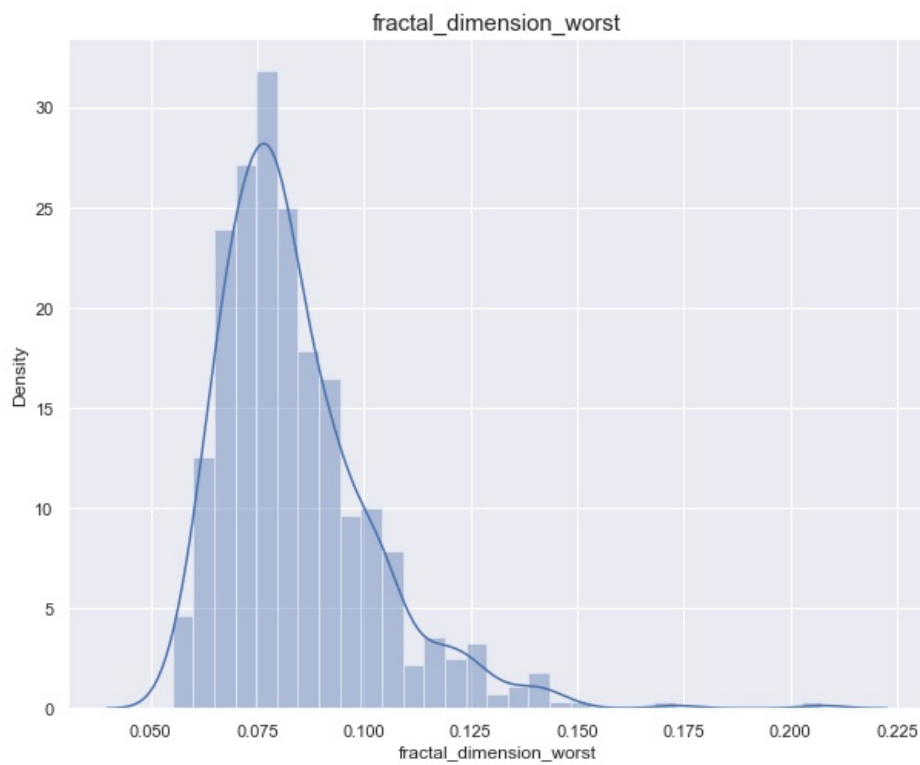




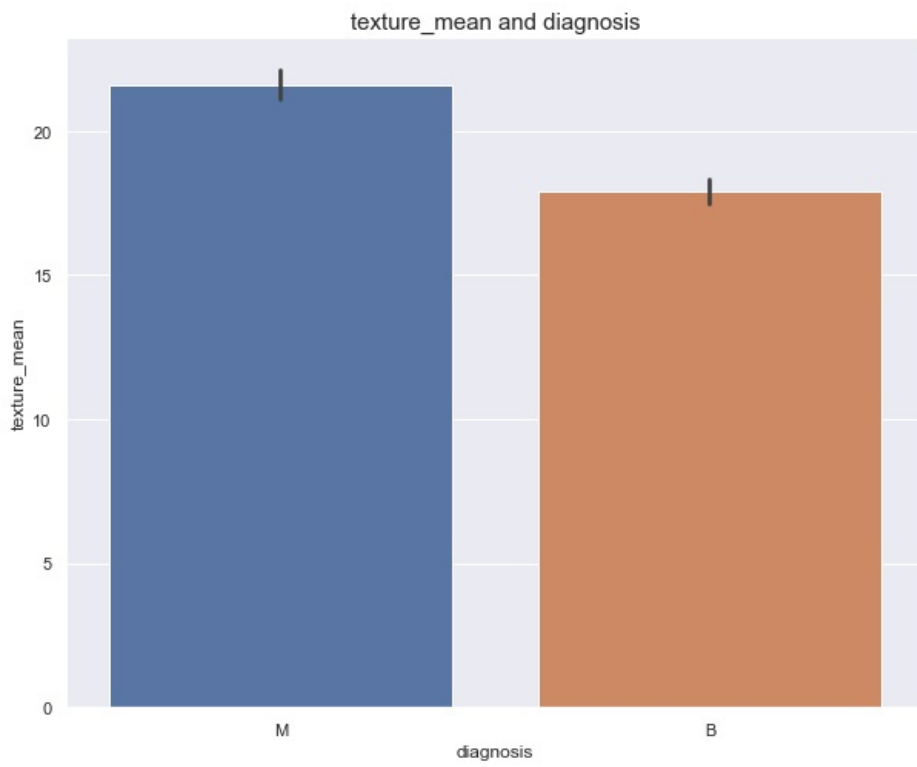
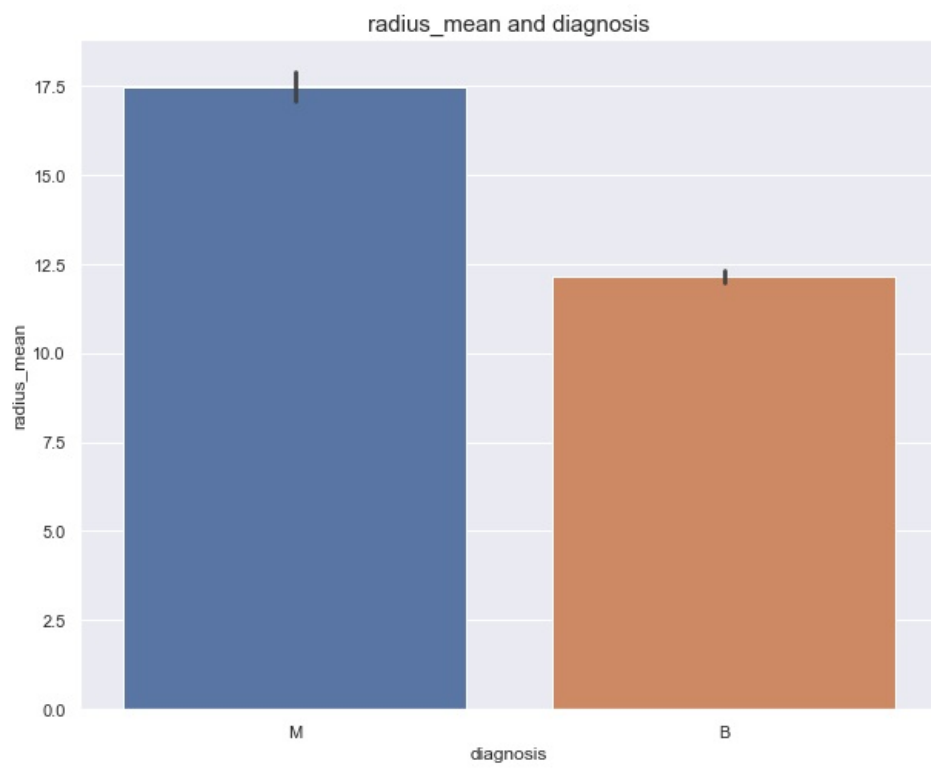


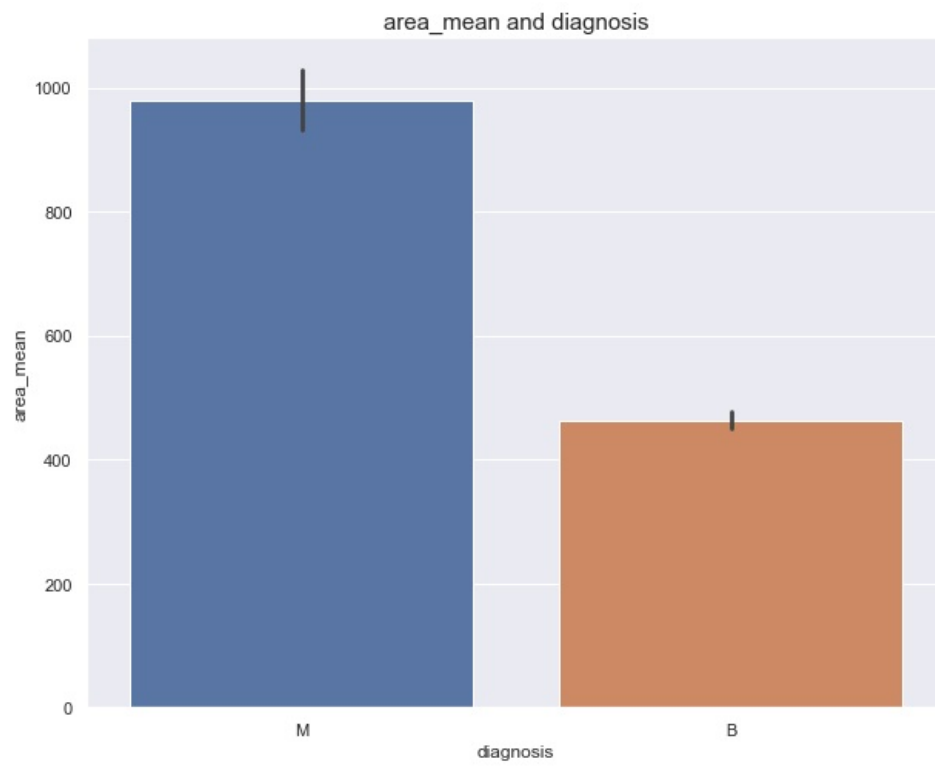
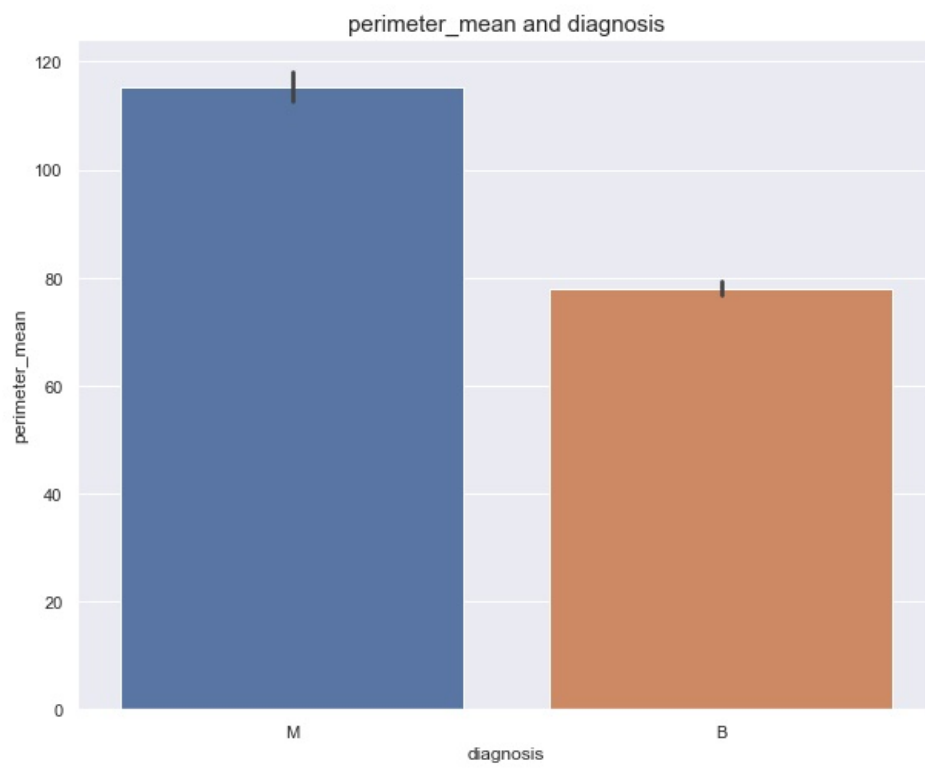


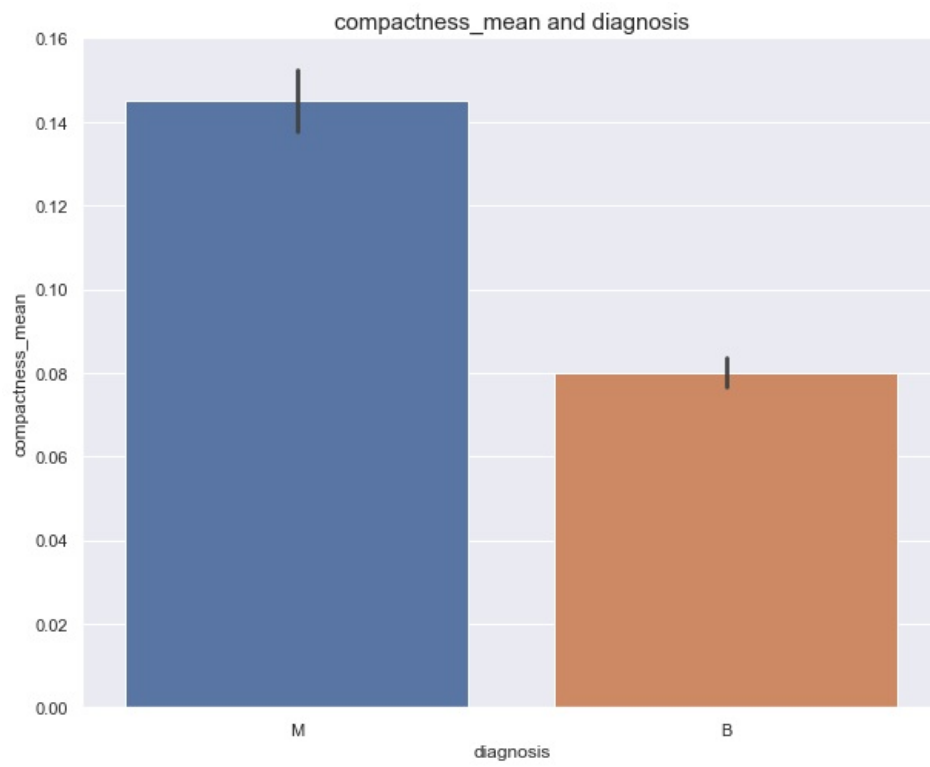
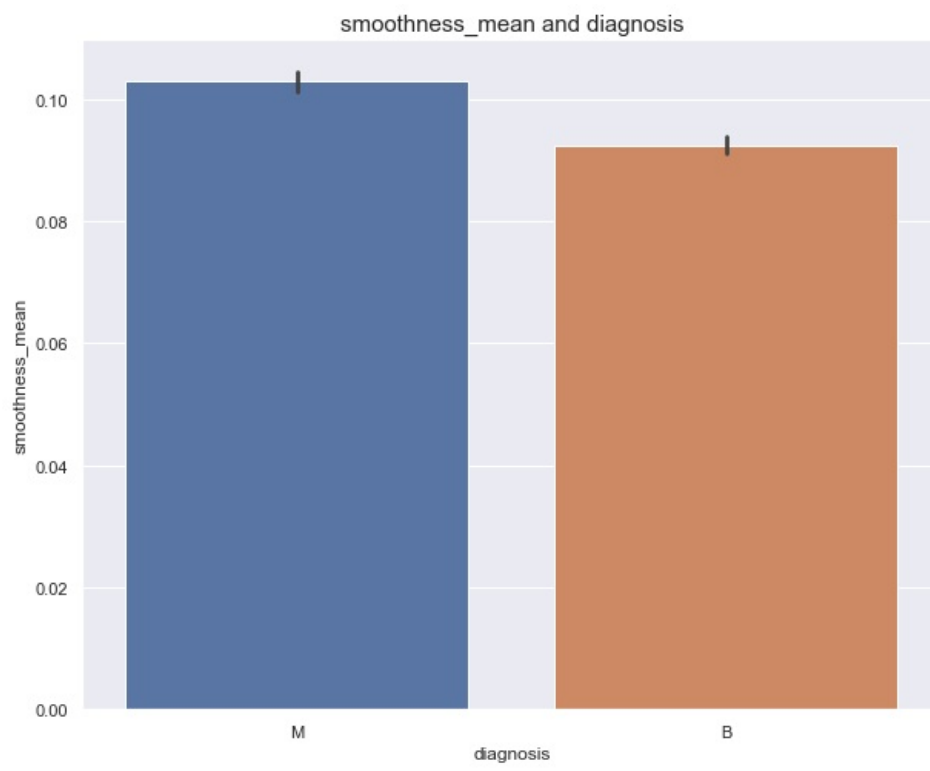


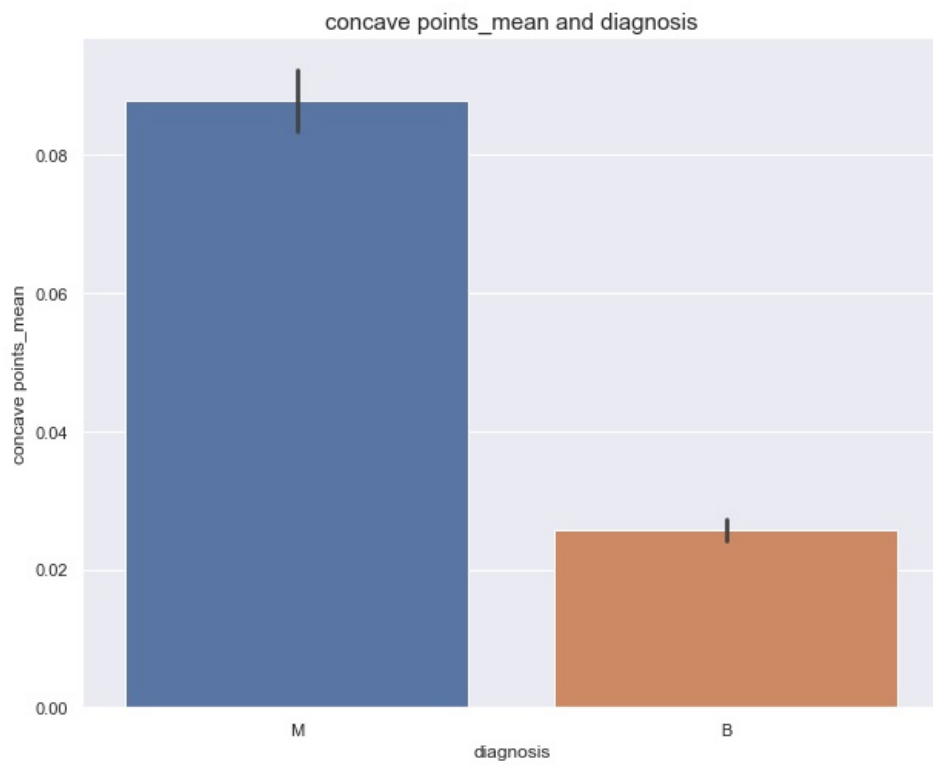
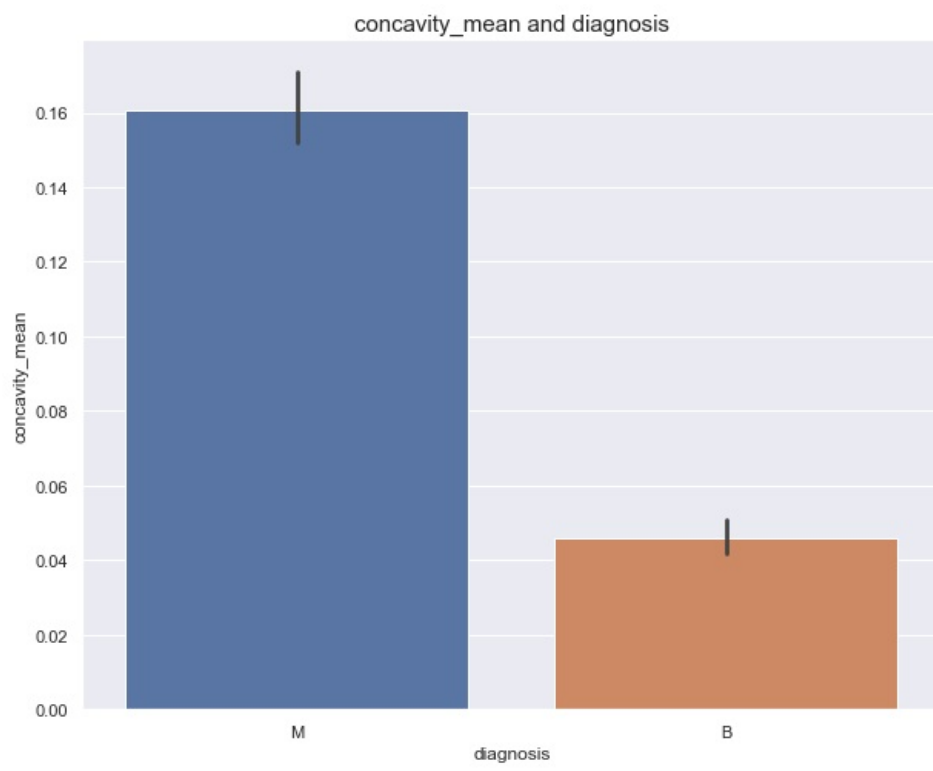


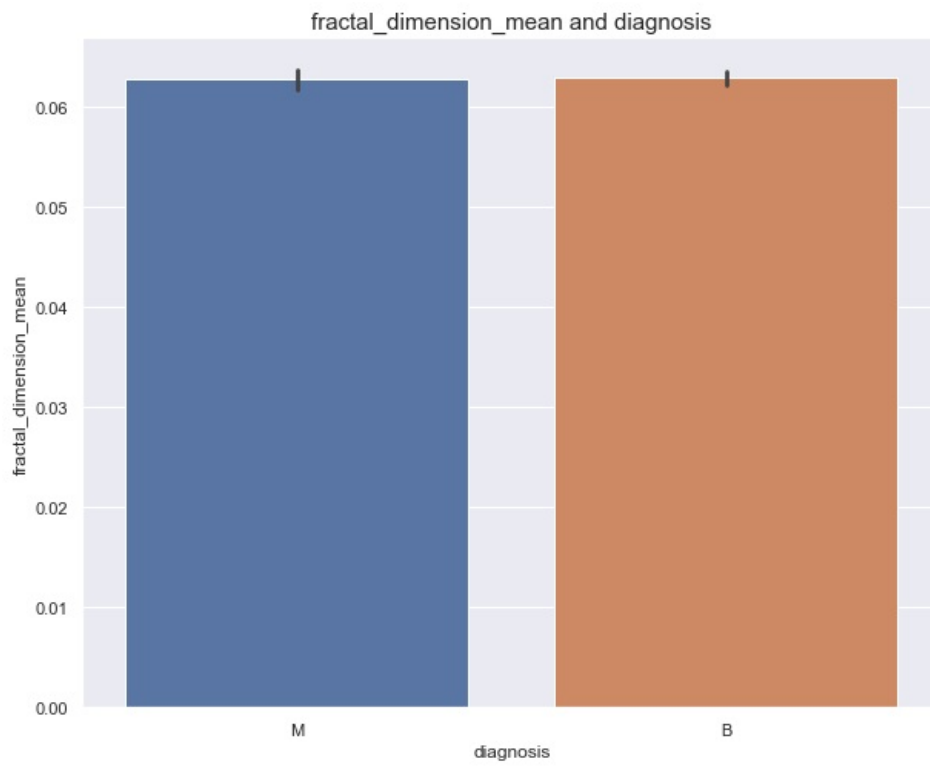
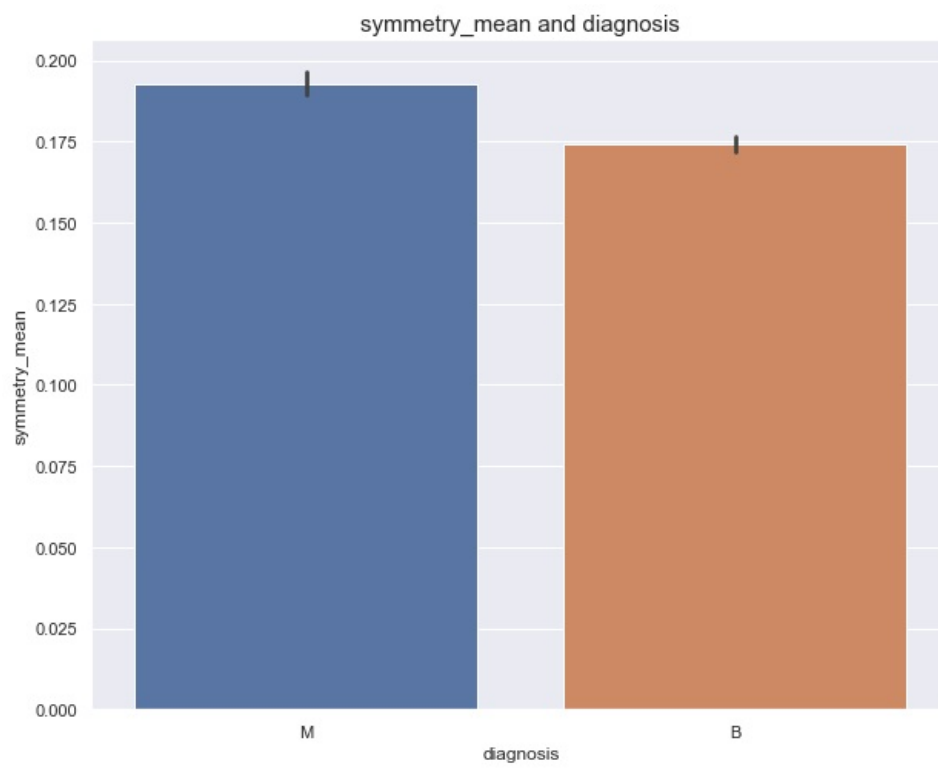
```
In [62]: for col in df.drop("diagnosis", axis=1).columns:
plt.figure(figsize=(10,8))
sns.barplot(x=df["diagnosis"], y=df[col])
plt.title(f"{col} and diagnosis", size=15)
plt.show()
```

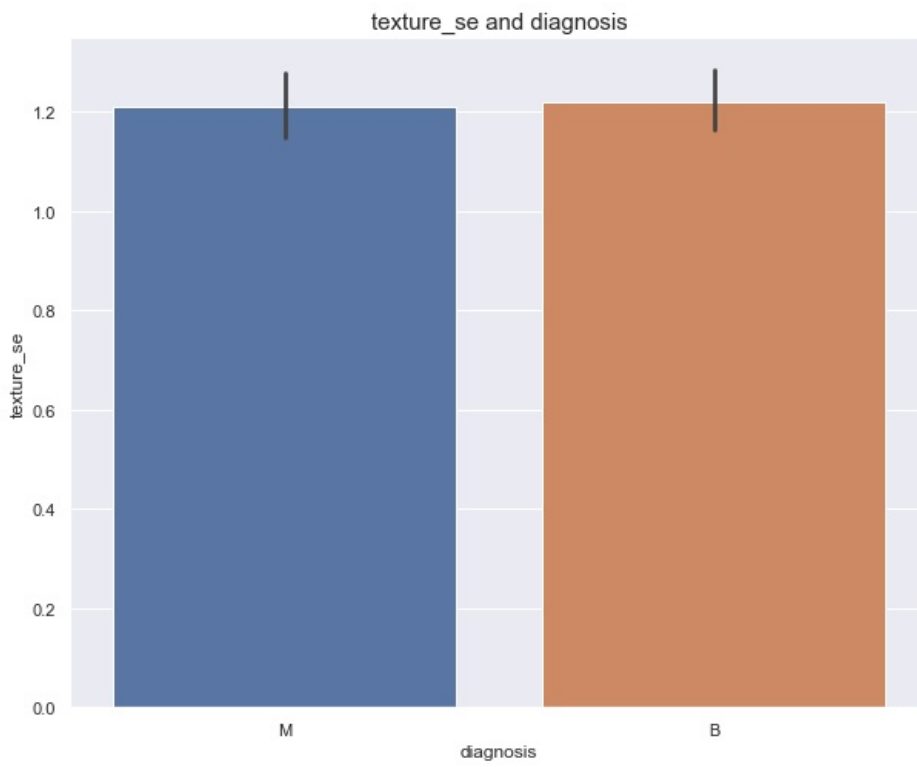
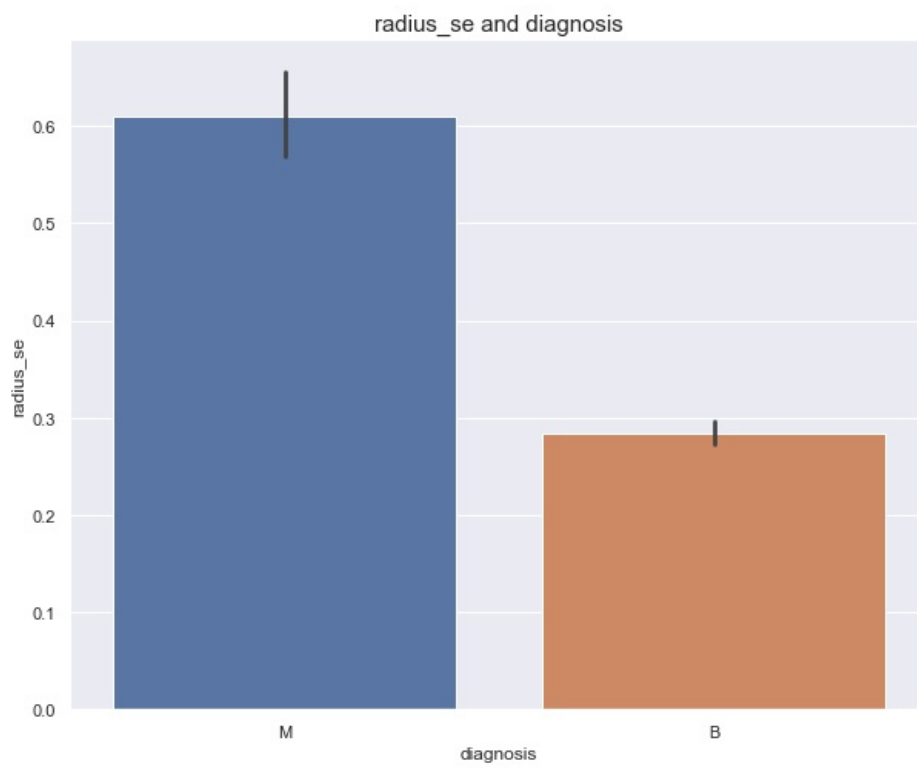




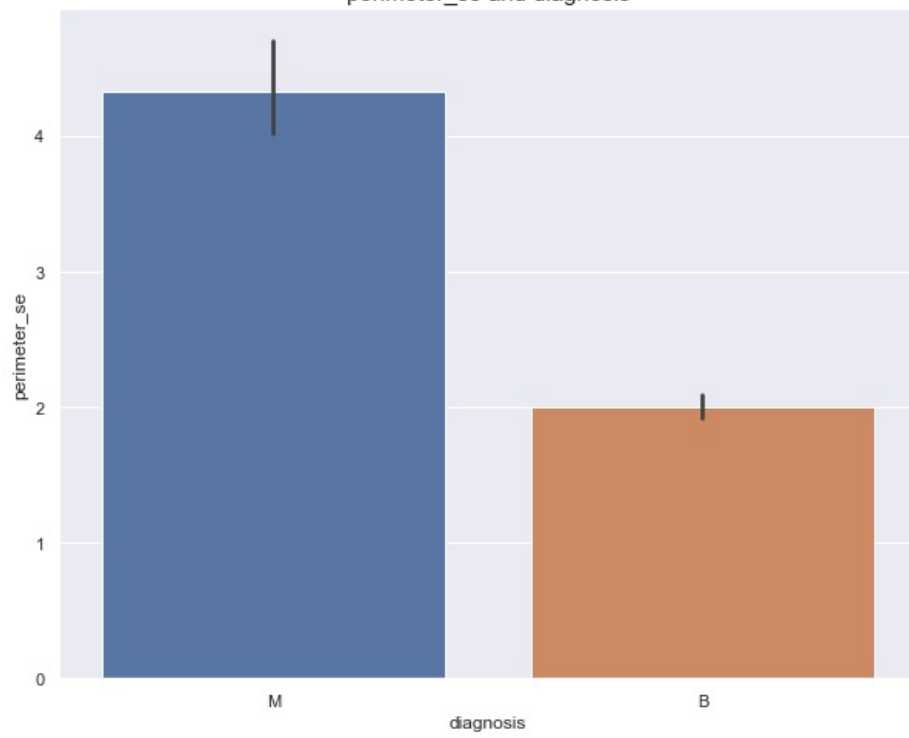




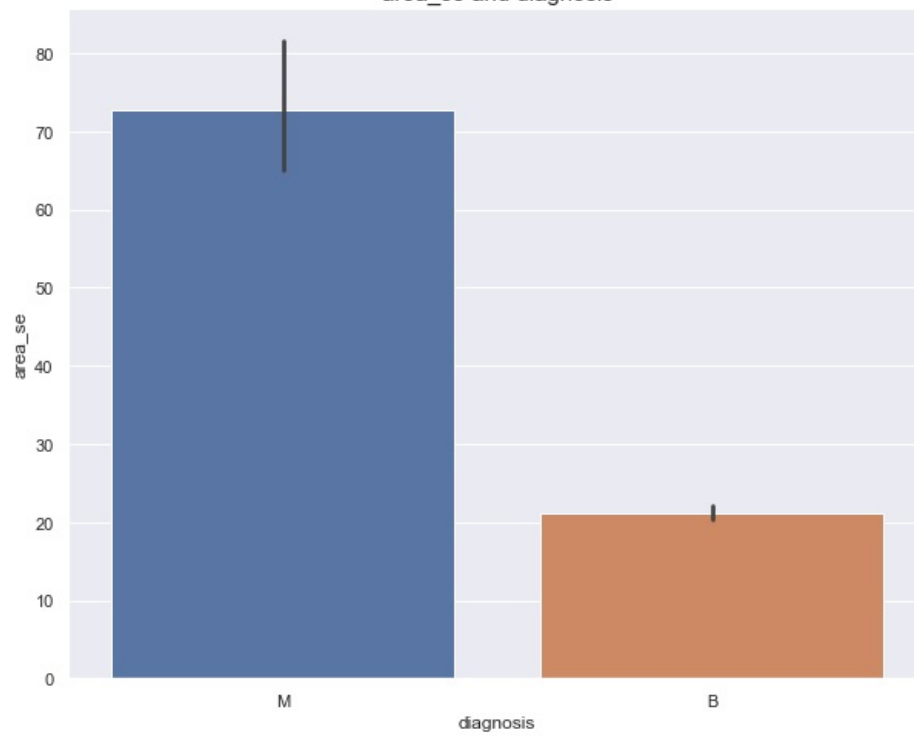


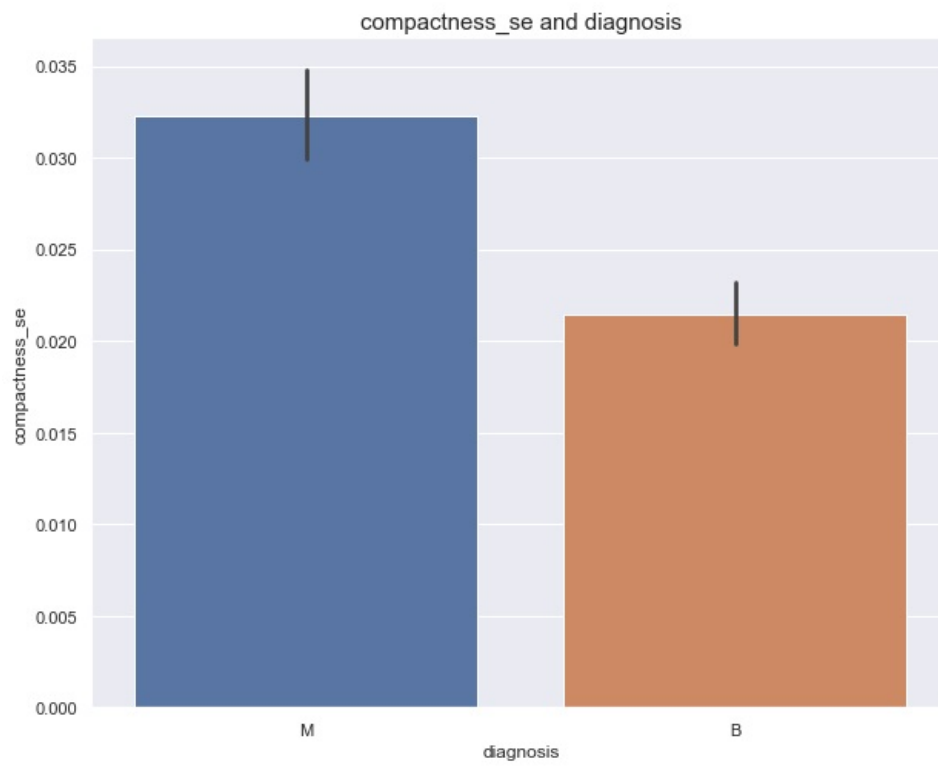
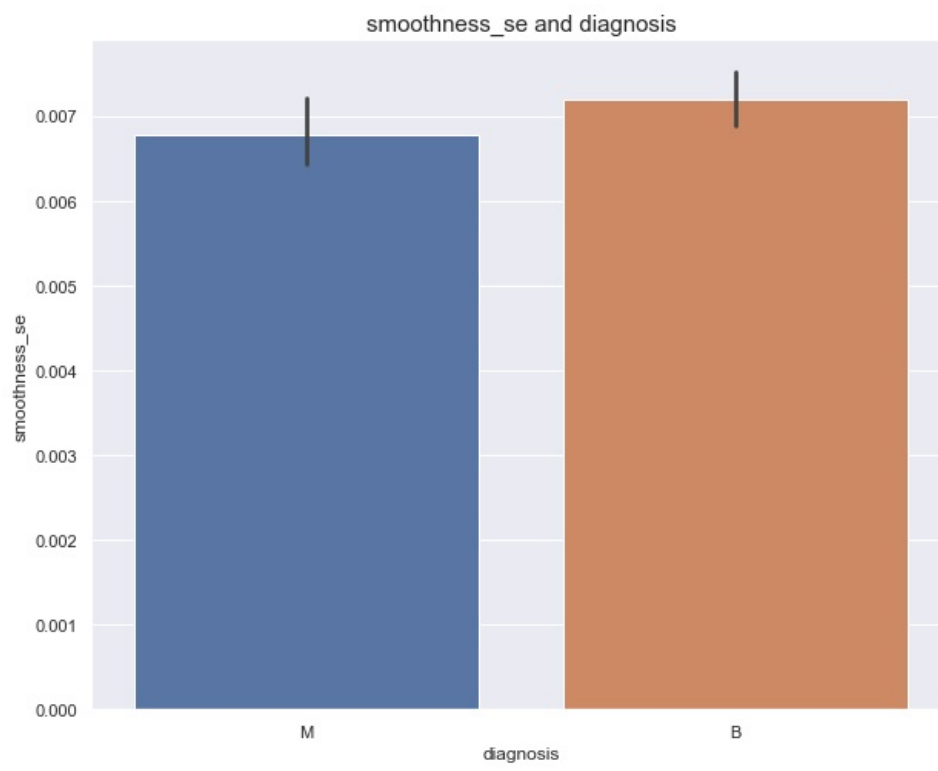


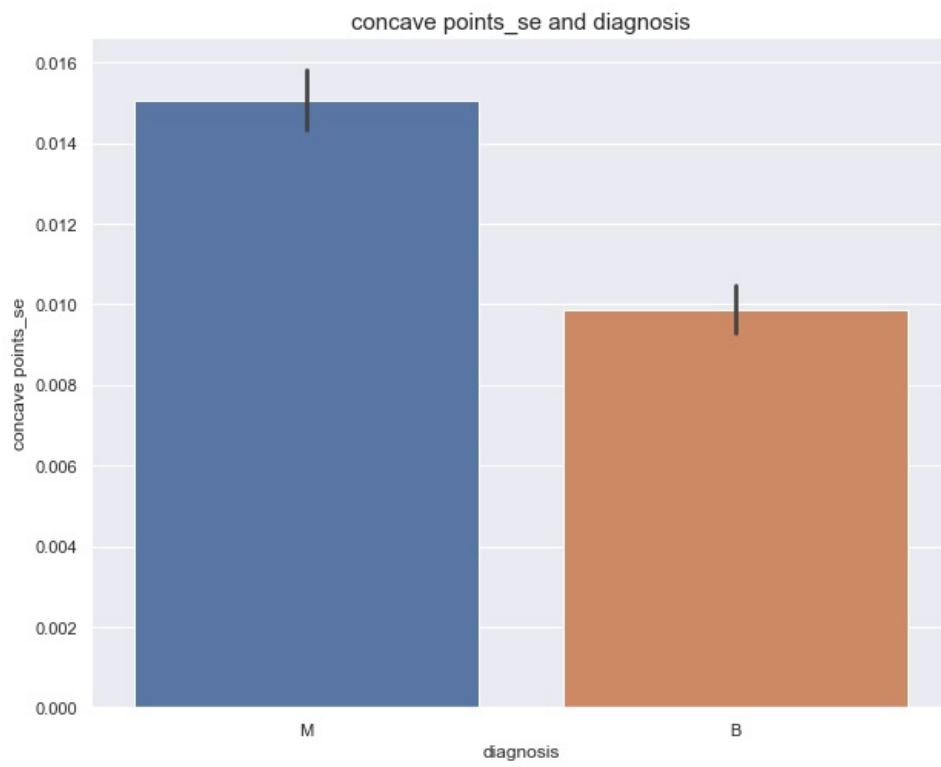
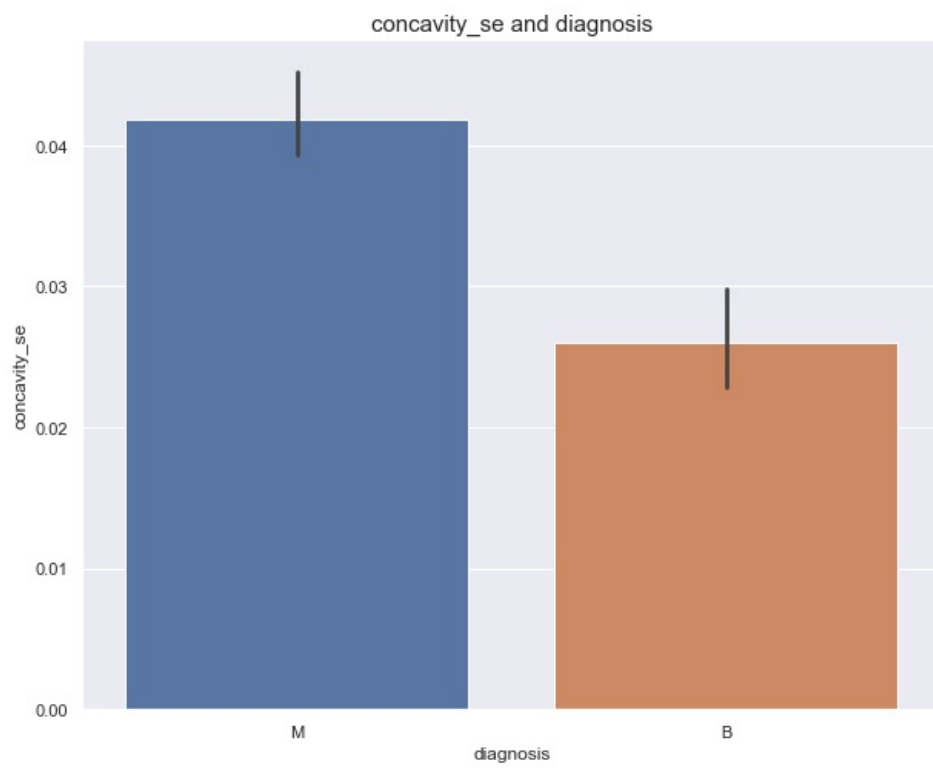
perimeter_se and diagnosis

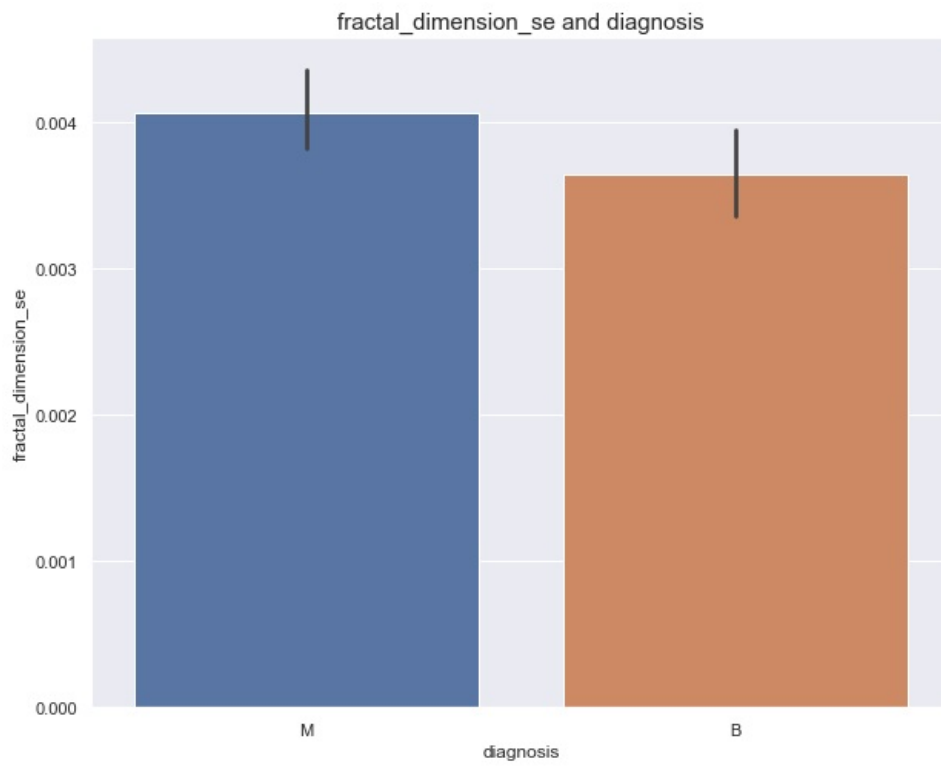
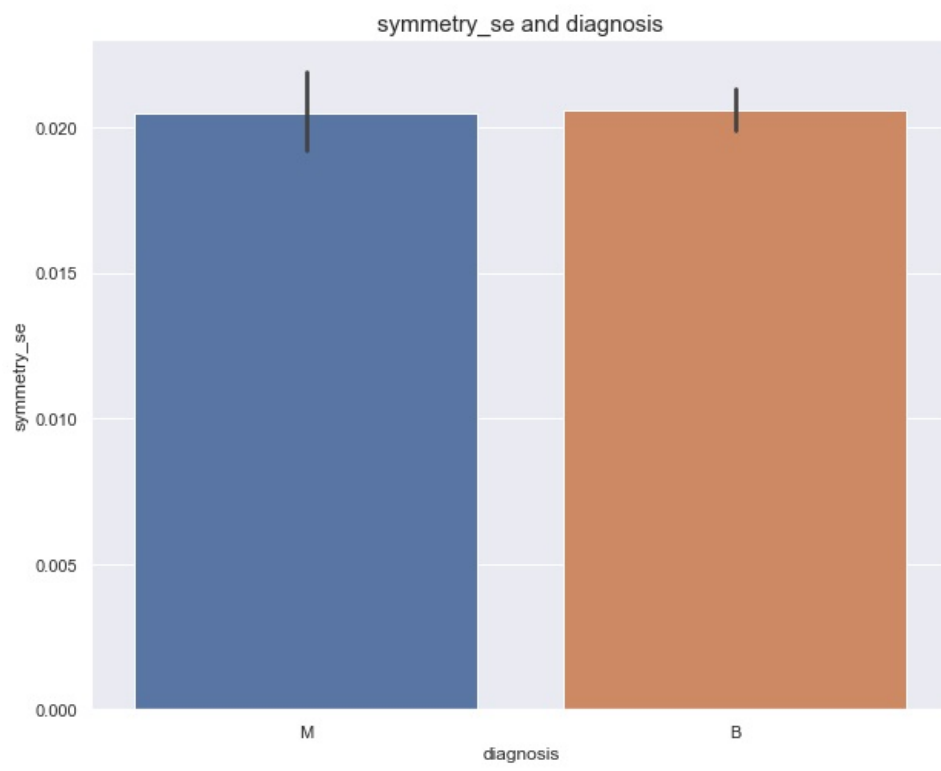


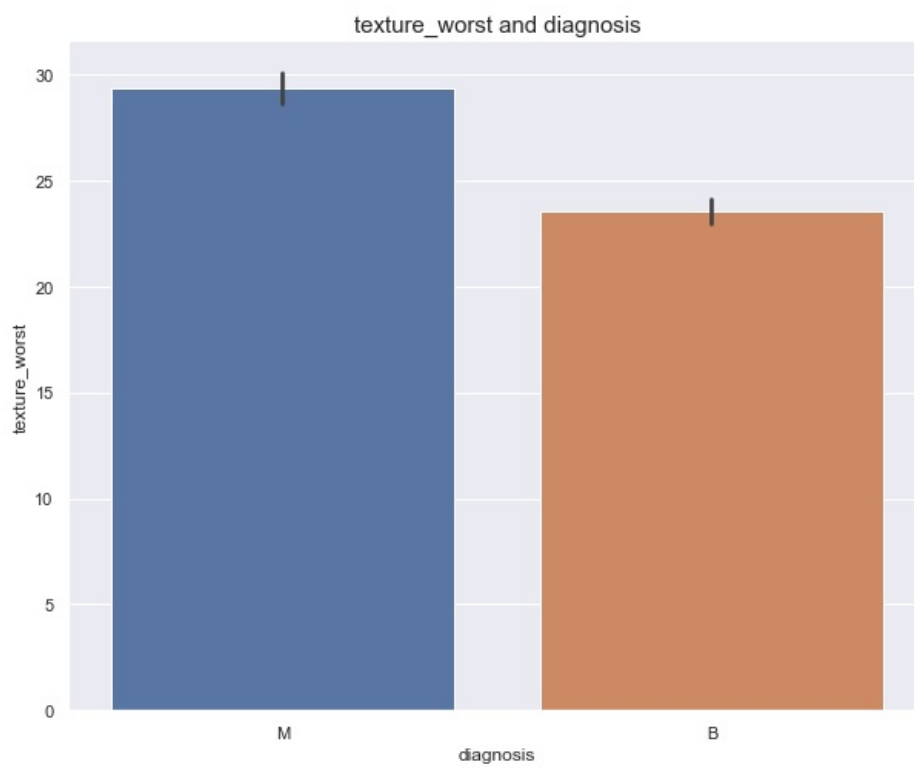
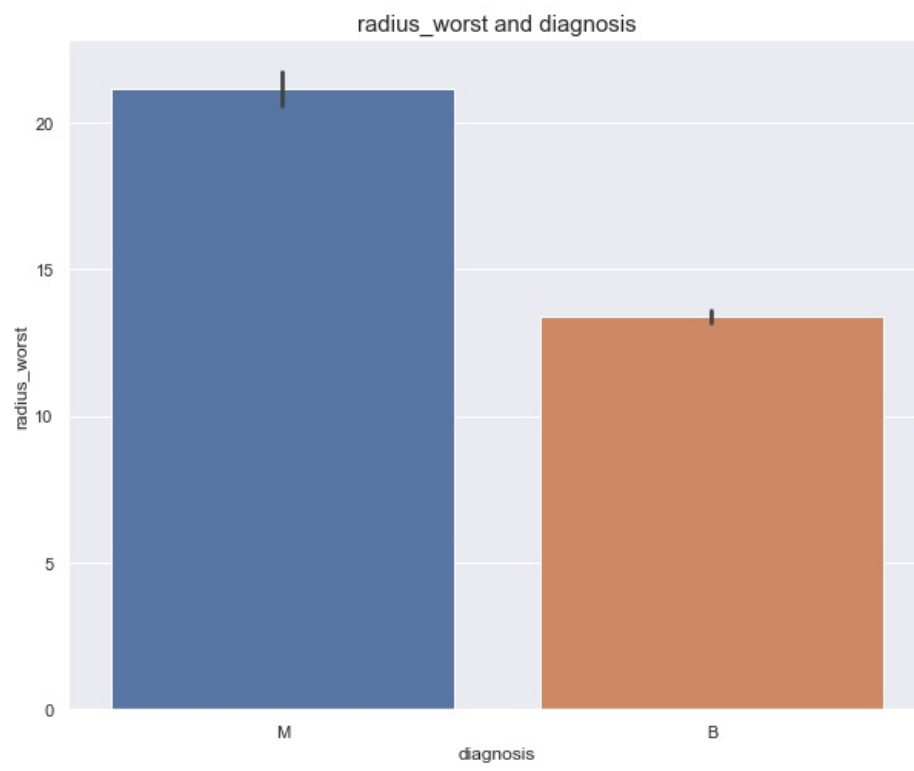
area_se and diagnosis

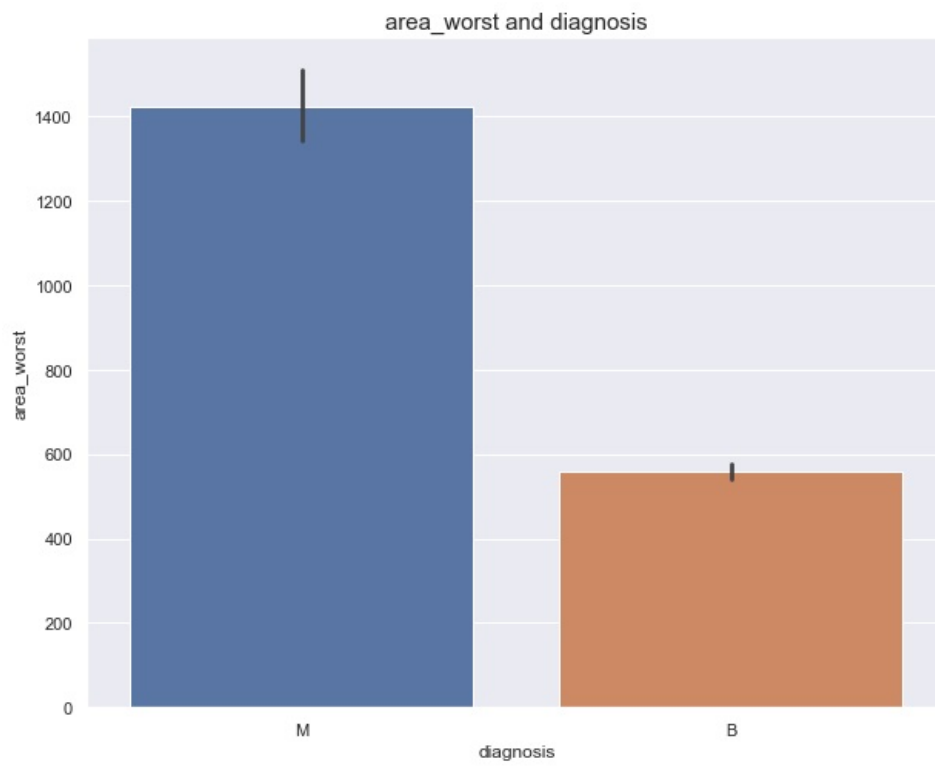
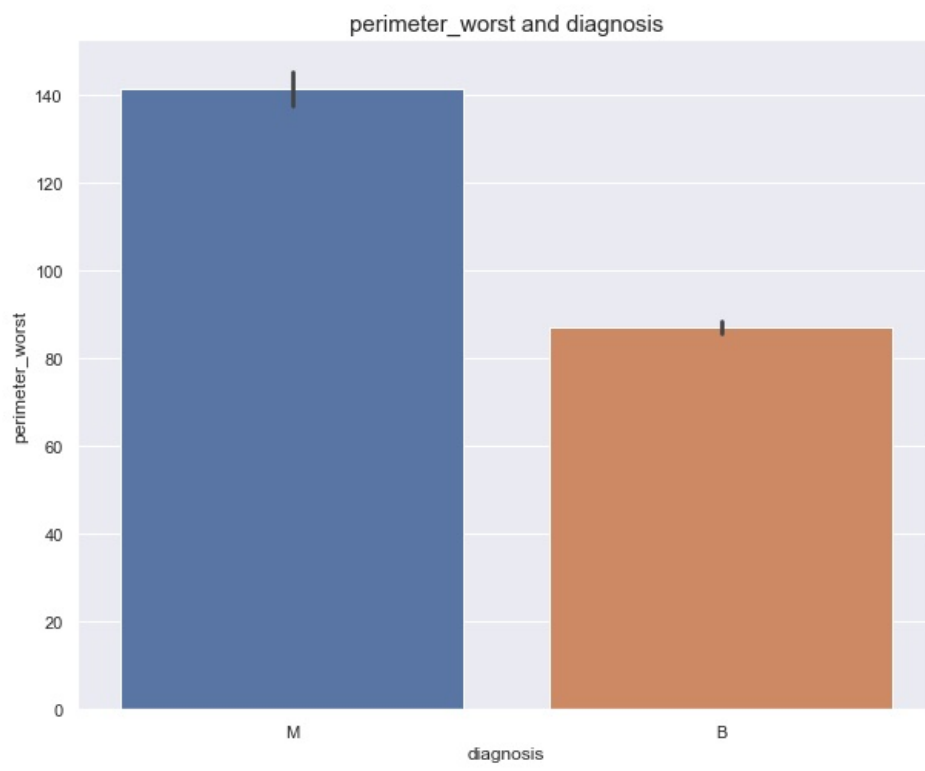


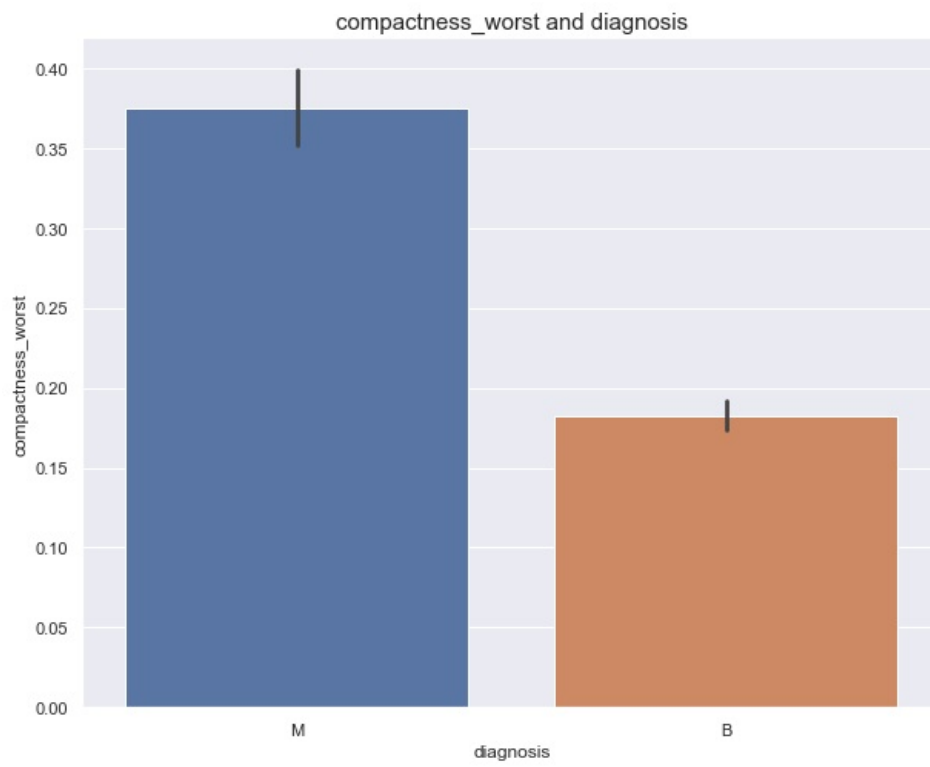
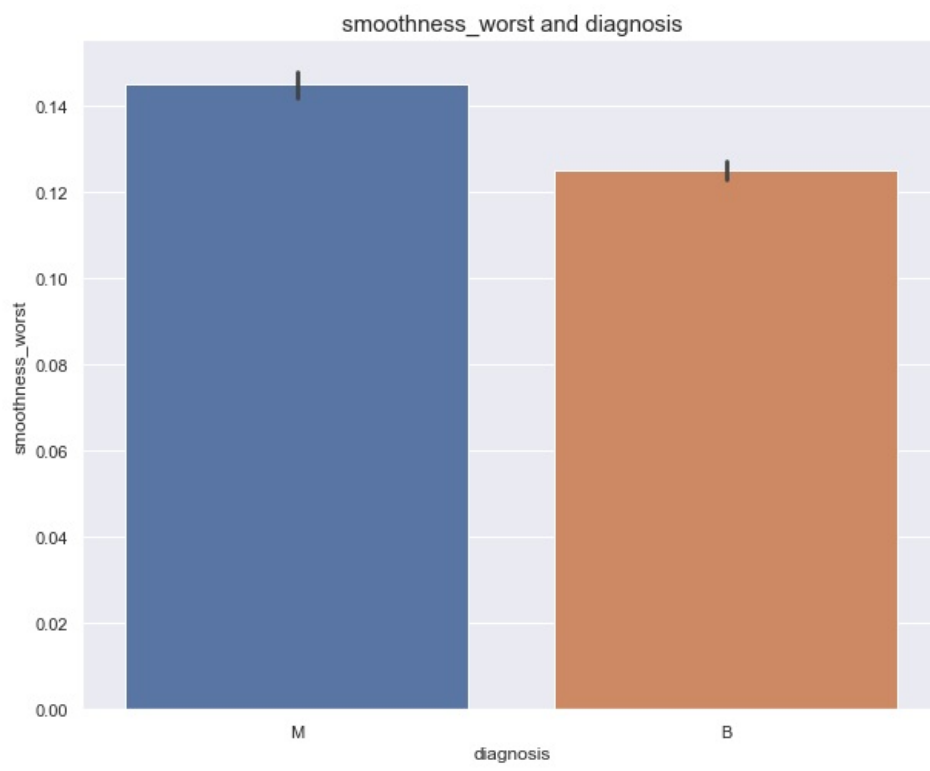


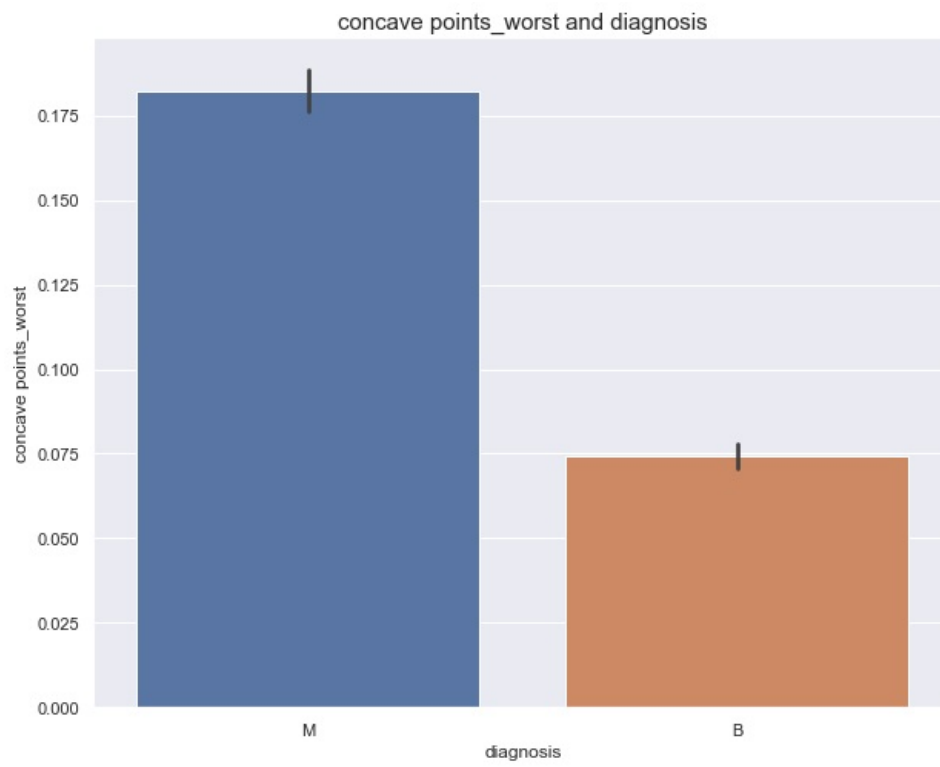
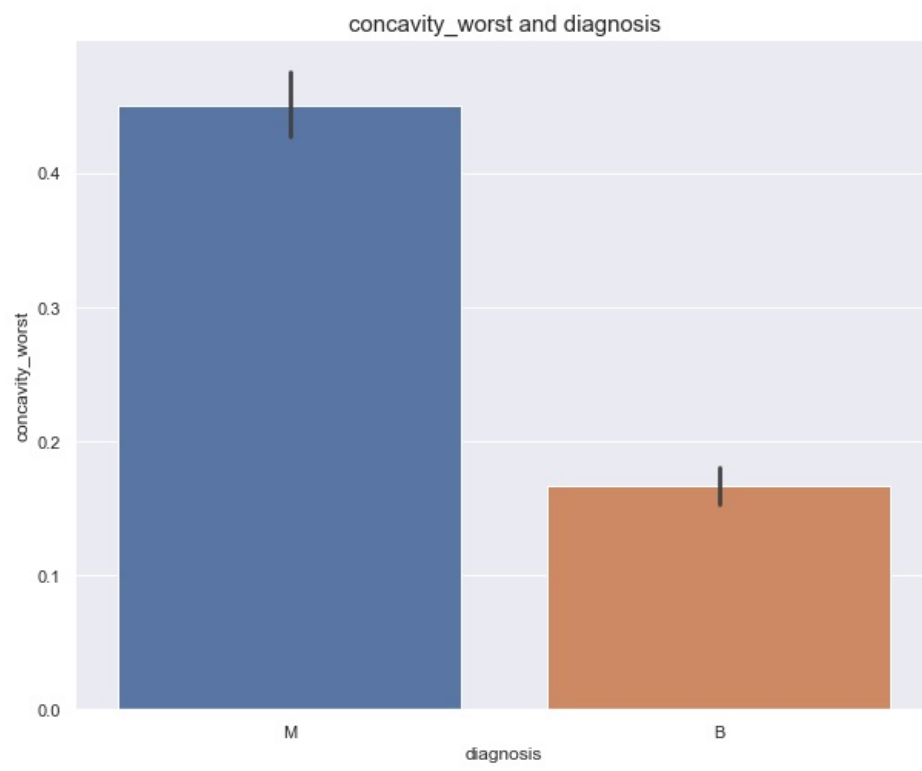


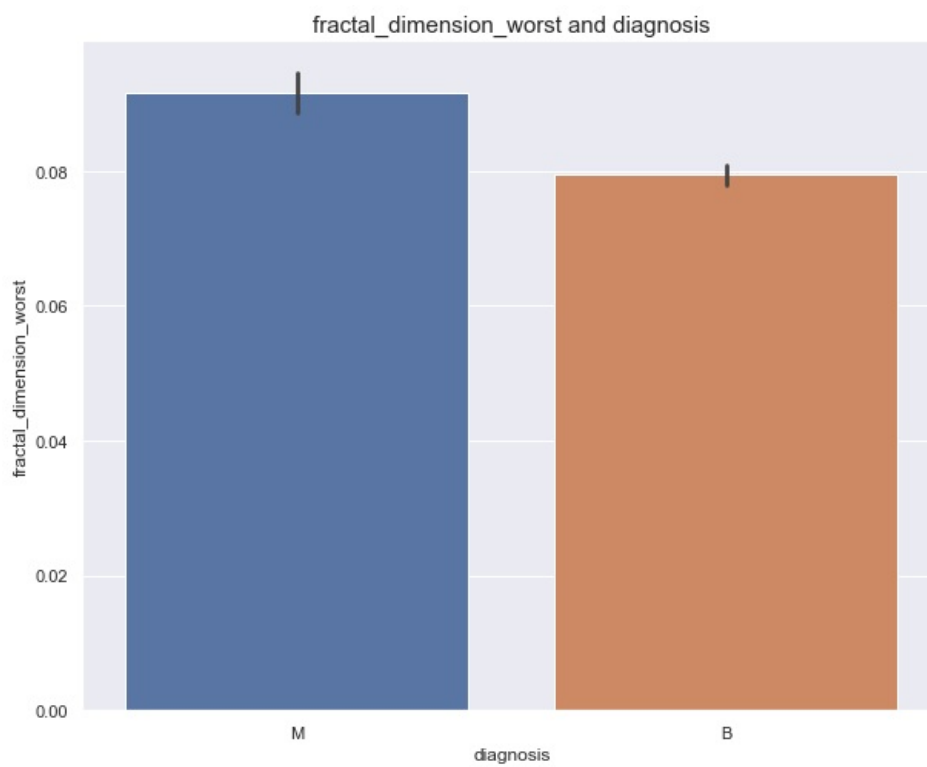
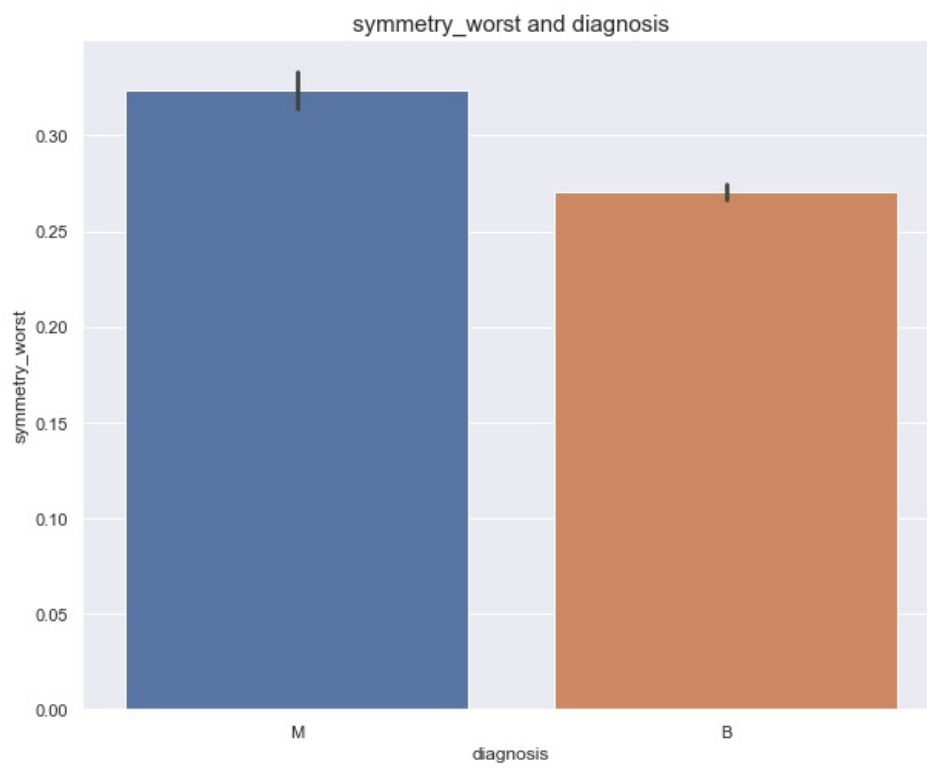




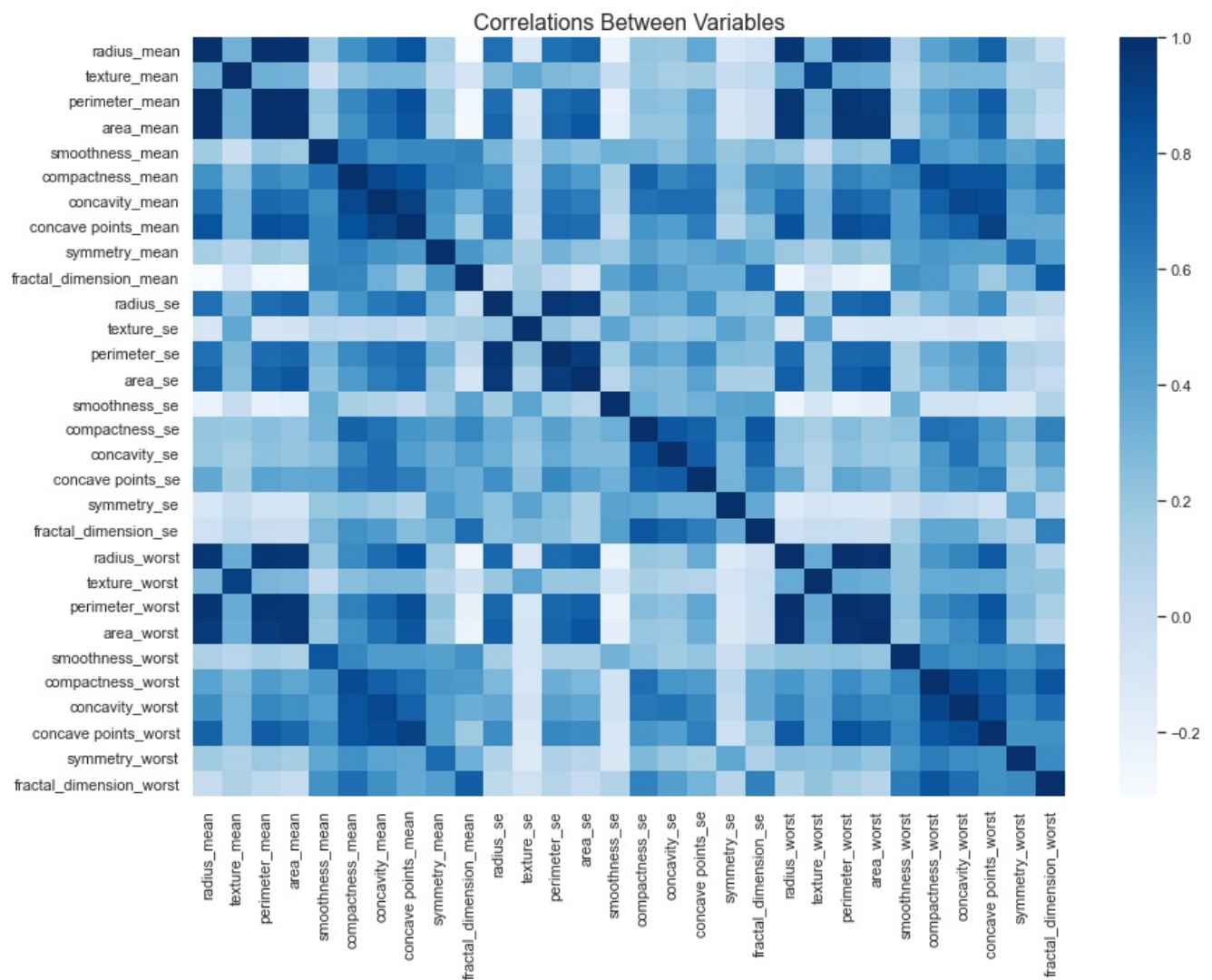








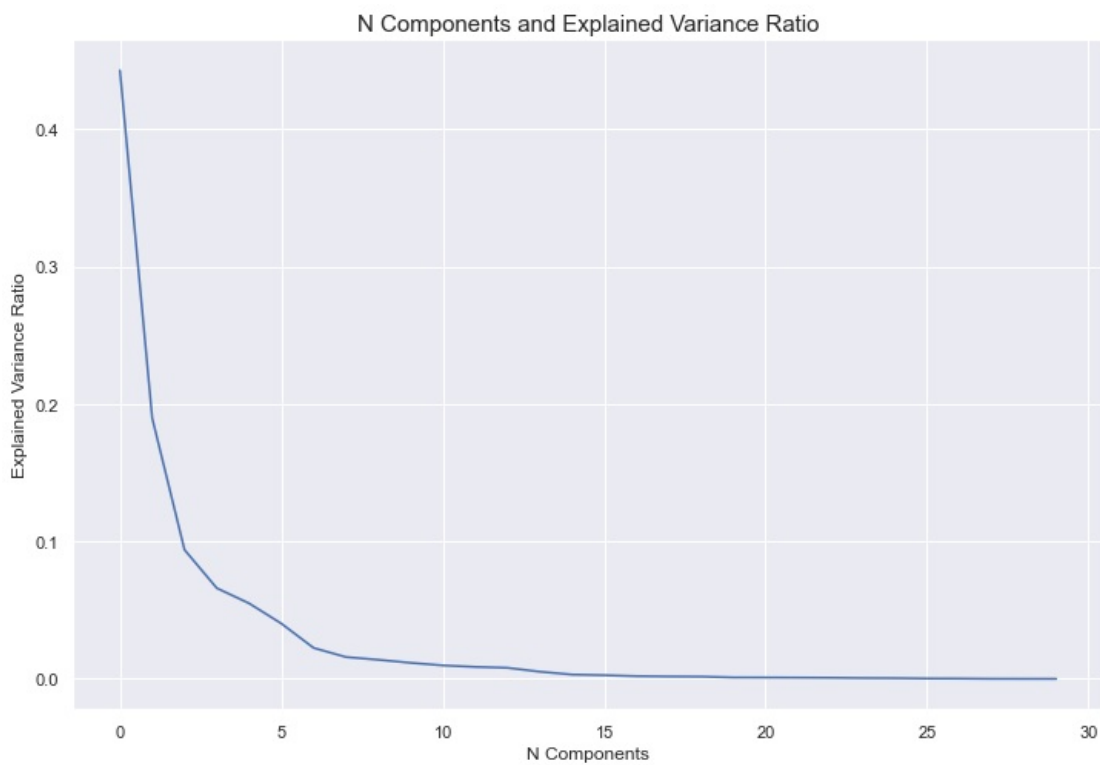
```
In [63]: plt.figure(figsize=(14,10))
sns.heatmap(df.corr(), cmap="Blues")
plt.title("Correlations Between Variables", size=16)
plt.show()
```



```
In [64]: X = df.drop("diagnosis", axis=1)
y = df["diagnosis"].replace({"B": 0, "M": 1})
```

```
In [65]: scaler = StandardScaler()
X = scaler.fit_transform(X)
```

```
In [66]: pca = PCA()
pca.fit(X)
plt.figure(figsize=(12,8))
plt.plot(pca.explained_variance_ratio_)
plt.title("N Components and Explained Variance Ratio", size=15)
plt.xlabel("N Components")
plt.ylabel("Explained Variance Ratio")
plt.show()
```



```
In [67]: pca = PCA(n_components = 5)
X = pca.fit_transform(X)
pca.explained_variance_ratio_.sum()
```

```
Out[67]: 0.8473427431679001
```

```
In [68]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
```

```
In [69]: models = pd.DataFrame(columns=["Model", "Accuracy Score"])
```

```
In [70]: log_reg = LogisticRegression()
log_reg.fit(X_train, y_train)
predictions = log_reg.predict(X_test)
score = accuracy_score(y_test, predictions)
print("Accuracy Score:", score)
new_row = {"Model": "LogisticRegression", "Accuracy Score": score}
models = models.append(new_row, ignore_index=True)
```

```
Accuracy Score: 0.9883040935672515
```

```
In [71]: rfc = RandomForestClassifier()
rfc.fit(X_train, y_train)
predictions = rfc.predict(X_test)
score = accuracy_score(y_test, predictions)
print("Accuracy Score:", score)
new_row = {"Model": "RandomForestClassifier", "Accuracy Score": score}
models = models.append(new_row, ignore_index=True)
```

```
Accuracy Score: 0.9707602339181286
```

```
In [72]: gbc = GradientBoostingClassifier()
gbc.fit(X_train, y_train)
predictions = gbc.predict(X_test)
score = accuracy_score(y_test, predictions)
print("Accuracy Score:", score)
new_row = {"Model": "GradientBoostingClassifier", "Accuracy Score": score}
models = models.append(new_row, ignore_index=True)
```

```
Accuracy Score: 0.9766081871345029
```

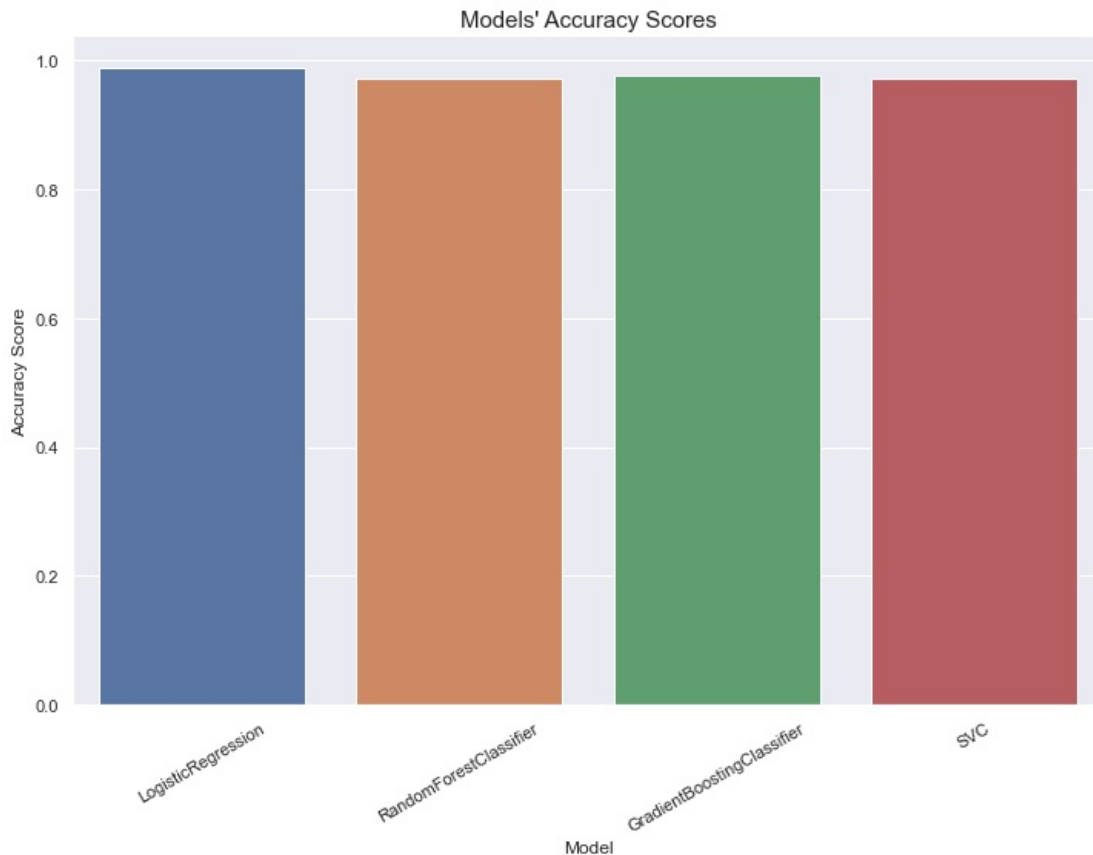
```
In [73]: svc = SVC()
svc.fit(X_train, y_train)
predictions = svc.predict(X_test)
score = accuracy_score(y_test, predictions)
print("Accuracy Score:", score)
new_row = {"Model": "SVC", "Accuracy Score": score}
models = models.append(new_row, ignore_index=True)
models.sort_values(by="Accuracy Score", ascending=False)
```

```
Accuracy Score: 0.9707602339181286
```

Out[73]:

	Model	Accuracy Score
0	LogisticRegression	0.988304
2	GradientBoostingClassifier	0.976608
1	RandomForestClassifier	0.97076
3	SVC	0.97076

```
In [74]: plt.figure(figsize=(12,8))
sns.barplot(x=models["Model"], y=models["Accuracy Score"])
plt.title("Models' Accuracy Scores", size=15)
plt.xticks(rotation=30)
plt.show()
```



```
In [75]: def visualize_roc_auc_curve(model, model_name):
pred_prob = model.predict_proba(X_test)
fpr, tpr, thresh = roc_curve(y_test, pred_prob[:,1], pos_label=1)

score = roc_auc_score(y_test, pred_prob[:, 1])

plt.figure(figsize=(10,8))
plt.plot(fpr, tpr, linestyle="--", color="orange", label="ROC curve (area = %0.5f)" % score)
plt.plot([0, 1], [0, 1], color="navy", lw=2, linestyle="--")

plt.title(f"{model_name} ROC Curve", size=15)
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.legend(loc="lower right", prop={'size': 15})
plt.show()
```

```
In [76]: tuned_models = pd.DataFrame(columns=["Model", "Accuracy Score"])
param_grid_log_reg = {"C": [0.0001, 0.001, 0.01, 0.1, 1, 10]}
```

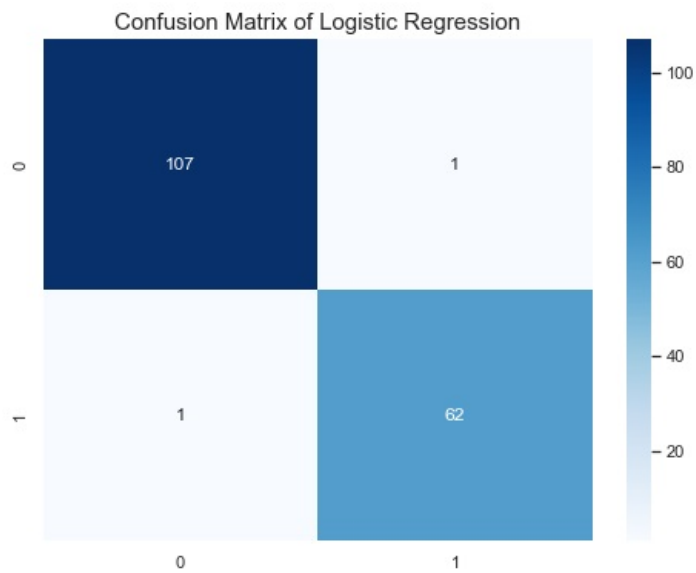
```
In [77]: grid_log_reg = GridSearchCV(LogisticRegression(), param_grid_log_reg, scoring="accuracy", cv=5, verbose=0, n_jobs=-1)
grid_log_reg.fit(X_train, y_train)
log_reg_params = grid_log_reg.best_params_
log_reg = LogisticRegression(**log_reg_params)
log_reg.fit(X_train, y_train)
predictions = log_reg.predict(X_test)
score = accuracy_score(y_test, predictions)
print("Accuracy Score:", score)

new_row = {"Model": "LogisticRegression", "Accuracy Score": score}
tuned_models = tuned_models.append(new_row, ignore_index=True)
```

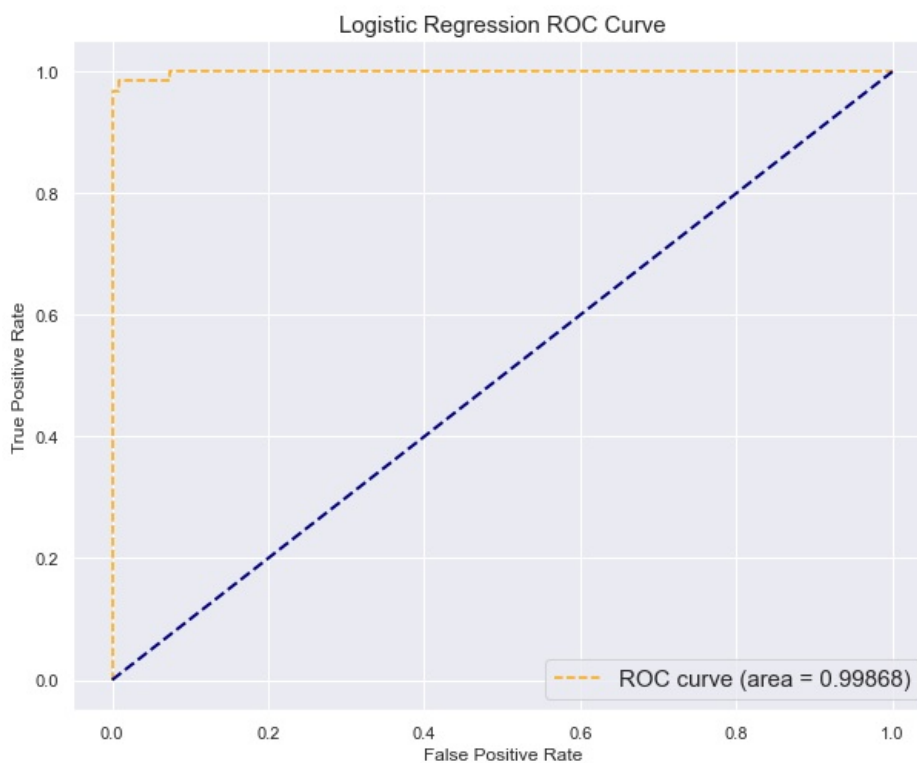
Accuracy Score: 0.9883040935672515

```
In [78]: plt.figure(figsize=(8,6))
sns.heatmap(confusion_matrix(y_test, predictions), annot=True, cmap="Blues", fmt="d")
plt.title("Confusion Matrix of Logistic Regression", size=15)
```

```
plt.show()
```



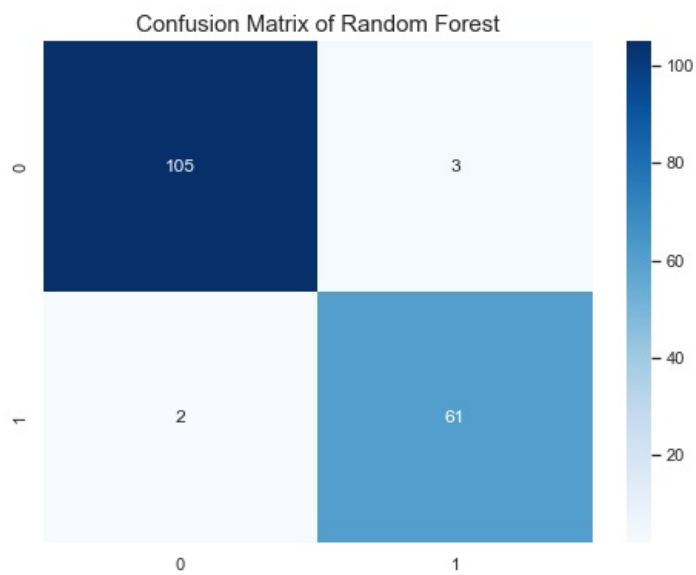
```
In [79]: visualize_roc_auc_curve(log_reg, "Logistic Regression")
```



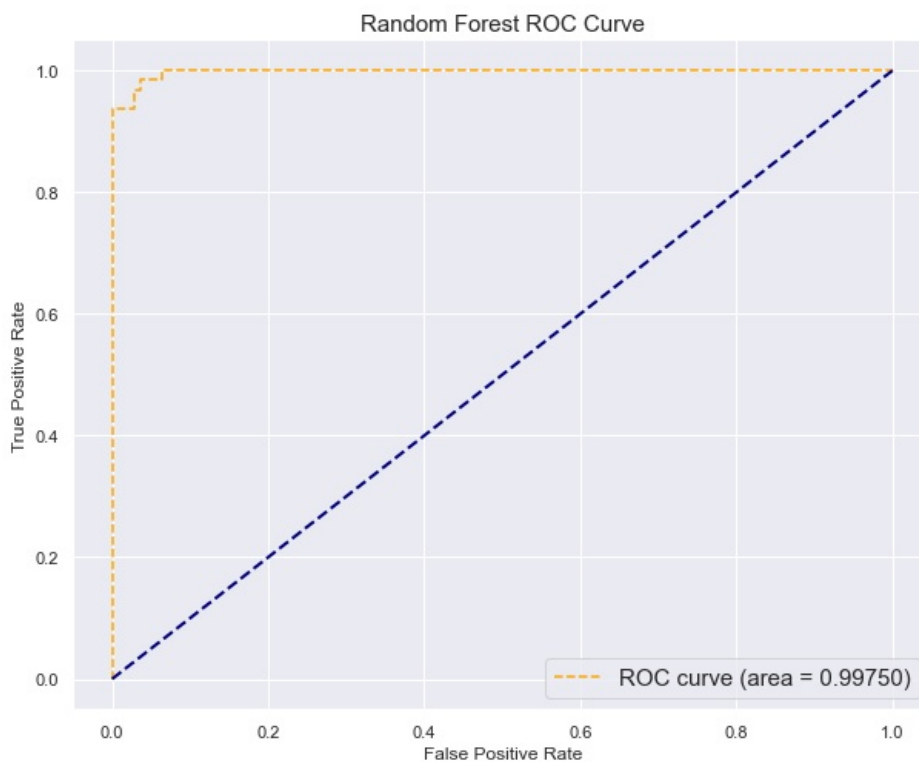
```
In [80]: param_grid_rfc = {"min_samples_split": [2, 3, 10],
                           "min_samples_leaf": [1, 3, 10],
                           "n_estimators": [100, 200, 500]}
grid_rfc = GridSearchCV(RandomForestClassifier(), param_grid_rfc, scoring="accuracy", cv=5, verbose=0, n_jobs=-1)
grid_rfc.fit(X_train, y_train)
rfc_params = grid_rfc.best_params_
rfc = RandomForestClassifier(**rfc_params)
rfc.fit(X_train, y_train)
predictions = rfc.predict(X_test)
score = accuracy_score(y_test, predictions)
print("Accuracy Score:", score)
new_row = {"Model": "RandomForestClassifier", "Accuracy Score": score}
tuned_models = tuned_models.append(new_row, ignore_index=True)
```

Accuracy Score: 0.9707602339181286

```
In [81]: plt.figure(figsize=(8,6))
sns.heatmap(confusion_matrix(y_test, predictions), annot=True, cmap="Blues", fmt='d')
plt.title("Confusion Matrix of Random Forest", size=15)
plt.show()
```



```
In [82]: visualize_roc_auc_curve(rfc, "Random Forest")
```



```
In [83]: param_grid_gbc = {'n_estimators': [100, 200, 500],
                           'learning_rate': [0.1, 0.05, 0.01],
                           'max_depth': [2, 3, 6],
                           'min_samples_leaf': [1, 2, 5]}
```

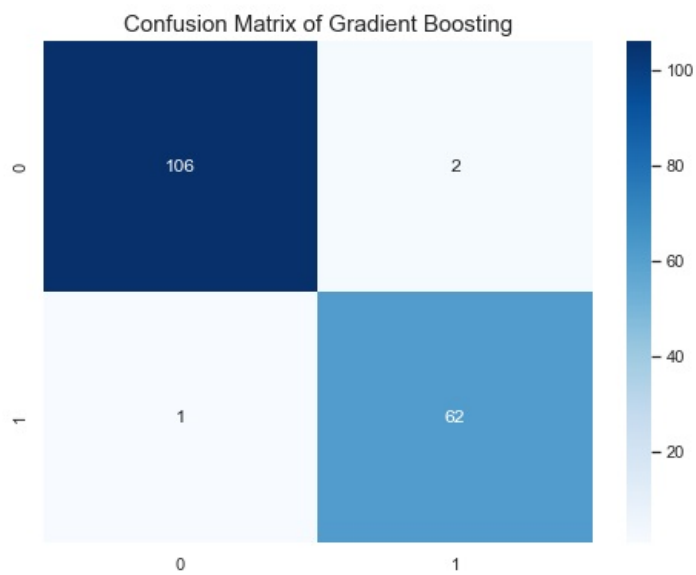
```
In [84]: grid_gbc = GridSearchCV(GradientBoostingClassifier(), param_grid_gbc, scoring="accuracy", cv=5, verbose=0, n_jobs=-1)

grid_gbc.fit(X_train, y_train)
gbc_params = grid_gbc.best_params_
gbc = GradientBoostingClassifier(**gbc_params)
gbc.fit(X_train, y_train)
predictions = gbc.predict(X_test)
score = accuracy_score(y_test, predictions)
print("Accuracy Score:", score)
```

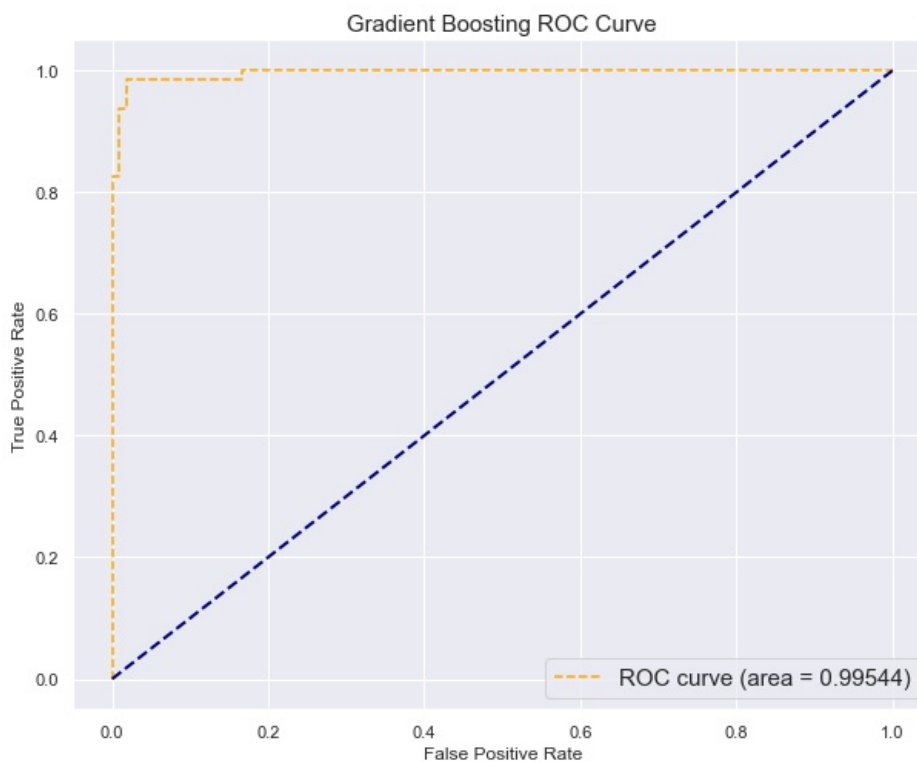
```
new_row = {"Model": "GradientBoostingClassifier", "Accuracy Score": score}
tuned_models = tuned_models.append(new_row, ignore_index=True)
```

Accuracy Score: 0.9824561403508771

```
In [85]: plt.figure(figsize=(8,6))
sns.heatmap(confusion_matrix(y_test, predictions), annot=True, cmap="Blues", fmt='d')
plt.title("Confusion Matrix of Gradient Boosting", size=15)
plt.show()
```



```
In [86]: visualize_roc_auc_curve(gbc, "Gradient Boosting")
```



```
In [87]: param_grid_svc = {'gamma': [ 0.001, 0.01, 0.1, 1, 10],
                          'C': [1, 10, 50, 100, 200, 300, 500, 1000]}
grid_svc = GridSearchCV(SVC(), param_grid_svc, scoring="accuracy", cv=5, verbose=0, n_jobs=-1)
grid_svc.fit(X_train, y_train)
```

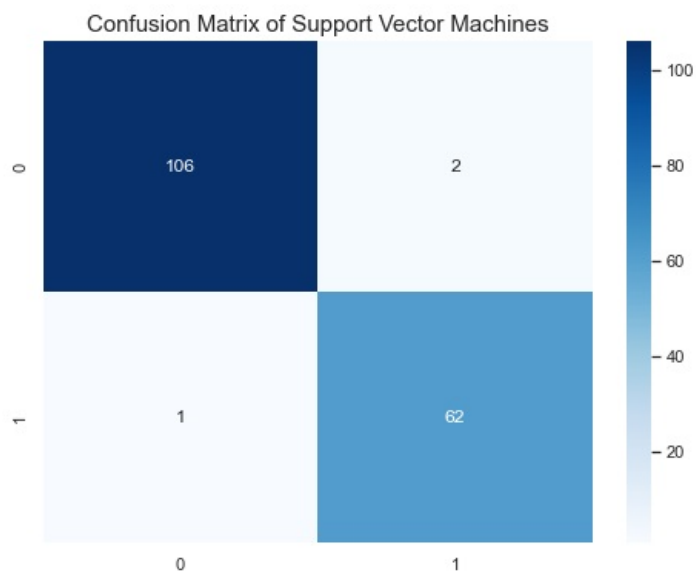
```
Out[87]: GridSearchCV(cv=5, estimator=SVC(), n_jobs=-1,
                      param_grid={'C': [1, 10, 50, 100, 200, 300, 500, 1000],
                                   'gamma': [0.001, 0.01, 0.1, 1, 10]},
                      scoring='accuracy')
```

```
In [88]: svc_params = grid_svc.best_params_
svc = SVC(**svc_params)
svc.fit(X_train, y_train)
predictions = svc.predict(X_test)
score = accuracy_score(y_test, predictions)
print("Accuracy Score:", score)

new_row = {"Model": "SVC", "Accuracy Score": score}
tuned_models = tuned_models.append(new_row, ignore_index=True)
```

Accuracy Score: 0.9824561403508771

```
In [89]: plt.figure(figsize=(8,6))
sns.heatmap(confusion_matrix(y_test, predictions), annot=True, cmap="Blues", fmt='d')
plt.title("Confusion Matrix of Support Vector Machines", size=15)
plt.show()
```



```
In [90]: tuned_models.sort_values(by="Accuracy Score", ascending=False)
```

Out[90]:

	Model	Accuracy Score
0	LogisticRegression	0.988304
2	GradientBoostingClassifier	0.982456
3	SVC	0.982456
1	RandomForestClassifier	0.97076

```
In [91]: plt.figure(figsize=(12, 8))
sns.barplot(x=tuned_models["Model"], y=tuned_models["Accuracy Score"])
plt.title("Models' Accuracy Scores After Hyperparameter Tuning", size=15)
plt.xticks(rotation=30)
plt.show()
```

