

Experiment: Effect of PCA on Classifiers (Breast Cancer Dataset) – Filled

September 18, 2025

Objective

Compare classifier performance (Accuracy, Precision, Recall, F1) and stability (fold-wise variance) on the Breast Cancer Wisconsin (WDBC) dataset under two settings: **No-PCA** and **With-PCA** (variance-target PCA, 95% explained variance). Evaluate per-model hyperparameter tuning, 5-fold CV, and stacking behavior under both settings to determine when PCA helps.

Libraries used

- Python 3.x, numpy, pandas, matplotlib, seaborn
- scikit-learn: preprocessing, model_selection, decomposition (PCA), metrics, pipeline, ensemble, classifiers
- xgboost (optional)
- scipy (optional, for paired t-test)
- LaTeX packages: pdfpages, listings, csvsimple, booktabs, graphicx, caption, subcaption

Contents

1	Code	2
2	Output	12
3	Results summary (test set)	30
4	5-Fold Cross-Validation Results (All Models)	31
5	Observation Questions	32
6	Conclusion	32

1 Code

Below is the full Python script used for the experiment (no omissions). Save this as a ‘.py’ if you want to re-run it exactly.

```
1 # breast_cancer_with_pca_comparison_full.py
2 import time
3 import numpy as np
4 import pandas as pd
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7
8 from sklearn.preprocessing import StandardScaler, LabelEncoder
9 from sklearn.model_selection import train_test_split, GridSearchCV,
   KFold, cross_val_score, StratifiedKFold
10 from sklearn.decomposition import PCA
11 from sklearn.pipeline import Pipeline
12 from sklearn.metrics import (
13     accuracy_score, precision_score, recall_score, f1_score,
14     roc_curve, auc, confusion_matrix, ConfusionMatrixDisplay
15 )
16 from sklearn.base import clone
17
18 # Models
19 from sklearn.tree import DecisionTreeClassifier
20 from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier
   , GradientBoostingClassifier, StackingClassifier
21 from sklearn.naive_bayes import GaussianNB
22 from sklearn.svm import SVC
23 from sklearn.neighbors import KNeighborsClassifier
24 from sklearn.linear_model import LogisticRegression
25
26 # Optional libraries
27 try:
28     import xgboost as xgb
29     XGB_AVAILABLE = True
30 except Exception:
31     XGB_AVAILABLE = False
32
33 try:
34     from scipy.stats import ttest_rel
35     SCIPY_AVAILABLE = True
36 except Exception:
37     SCIPY_AVAILABLE = False
38
39 # ----- User settings -----
40 DATA_PATH = "wdbc.data" # change if your file is at
   different path
41 PCA_VARIANCE = 0.95 # keep 95% explained variance when
   PCA is used
42 TEST_SIZE = 0.20
43 RANDOM_STATE = 42
44 CV_FOLDS = 5
45 # -----
46
47 # ----- Load & preprocess -----
48 cols = ["ID", "Diagnosis"] + [f"feature_{i}" for i in range(1, 31)]
49 df = pd.read_csv(DATA_PATH, header=None, names=cols)
```

```

50 df.drop(columns=["ID"], inplace=True)
51 df["Diagnosis"] = LabelEncoder().fit_transform(df["Diagnosis"]) # M=1,
    B=0
52
53 X_raw = df.drop(columns=["Diagnosis"])
54 y = df["Diagnosis"].values
55
56 scaler = StandardScaler()
57 X_scaled = scaler.fit_transform(X_raw)
58
59 X_train, X_test, y_train, y_test = train_test_split(
60     X_scaled, y, test_size=TEST_SIZE, random_state=RANDOM_STATE,
        stratify=y
61 )
62
63 print("\nClass Distribution (0 = Benign, 1 = Malignant):")
64 print(df["Diagnosis"].value_counts())
65
66 plt.figure(figsize=(10, 8))
67 sns.heatmap(df.drop(columns=["Diagnosis"]).corr(), cmap="coolwarm",
    cbar=True)
68 plt.title("Feature Correlation Heatmap")
69 plt.tight_layout()
70 plt.show()
71
72 # ----- Evaluation helper -----
73 def evaluate(name, model, X_test, y_test, plot=True):
74     """Evaluate model: prints metrics, plots ROC + confusion matrix (if
        plot=True)."""
75     if plot:
76         fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 5))
77         fig.suptitle(f"{name} Performance", fontsize=16)
78
79         # Probabilities for ROC
80         y_pred = model.predict(X_test)
81         if hasattr(model, "predict_proba"):
82             probs = model.predict_proba(X_test)[:, 1]
83         else:
84             # if no predict_proba, try decision_function
85             try:
86                 dfun = model.decision_function(X_test)
87                 # normalize to [0,1]
88                 probs = (dfun - dfun.min()) / (dfun.max() - dfun.min())
89             except Exception:
90                 probs = None
91
92         if plot and probs is not None:
93             fpr, tpr, _ = roc_curve(y_test, probs)
94             roc_auc = auc(fpr, tpr)
95             ax1.plot(fpr, tpr, lw=2, label=f"ROC AUC = {roc_auc:.3f}")
96             ax1.plot([0, 1], [0, 1], lw=2, linestyle="--")
97             ax1.set_xlim([0.0, 1.0])
98             ax1.set_ylim([0.0, 1.05])
99             ax1.set_xlabel("False Positive Rate")
100             ax1.set_ylabel("True Positive Rate")
101             ax1.set_title("ROC Curve")
102             ax1.legend(loc="lower right")
103         elif plot:

```

```

104         ax1.text(0.5, 0.5, "No probability scores available", ha="
           center")
105         ax1.set_title("ROC Curve")
106
107     if plot:
108         cm = confusion_matrix(y_test, y_pred)
109         disp = ConfusionMatrixDisplay(confusion_matrix=cm)
110         disp.plot(ax=ax2, cmap=None)
111         ax2.set_title("Confusion Matrix")
112         plt.tight_layout()
113         plt.show()
114
115     acc = accuracy_score(y_test, y_pred)
116     prec = precision_score(y_test, y_pred)
117     rec = recall_score(y_test, y_pred)
118     f1 = f1_score(y_test, y_pred)
119     print(f"\n{name}")
120     print("Accuracy :", acc)
121     print("Precision:", prec)
122     print("Recall   :", rec)
123     print("F1 Score :", f1)
124     return acc, prec, rec, f1
125
126 # ----- Models + hyperparameter grids -----
127 models_params = {
128     "Decision Tree": (
129         DecisionTreeClassifier(random_state=RANDOM_STATE),
130         {"criterion": ["gini", "entropy"], "max_depth": [3, 5, 10, None]
131         }],
132     "Random Forest": (
133         RandomForestClassifier(random_state=RANDOM_STATE, n_jobs=-1),
134         {"n_estimators": [50, 100, 200], "max_depth": [3, 5, 10, None],
135         "criterion": ["gini", "entropy"]}],
136     "AdaBoost": (
137         AdaBoostClassifier(random_state=RANDOM_STATE, estimator=
138         DecisionTreeClassifier(random_state=RANDOM_STATE)),
139         {"n_estimators": [50, 100, 200], "learning_rate": [0.01, 0.1,
140         1], "estimator__max_depth": [1, 3, 5]}],
141     "Gradient Boosting": (
142         GradientBoostingClassifier(random_state=RANDOM_STATE),
143         {"n_estimators": [50, 100, 200], "learning_rate": [0.01, 0.1,
144         0.5], "max_depth": [3, 5, 7]}],
145     "KNN": (
146         KNeighborsClassifier(),
147         {"n_neighbors": [3, 5, 7], "weights": ["uniform", "distance"],
148         "p": [1, 2]}],
149     "SVM": (
150         SVC(probability=True, random_state=RANDOM_STATE),
151         {"kernel": ["rbf", "linear"], "C": [0.1, 1, 10], "gamma": ["
152         scale", "auto"]}],
153     "Naive Bayes": (
154         GaussianNB(),

```

```

154         {"var_smoothing": [1e-9, 1e-8, 1e-7]}
155     ),
156     "Logistic Regression": (
157         LogisticRegression(max_iter=5000, solver="saga", random_state=
158             RANDOM_STATE),
159         {"C": [0.01, 0.1, 1, 10], "penalty": ["l2", "none"]}
160     )
161 }
162 if XGB_AVAILABLE:
163     models_params["XGBoost"] = (
164         xgb.XGBClassifier(use_label_encoder=False, eval_metric="logloss
165             ", random_state=RANDOM_STATE, n_jobs=-1),
166         {"n_estimators": [50, 100, 200], "learning_rate": [0.01, 0.1,
167             0.3], "max_depth": [3, 5, 7], "gamma": [0, 0.1, 0.3]}
168     )
169 # ----- Results containers -----
170 results_table = [] # rows: (name, variant, best_params, acc, prec,
171     rec, f1, elapsed)
172 trial_tables = {} # name -> { 'no_pca': df, 'with_pca': df }
173 best_estimators = {} # name -> { 'no_pca': estimator, 'with_pca':
174     estimator}
175 # ----- Helper to prefix param grid keys for pipeline ('clf__'
176     prefix) -----
177 def prefix_param_grid(param_grid, prefix="clf__"):
178     new_grid = {}
179     for k, v in param_grid.items():
180         new_key = prefix + k
181         new_grid[new_key] = v
182     return new_grid
183 # ----- GridSearch: No-PCA and With-PCA -----
184 for name, (model, params) in models_params.items():
185     print(f"\n--- Grid Search for {name} ---")
186     # Make two pipelines
187     pipe_no_pca = Pipeline([("clf", model)]) # using scaled X (
188         X_scaled), so no scaler in pipeline
189     pipe_with_pca = Pipeline([("pca", PCA(n_components=PCA_VARIANCE,
190         svd_solver="full")), ("clf", model)])
191
192     # Prefix param grid keys for pipeline
193     grid_params = prefix_param_grid(params) # adds 'clf__' prefix
194
195     # GridSearch No-PCA
196     print(" GridSearch (No-PCA)...")
197     gs_no = GridSearchCV(pipe_no_pca, param_grid=grid_params, cv=
198         CV_FOLDS, scoring="accuracy", n_jobs=-1, return_train_score=
199         False)
200     gs_no.fit(X_train, y_train)
201     best_no = gs_no.best_estimator_
202     best_estimators.setdefault(name, {})[ 'no_pca' ] = best_no
203     print(" Best params (No-PCA):", gs_no.best_params_)
204     print(" Best CV score (No-PCA):", gs_no.best_score_)
205
206     # Record top-5 trials (no-pca)
207     trials_no = []

```

```

202     for i in range(len(gs_no.cv_results_['params'])):
203         trials_no.append({
204             **gs_no.cv_results_['params'][i],
205             "CV Accuracy": gs_no.cv_results_['mean_test_score'][i]
206         })
207     trial_df_no = pd.DataFrame(trials_no).sort_values("CV Accuracy",
208         ascending=False).head(5)
209     y_pred_no = best_no.predict(X_test)
210     trial_df_no["F1 Score (Test)"] = f1_score(y_test, y_pred_no)
211     trial_tables.setdefault(name, {})[ 'no_pca' ] = trial_df_no
212
213     # Evaluate best_no on test set
214     start = time.time()
215     best_no.fit(X_train, y_train)
216     elapsed = time.time() - start
217     res_no = evaluate(f"{name} (No-PCA)", best_no, X_test, y_test, plot
218         =True)
219     results_table.append((name, "No-PCA", gs_no.best_params_, *res_no,
220         elapsed))
221
222     # GridSearch With-PCA
223     print(" GridSearch (With-PCA)...")
224     gs_pca = GridSearchCV(pipe_with_pca, param_grid=grid_params, cv=
225         CV_FOLDS, scoring="accuracy", n_jobs=-1, return_train_score=
226         False)
227     gs_pca.fit(X_train, y_train)
228     best_pca = gs_pca.best_estimator_
229     best_estimators[name][ 'with_pca' ] = best_pca
230     print(" Best params (With-PCA):", gs_pca.best_params_)
231     print(" Best CV score (With-PCA):", gs_pca.best_score_)
232
233     # Record top-5 trials (with-pca)
234     trials_pca = []
235     for i in range(len(gs_pca.cv_results_['params'])):
236         trials_pca.append({
237             **gs_pca.cv_results_['params'][i],
238             "CV Accuracy": gs_pca.cv_results_['mean_test_score'][i]
239         })
240     trial_df_pca = pd.DataFrame(trials_pca).sort_values("CV Accuracy",
241         ascending=False).head(5)
242     y_pred_pca = best_pca.predict(X_test)
243     trial_df_pca["F1 Score (Test)"] = f1_score(y_test, y_pred_pca)
244     trial_tables[name][ 'with_pca' ] = trial_df_pca
245
246     # Evaluate best_pca on test set
247     start = time.time()
248     best_pca.fit(X_train, y_train)
249     elapsed = time.time() - start
250     res_pca = evaluate(f"{name} (With-PCA)", best_pca, X_test, y_test,
251         plot=True)
252     results_table.append((name, "With-PCA", gs_pca.best_params_, *
253         res_pca, elapsed))
254
255     # ----- Stacking classifiers (3 original variants) -----
256     print("\n--- Stacking Classifiers (original single-run variants) ---")
257     stacking_variants = {
258         "Stacking (SVM + NB + DT -> LogisticRegression)":
259             StackingClassifier(

```

```

251     estimators=[
252         ("svm", SVC(probability=True, kernel="rbf", C=1, gamma="
253             scale")),
254         ("nb", GaussianNB()),
255         ("dt", DecisionTreeClassifier(max_depth=5, random_state=
256             RANDOM_STATE))
257     ],
258     final_estimator=LogisticRegression(max_iter=500, random_state=
259         RANDOM_STATE),
260     n_jobs=-1
261 ),
262 "Stacking (SVM + NB + DT -> RandomForest)": StackingClassifier(
263     estimators=[
264         ("svm", SVC(probability=True, kernel="rbf", C=1, gamma="
265             scale")),
266         ("nb", GaussianNB()),
267         ("dt", DecisionTreeClassifier(max_depth=5, random_state=
268             RANDOM_STATE))
269     ],
270     final_estimator=RandomForestClassifier(n_estimators=100,
271         random_state=RANDOM_STATE),
272     n_jobs=-1
273 ),
274 "Stacking (SVM + DT + KNN -> LogisticRegression)":
275     StackingClassifier(
276         estimators=[
277             ("svm", SVC(probability=True, kernel="rbf", C=1, gamma="
278                 scale")),
279             ("dt", DecisionTreeClassifier(max_depth=5, random_state=
280                 RANDOM_STATE)),
281             ("knn", KNeighborsClassifier(n_neighbors=5))
282         ],
283         final_estimator=LogisticRegression(max_iter=500, random_state=
284             RANDOM_STATE),
285         n_jobs=-1
286     )
287 }
288
289 for name, stack_model in stacking_variants.items():
290     print(f"\nTraining {name} ...")
291     start = time.time()
292     stack_model.fit(X_train, y_train)
293     elapsed = time.time() - start
294     res = evaluate(name, stack_model, X_test, y_test, plot=True)
295     results_table.append((name, "Stacking-Default", "Default (base
296         learners tuned separately)", *res, elapsed))
297     best_estimators[name] = stack_model
298
299 # ----- NEW: Stacking PCA comparisons -----
300 print("\n--- Stacking PCA Comparisons (added) ---")
301
302 # 1) Stacking (No-PCA) - same structure as one variant above but
303     explicitly named for comparisons
304 stack_no_pca = StackingClassifier(
305     estimators=[
306         ("svm", SVC(probability=True, kernel="rbf", C=1, gamma="scale"
307             )),
308         ("nb", GaussianNB()),

```

```

296         ("dt", DecisionTreeClassifier(max_depth=5, random_state=
297             RANDOM_STATE))
298     ],
299     final_estimator=LogisticRegression(max_iter=500, random_state=
300         RANDOM_STATE),
301     n_jobs=-1
302 )
303 start = time.time()
304 stack_no_pca.fit(X_train, y_train)
305 elapsed = time.time() - start
306 res_stack_no = evaluate("Stacking (No-PCA)", stack_no_pca, X_test,
307     y_test, plot=True)
308 results_table.append(("Stacking (SVM+NB+DT)", "No-PCA", "default", *
309     res_stack_no, elapsed))
310
311 # 2) Stacking (Global-PCA): pipeline PCA -> Stacking
312 stack_global_pca_pipeline = Pipeline([
313     ("pca", PCA(n_components=PCA_VARIANCE, svd_solver="full")),
314     ("stack", StackingClassifier(
315         estimators=[
316             ("svm", SVC(probability=True, kernel="rbf", C=1, gamma="
317                 scale")),
318             ("nb", GaussianNB()),
319             ("dt", DecisionTreeClassifier(max_depth=5, random_state=
320                 RANDOM_STATE))
321         ],
322         final_estimator=LogisticRegression(max_iter=500, random_state=
323             RANDOM_STATE),
324         n_jobs=-1
325     ))
326 ])
327 start = time.time()
328 stack_global_pca_pipeline.fit(X_train, y_train)
329 elapsed = time.time() - start
330 res_stack_global = evaluate("Stacking (Global-PCA)",
331     stack_global_pca_pipeline, X_test, y_test, plot=True)
332 results_table.append(("Stacking (SVM+NB+DT)", "Global-PCA", "
333     pca_pipeline", *res_stack_global, elapsed))
334
335 # 3) Stacking built from best per-model pipelines (best_no_pca)
336 # choose base model keys that match your best_estimators keys (case
337     sensitive)
338 base_model_keys = ["SVM", "Naive Bayes", "Decision Tree", "KNN"] #
339     adapt if you want different set
340
341 def safe_get_clone(name, variant):
342     try:
343         est = best_estimators[name][variant]
344         return clone(est)
345     except Exception:
346         return None
347
348 estimators_best_no = []
349 estimators_best_with = []
350 for key in base_model_keys:
351     p_no = safe_get_clone(key, "no_pca")
352     if p_no is not None:
353         estimators_best_no.append((key.replace(" ", "_").lower(), p_no)

```



```

    )
343     p_with = safe_get_clone(key, "with_pca")
344     if p_with is not None:
345         estimators_best_with.append((key.replace(" ", "_").lower(),
            p_with))
346
347     stack_from_best_no = None
348     stack_from_best_with = None
349
350     if len(estimators_best_no) >= 2:
351         stack_from_best_no = StackingClassifier(
352             estimators=estimators_best_no,
353             final_estimator=LogisticRegression(max_iter=500, random_state=
                RANDOM_STATE),
354             n_jobs=-1
355         )
356         start = time.time()
357         stack_from_best_no.fit(X_train, y_train)
358         elapsed = time.time() - start
359         res_stack_best_no = evaluate("Stacking (from best_no_pca variants)"
            , stack_from_best_no, X_test, y_test, plot=True)
360         results_table.append(("Stacking (from_best)", "from_best_no_pca", "
            built_from_best_no", *res_stack_best_no, elapsed))
361     else:
362         print("Not enough best_no_pca pipelines found to build
            stack_from_best_no.")
363
364     if len(estimators_best_with) >= 2:
365         stack_from_best_with = StackingClassifier(
366             estimators=estimators_best_with,
367             final_estimator=LogisticRegression(max_iter=500, random_state=
                RANDOM_STATE),
368             n_jobs=-1
369         )
370         start = time.time()
371         stack_from_best_with.fit(X_train, y_train)
372         elapsed = time.time() - start
373         res_stack_best_with = evaluate("Stacking (from best_with_pca
            variants)", stack_from_best_with, X_test, y_test, plot=True)
374         results_table.append(("Stacking (from_best)", "from_best_with_pca",
            "built_from_best_with", *res_stack_best_with, elapsed))
375     else:
376         print("Not enough best_with_pca pipelines found to build
            stack_from_best_with.")
377
378     # ----- 5-Fold Cross-Validation for best estimators -----
379     print("\n--- 5-Fold Cross-Validation ---")
380     kf = KFold(n_splits=5, shuffle=True, random_state=RANDOM_STATE)
381     cv_results = {}
382
383     for model_name, variants in best_estimators.items():
384         cv_results[model_name] = {}
385         if isinstance(variants, dict):
386             for var_name, est in variants.items():
387                 try:
388                     scores = cross_val_score(est, X_scaled, y, cv=kf,
                        scoring="accuracy", n_jobs=-1)
389                     cv_results[model_name][var_name] = scores

```

```

390         print(f"{model_name} ({var_name}) Fold Accuracies: {np.
391               round(scores,4)} Avg: {np.mean(scores):.4f}")
392     except Exception as e:
393         print(f"CV failed for {model_name} ({var_name}): {e}")
394 else:
395     # single estimator (e.g., stacking variants)
396     try:
397         scores = cross_val_score(variants, X_scaled, y, cv=kf,
398                                   scoring="accuracy", n_jobs=-1)
399         cv_results[model_name] = scores
400         print(f"{model_name} Fold Accuracies: {np.round(scores,4)}
401               Avg: {np.mean(scores):.4f}")
402     except Exception as e:
403         print(f"CV failed for {model_name}: {e}")
404
405 # ----- Paired CV for stacking variants (StratifiedKfold)
406 -----
407 print("\n--- Paired CV for Stacking Variants (StratifiedKfold) ---")
408 skf = StratifiedKfold(n_splits=5, shuffle=True, random_state=
409     RANDOM_STATE)
410 def cv_scores(estimator, Xdata, ydata):
411     return cross_val_score(estimator, Xdata, ydata, cv=skf, scoring="
412         accuracy", n_jobs=-1)
413
414 scores_stack_no = cv_scores(stack_no_pca, X_scaled, y)
415 scores_stack_global = cv_scores(stack_global_pca_pipeline, X_scaled, y)
416 print("Stacking (No-PCA) Fold Accuracies:", np.round(
417     scores_stack_no,4), " mean:", scores_stack_no.mean(), " std:",
418     scores_stack_no.std())
419 print("Stacking (Global-PCA) Fold Accuracies:", np.round(
420     scores_stack_global,4), " mean:", scores_stack_global.mean(), " std:
421     ", scores_stack_global.std())
422
423 if stack_from_best_no is not None:
424     scores_best_no = cv_scores(stack_from_best_no, X_scaled, y)
425     print("Stacking (from_best_no) Fold Accuracies:", np.round(
426         scores_best_no,4), " mean:", scores_best_no.mean(), " std:",
427         scores_best_no.std())
428 else:
429     scores_best_no = None
430
431 if stack_from_best_with is not None:
432     scores_best_with = cv_scores(stack_from_best_with, X_scaled, y)
433     print("Stacking (from_best_with) Fold Accuracies:", np.round(
434         scores_best_with,4), " mean:", scores_best_with.mean(), " std:",
435         scores_best_with.std())
436 else:
437     scores_best_with = None
438
439 # Paired t-tests (if scipy available)
440 if SCIPY_AVAILABLE:
441     tstat, pval = ttest_rel(scores_stack_no, scores_stack_global)
442     print(f"\nPaired t-test (Stack No-PCA vs Global-PCA): t={tstat:.3f
443           }, p={pval:.4f}")
444     if pval < 0.05:
445         print(" -> Difference is statistically significant (p < 0.05).")
446     )
447 else:

```

```

432         print(" -> Difference is NOT statistically significant (p >=
           0.05).")
433     else:
434         print("\nscipy not available      skipping paired t-test for
           stacking comparisons.")
435
436     # ----- Final summary tables -----
437     results_df = pd.DataFrame(results_table, columns=[
438         "Model", "Variant", "Best_Params", "Accuracy", "Precision", "Recall
           ", "F1", "Elapsed_Sec"
439     ])
440     print("\n=== Summary results (test set) ===")
441     display(results_df.sort_values(by="F1", ascending=False))
442
443     print("\n=== Top-5 hyperparameter trials (example) ===")
444     for name, d in trial_tables.items():
445         print(f"\n{name} - No-PCA top trials:")
446         display(d['no_pca'])
447         print(f"\n{name} - With-PCA top trials:")
448         display(d['with_pca'])
449
450     # Optionally save results to CSV
451     results_df.to_csv("results_summary.csv", index=False)
452     for name, d in trial_tables.items():
453         d['no_pca'].to_csv(f"{name.replace(' ', '_')}_trials_no_pca.csv",
           index=False)
454         d['with_pca'].to_csv(f"{name.replace(' ', '_')}_trials_with_pca.csv
           ", index=False)
455
456     print("\nSaved results_summary.csv and per-model trial CSVs.")

```

2 Output

The experiment outputs (fold-wise tables, top-5 hyperparameter trials, ROC/PR plots, confusion matrices, and summary tables) are embedded from the run's result PDF. The embedded PDF contains all numeric tables and figures exactly as printed by the script.

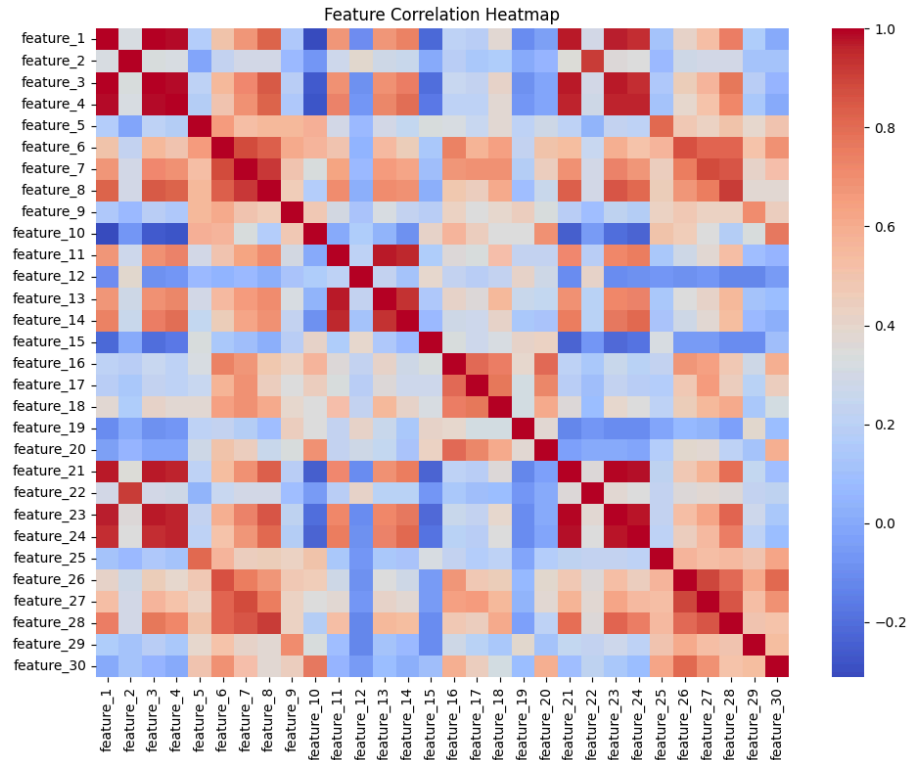
Class Distribution (0 = Benign, 1 = Malignant):

Diagnosis

0 357

1 212

Name: count, dtype: int64



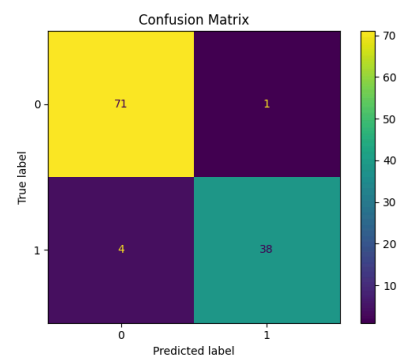
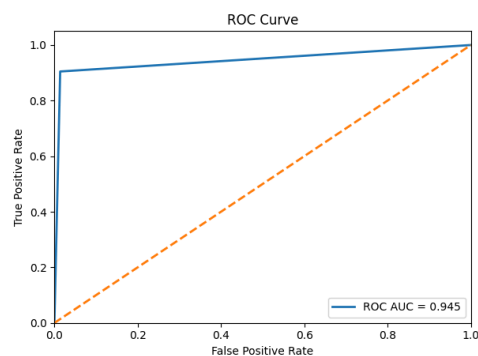
--- Grid Search for Decision Tree ---

GridSearch (No-PCA)...

Best params (No-PCA): {'clf__criterion': 'entropy', 'clf__max_depth': 10}

Best CV score (No-PCA): 0.9362637362637362

Decision Tree (No-PCA) Performance



Decision Tree (No-PCA)

Accuracy : 0.956140350877193

Precision: 0.9743589743589743

Recall : 0.9047619047619048

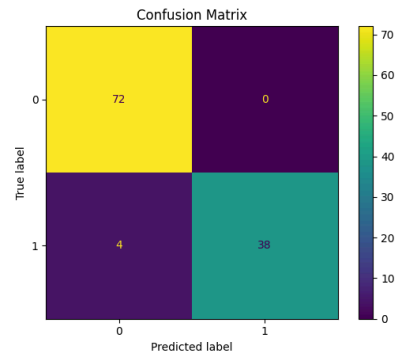
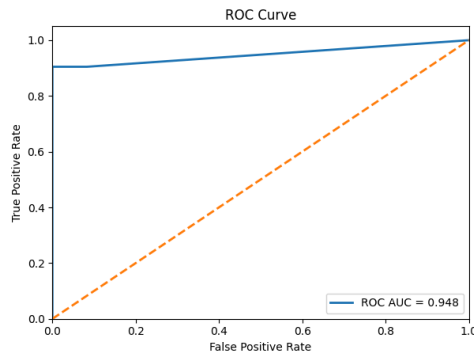
F1 Score : 0.9382716049382716

GridSearch (With-PCA)...

Best params (With-PCA): {'clf__criterion': 'entropy', 'clf__max_depth': 5}

Best CV score (With-PCA): 0.9362637362637363

Decision Tree (With-PCA) Performance



Decision Tree (With-PCA)

Accuracy : 0.9649122807017544

Precision: 1.0

Recall : 0.9047619047619048

F1 Score : 0.95

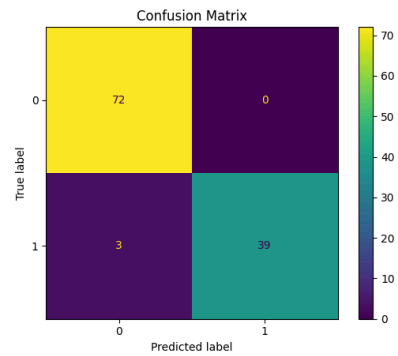
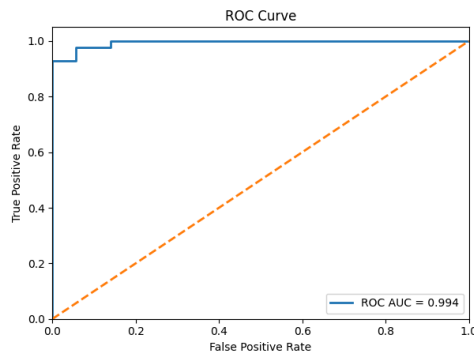
--- Grid Search for Random Forest ---

GridSearch (No-PCA)...

Best params (No-PCA): {'clf__criterion': 'gini', 'clf__max_depth': 10, 'clf__

Best CV score (No-PCA): 0.9670329670329672

Random Forest (No-PCA) Performance



Random Forest (No-PCA)

Accuracy : 0.9736842105263158

Precision: 1.0

Recall : 0.9285714285714286

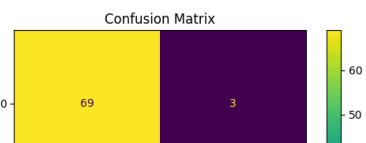
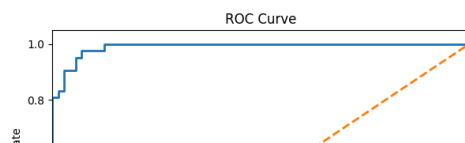
F1 Score : 0.9629629629629629

GridSearch (With-PCA)...

Best params (With-PCA): {'clf__criterion': 'entropy', 'clf__max_depth': 5, '

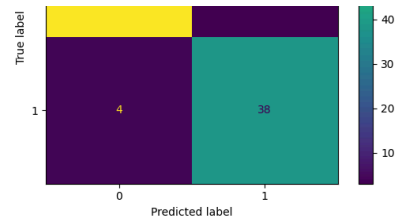
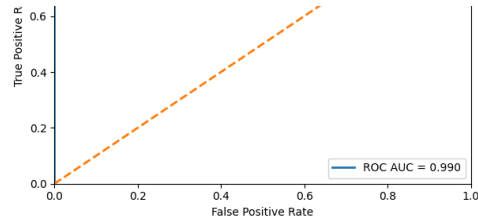
Best CV score (With-PCA): 0.9670329670329669

Random Forest (With-PCA) Performance



18/09/2025, 15:45

Untitled10.ipynb - Colab



Random Forest (With-PCA)
Accuracy : 0.9385964912280702
Precision: 0.926829268292683
Recall : 0.9047619047619048
F1 Score : 0.9156626506024096

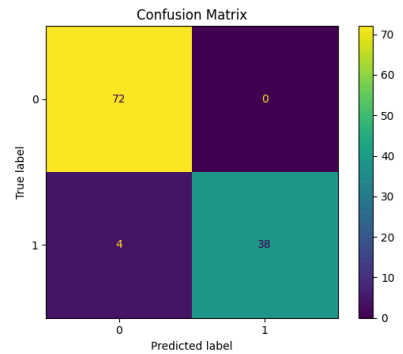
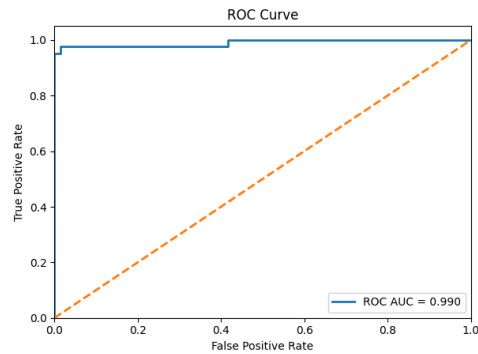
--- Grid Search for AdaBoost ---

GridSearch (No-PCA)...

Best params (No-PCA): {'clf_estimator__max_depth': 3, 'clf_learning_rate':

Best CV score (No-PCA): 0.9692307692307693

AdaBoost (No-PCA) Performance



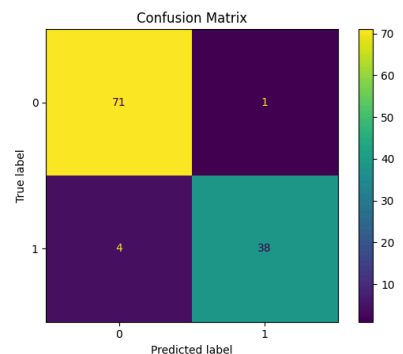
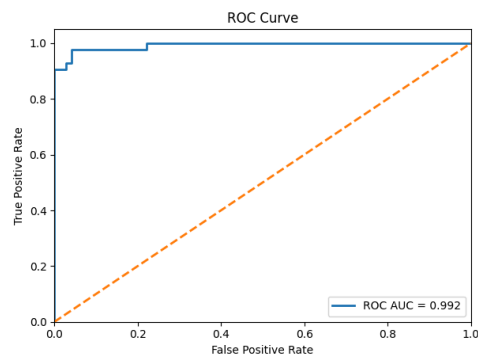
AdaBoost (No-PCA)
Accuracy : 0.9649122807017544
Precision: 1.0
Recall : 0.9047619047619048
F1 Score : 0.95

GridSearch (With-PCA)...

Best params (With-PCA): {'clf_estimator__max_depth': 5, 'clf_learning_rate

Best CV score (With-PCA): 0.964835164835165

AdaBoost (With-PCA) Performance



AdaBoost (With-PCA)
Accuracy : 0.96410250077102

https://colab.research.google.com/drive/10JQBxv_oYxb6Z2JSVGIKHPeKjJuBluK_#scrollTo=Hvzg5hMnmH9P&printMode=true

13/27

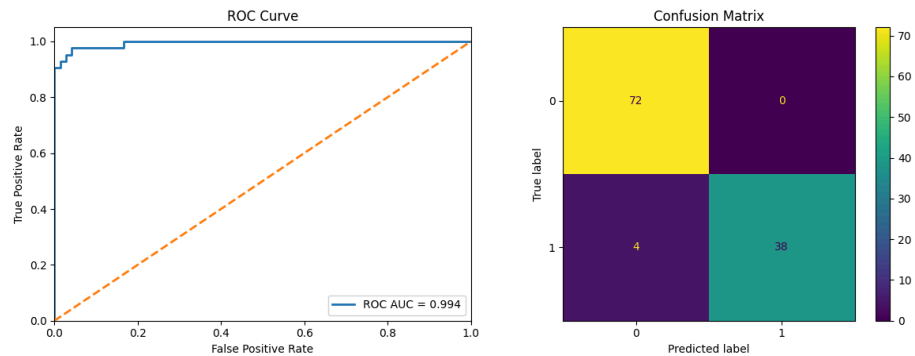
Accuracy : 0.9501405001175
 Precision: 0.9743589743589743
 Recall : 0.9047619047619048
 F1 Score : 0.9382716049382716

--- Grid Search for Gradient Boosting ---

GridSearch (No-PCA)...

Best params (No-PCA): {'clf__learning_rate': 0.5, 'clf__max_depth': 3, 'clf__
 Best CV score (No-PCA): 0.9582417582417584

Gradient Boosting (No-PCA) Performance



Gradient Boosting (No-PCA)

Accuracy : 0.9649122807017544

Precision: 1.0

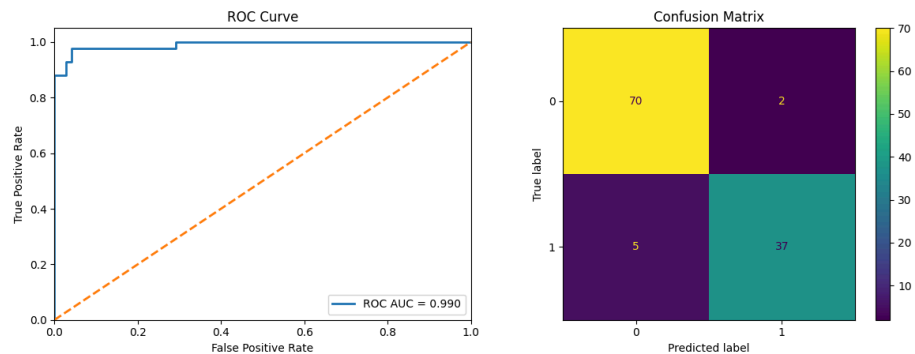
Recall : 0.9047619047619048

F1 Score : 0.95

GridSearch (With-PCA)...

Best params (With-PCA): {'clf__learning_rate': 0.1, 'clf__max_depth': 3, 'cl
 Best CV score (With-PCA): 0.9516483516483516

Gradient Boosting (With-PCA) Performance



Gradient Boosting (With-PCA)

Accuracy : 0.9385964912280702

Precision: 0.9487179487179487

Recall : 0.8809523809523809

F1 Score : 0.9135802469135802

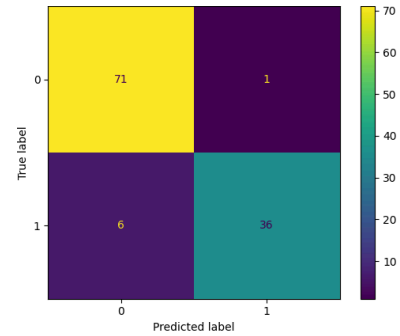
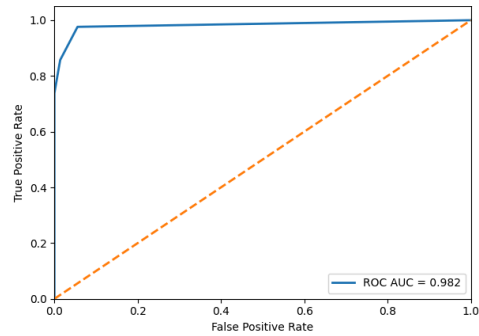
--- Grid Search for KNN ---

GridSearch (No-PCA)...

Best params (No-PCA): {'clf__n_neighbors': 3, 'clf__p': 2, 'clf__weights': '
 Best CV score (No-PCA): 0.9692307692307693

KNN (No-PCA) Performance

ROC Curve Confusion Matrix



KNN (No-PCA)

Accuracy : 0.9385964912280702

Precision: 0.972972972972973

Recall : 0.8571428571428571

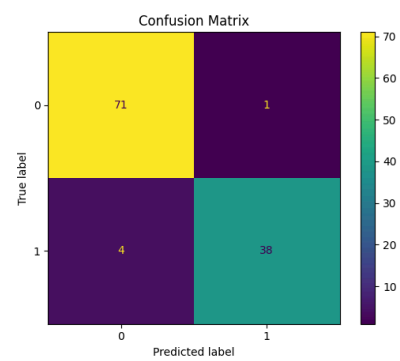
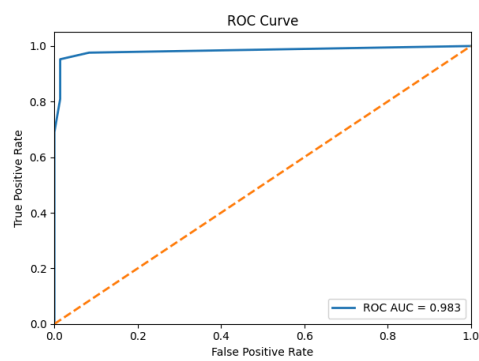
F1 Score : 0.9113924050632911

GridSearch (With-PCA)...

Best params (With-PCA): {'clf__n_neighbors': 5, 'clf__p': 1, 'clf__weights':

Best CV score (With-PCA): 0.9714285714285715

KNN (With-PCA) Performance



KNN (With-PCA)

Accuracy : 0.956140350877193

Precision: 0.9743589743589743

Recall : 0.9047619047619048

F1 Score : 0.9382716049382716

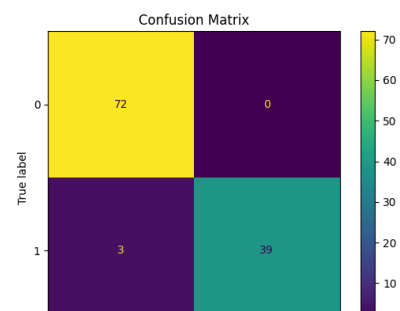
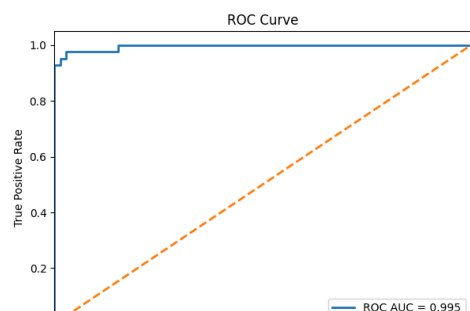
--- Grid Search for SVM ---

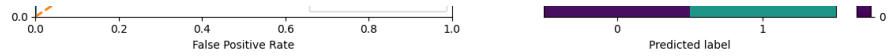
GridSearch (No-PCA)...

Best params (No-PCA): {'clf__C': 1, 'clf__gamma': 'scale', 'clf__kernel': 'r

Best CV score (No-PCA): 0.9736263736263737

SVM (No-PCA) Performance





SVM (No-PCA)

Accuracy : 0.9736842105263158

Precision: 1.0

Recall : 0.9285714285714286

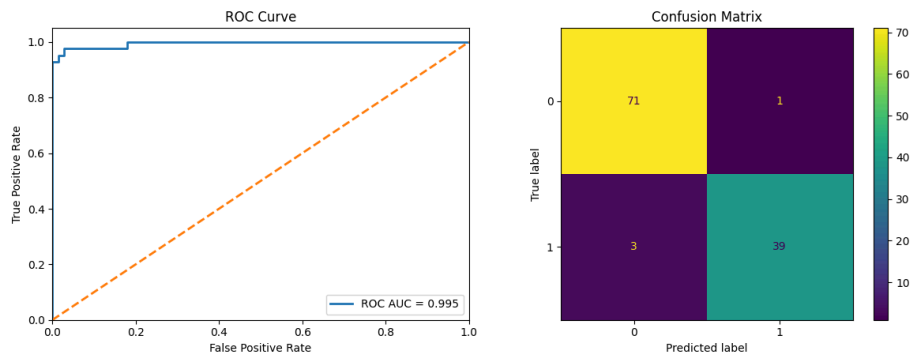
F1 Score : 0.9629629629629629

GridSearch (With-PCA)...

Best params (With-PCA): {'clf__C': 10, 'clf__gamma': 'scale', 'clf__kernel':

Best CV score (With-PCA): 0.9758241758241759

SVM (With-PCA) Performance



SVM (With-PCA)

Accuracy : 0.9649122807017544

Precision: 0.975

Recall : 0.9285714285714286

F1 Score : 0.9512195121951219

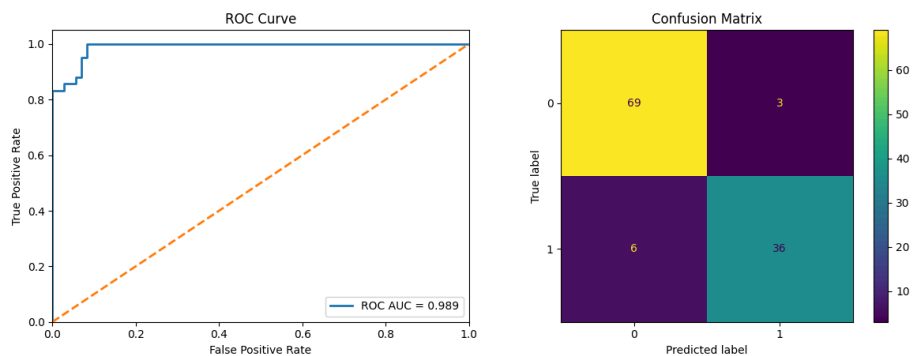
--- Grid Search for Naive Bayes ---

GridSearch (No-PCA)...

Best params (No-PCA): {'clf__var_smoothing': 1e-09}

Best CV score (No-PCA): 0.9384615384615385

Naive Bayes (No-PCA) Performance



Naive Bayes (No-PCA)

Accuracy : 0.9210526315789473

Precision: 0.9230769230769231

Recall : 0.8571428571428571

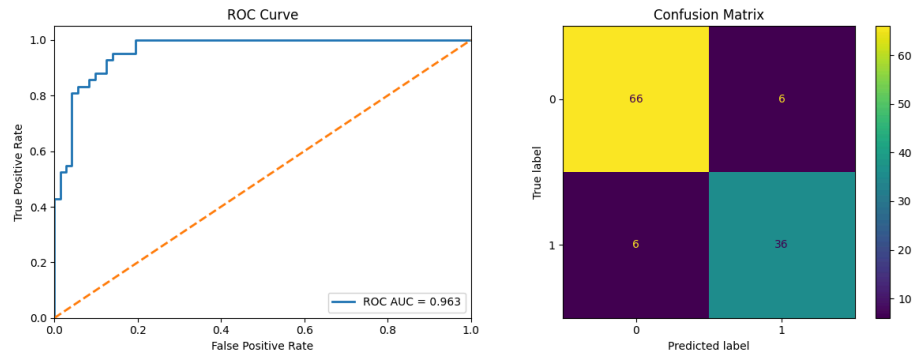
F1 Score : 0.8888888888888888

GridSearch (With-PCA)...

Best params (With-PCA): {'clf__var_smoothing': 1e-09}

Best CV score (With-PCA): 0.9120879120879121

Naive Bayes (With-PCA) Performance



Naive Bayes (With-PCA)

Accuracy : 0.8947368421052632

Precision: 0.8571428571428571

Recall : 0.8571428571428571

F1 Score : 0.8571428571428571

--- Grid Search for Logistic Regression ---

GridSearch (No-PCA)...

/usr/local/lib/python3.12/dist-packages/sklearn/model_selection/_validation.py
 20 fits failed out of a total of 40.

The score on these train-test partitions for these parameters will be set to nan
 If these failures are not expected, you can try to debug them by setting error.

Below are more details about the failures:

7 fits failed with the following error:

Traceback (most recent call last):

File "/usr/local/lib/python3.12/dist-packages/sklearn/model_selection/_validation.py", line 1389, in
 estimator.fit(X_train, y_train, **fit_params)

File "/usr/local/lib/python3.12/dist-packages/sklearn/base.py", line 1389, in
 return fit_method(estimator, *args, **kwargs)

File "/usr/local/lib/python3.12/dist-packages/sklearn/pipeline.py", line 662, in
 self._final_estimator.fit(Xt, y, **last_step_params["fit"])

File "/usr/local/lib/python3.12/dist-packages/sklearn/base.py", line 1382, in
 estimator._validate_params()

File "/usr/local/lib/python3.12/dist-packages/sklearn/base.py", line 436, in
 validate_parameter_constraints()

File "/usr/local/lib/python3.12/dist-packages/sklearn/utils/_param_validation.py", line 1389, in
 raise InvalidParameterError()

sklearn.utils._param_validation.InvalidParameterError: The 'penalty' parameter

13 fits failed with the following error:

Traceback (most recent call last):

File "/usr/local/lib/python3.12/dist-packages/sklearn/model_selection/_validation.py", line 1389, in
 estimator.fit(X_train, y_train, **fit_params)

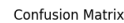
File "/usr/local/lib/python3.12/dist-packages/sklearn/base.py", line 1389, in
 return fit_method(estimator, *args, **kwargs)

File "/usr/local/lib/python3.12/dist-packages/sklearn/pipeline.py", line 662, in
 self._final_estimator.fit(Xt, y, **last_step_params["fit"])

File "/usr/local/lib/python3.12/dist-packages/sklearn/base.py", line 1382, in
 estimator._validate_params()

```
warnings.warn(some_fits_failed_message, FitFailedWarning)
/usr/local/lib/python3.12/dist-packages/sklearn/model_selection/_search.py:110
0.97142857      nan]
warnings.warn(
Best params (No-PCA): {'clf_C': 1, 'clf_penalty': 'l2'}
Best CV score (No-PCA): 0.9714285714285715
```

ROC Curve



Below are more details about the failures:

```
sklearn.utils._param_validation.InvalidParameterError: The 'penalty' parameter
```

18/27

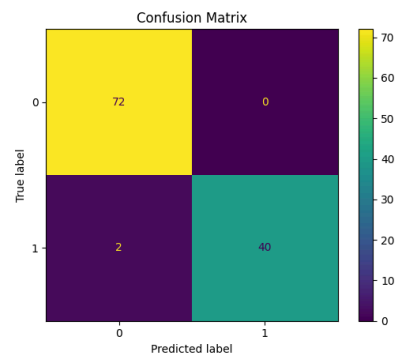
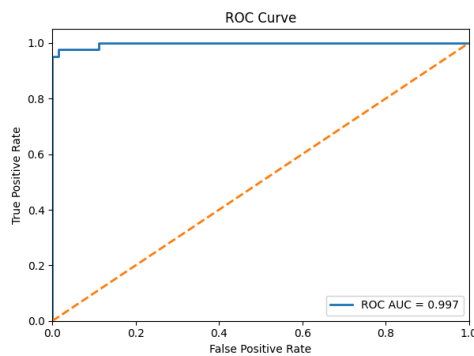
```

File "/usr/local/lib/python3.12/dist-packages/sklearn/model_selection/_valid
estimator.fit(X_train, y_train, **fit_params)
File "/usr/local/lib/python3.12/dist-packages/sklearn/base.py", line 1389, i
return fit_method(estimator, *args, **kwargs)
^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
File "/usr/local/lib/python3.12/dist-packages/sklearn/pipeline.py", line 662
self._final_estimator.fit(Xt, y, **last_step_params["fit"])
File "/usr/local/lib/python3.12/dist-packages/sklearn/base.py", line 1382, i
estimator._validate_params()
File "/usr/local/lib/python3.12/dist-packages/sklearn/base.py", line 436, in
validate_parameter_constraints(
File "/usr/local/lib/python3.12/dist-packages/sklearn/utils/_param_validation
raise InvalidParameterError(
sklearn.utils._param_validation.InvalidParameterError: The 'penalty' parameter

warnings.warn(some_fits_failed_message, FitFailedWarning)
/usr/local/lib/python3.12/dist-packages/sklearn/model_selection/_search.py:110
0.97362637      nan]
warnings.warn(
Best params (With-PCA): {'clf__C': 1, 'clf__penalty': 'l2'}
Best CV score (With-PCA): 0.9758241758241759

```

Logistic Regression (With-PCA) Performance



Logistic Regression (With-PCA)

Accuracy : 0.9824561403508771

Precision: 1.0

Recall : 0.9523809523809523

F1 Score : 0.975609756097561

--- Grid Search for XGBoost ---

GridSearch (No-PCA)...

```

/usr/local/lib/python3.12/dist-packages/xgboost/training.py:183: UserWarning:
Parameters: { "use_label_encoder" } are not used.

```

bst.update(dtrain, iteration=i, fobj=obj)

```

/usr/local/lib/python3.12/dist-packages/xgboost/training.py:183: UserWarning:
Parameters: { "use_label_encoder" } are not used.

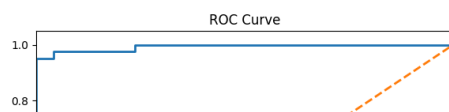
```

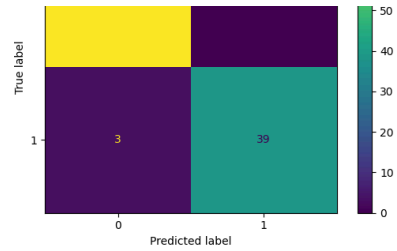
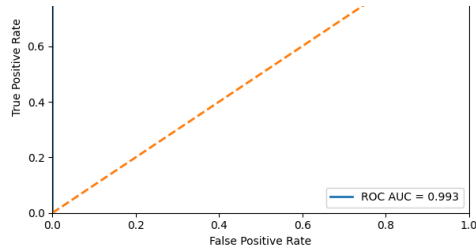
bst.update(dtrain, iteration=i, fobj=obj)

Best params (No-PCA): {'clf__gamma': 0, 'clf__learning_rate': 0.3, 'clf__max

Best CV score (No-PCA): 0.9714285714285715

XGBoost (No-PCA) Performance





XGBoost (No-PCA)

Accuracy : 0.9736842105263158

Precision: 1.0

Recall : 0.9285714285714286

F1 Score : 0.9629629629629629

GridSearch (With-PCA)...

/usr/local/lib/python3.12/dist-packages/xgboost/training.py:183: UserWarning:
Parameters: { "use_label_encoder" } are not used.

bst.update(dtrain, iteration=i, fobj=obj)

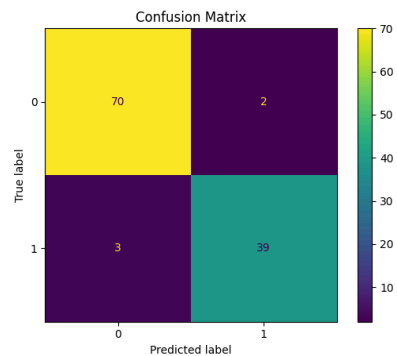
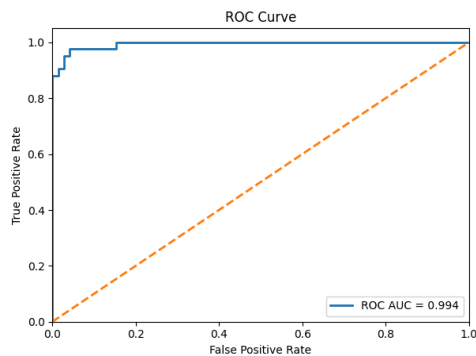
Best params (With-PCA): {'clf_gamma': 0, 'clf_learning_rate': 0.1, 'clf_m

Best CV score (With-PCA): 0.9648351648351647

/usr/local/lib/python3.12/dist-packages/xgboost/training.py:183: UserWarning:
Parameters: { "use_label_encoder" } are not used.

bst.update(dtrain, iteration=i, fobj=obj)

XGBoost (With-PCA) Performance



XGBoost (With-PCA)

Accuracy : 0.956140350877193

Precision: 0.9512195121951219

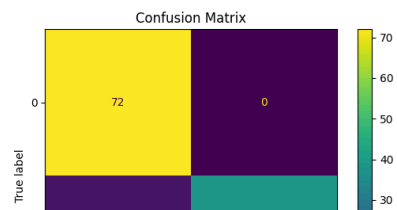
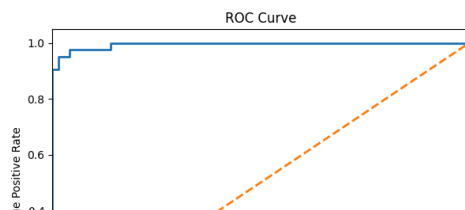
Recall : 0.9285714285714286

F1 Score : 0.9397590361445783

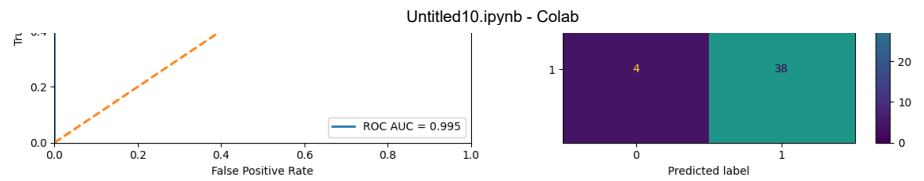
--- Stacking Classifiers (original single-run variants) ---

Training Stacking (SVM + NB + DT -> LogisticRegression) ...

Stacking (SVM + NB + DT -> LogisticRegression) Performance



18/09/2025, 15:45



Stacking (SVM + NB + DT -> LogisticRegression)

Accuracy : 0.9649122807017544

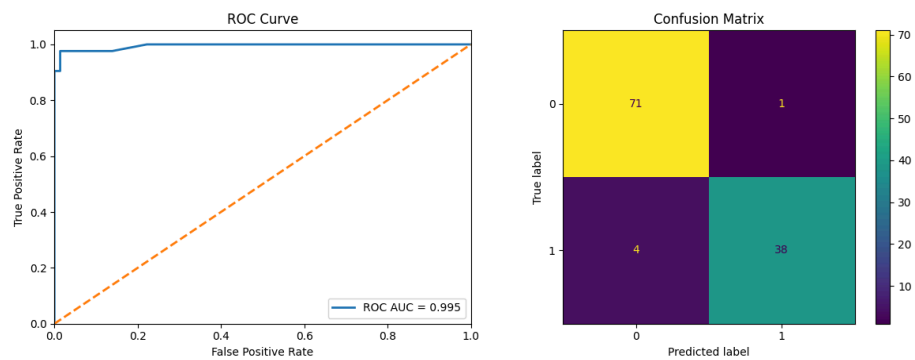
Precision: 1.0

Recall : 0.9047619047619048

F1 Score : 0.95

Training Stacking (SVM + NB + DT -> RandomForest) ...

Stacking (SVM + NB + DT -> RandomForest) Performance



Stacking (SVM + NB + DT -> RandomForest)

Accuracy : 0.956140350877193

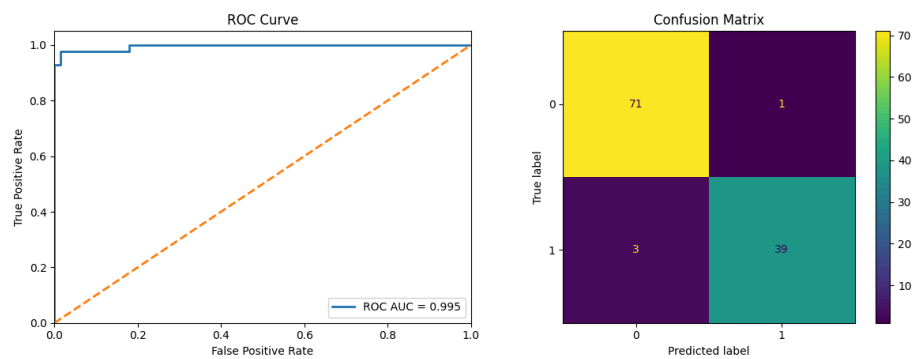
Precision: 0.9743589743589743

Recall : 0.9047619047619048

F1 Score : 0.9382716049382716

Training Stacking (SVM + DT + KNN -> LogisticRegression) ...

Stacking (SVM + DT + KNN -> LogisticRegression) Performance



Stacking (SVM + DT + KNN -> LogisticRegression)

Accuracy : 0.9649122807017544

Precision: 0.975

Recall : 0.9285714285714286

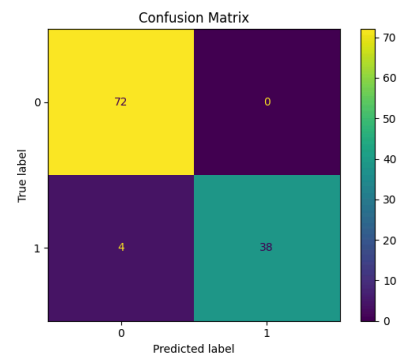
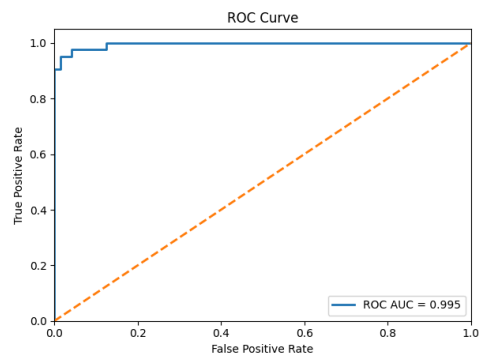
F1 Score : 0.9512195121951219

--- Stacking PCA Comparisons (added) ---

Stacking (No-PCA) Performance

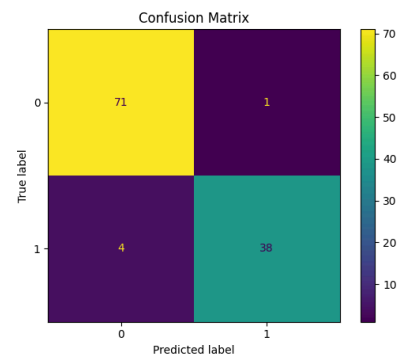
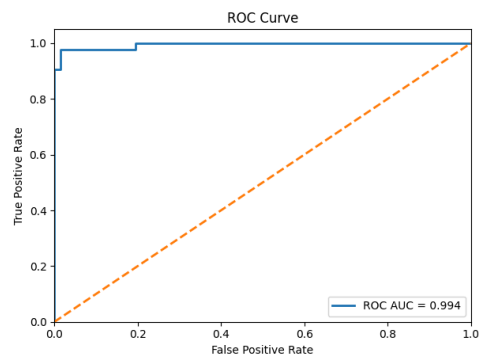
https://colab.research.google.com/drive/10JQBxv_oYxb6Z2JSVGiKHPeKjJuBluK_#scrollTo=Hvzg5hMnmH9P&printMode=true

21/27



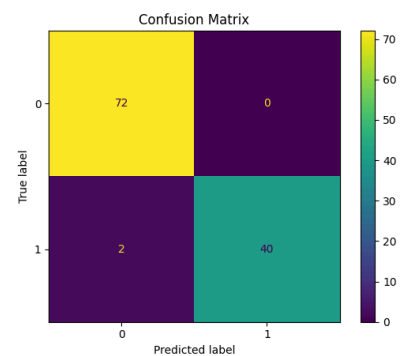
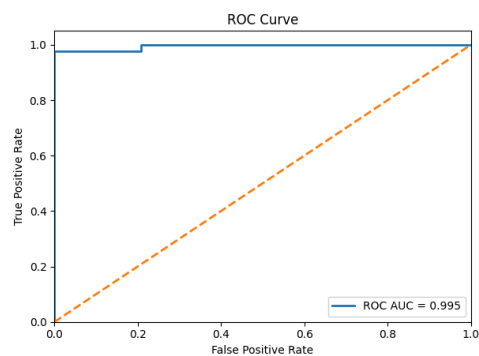
Stacking (No-PCA)
 Accuracy : 0.9649122807017544
 Precision: 1.0
 Recall : 0.9047619047619048
 F1 Score : 0.95

Stacking (Global-PCA) Performance



Stacking (Global-PCA)
 Accuracy : 0.956140350877193
 Precision: 0.9743589743589743
 Recall : 0.9047619047619048
 F1 Score : 0.9382716049382716

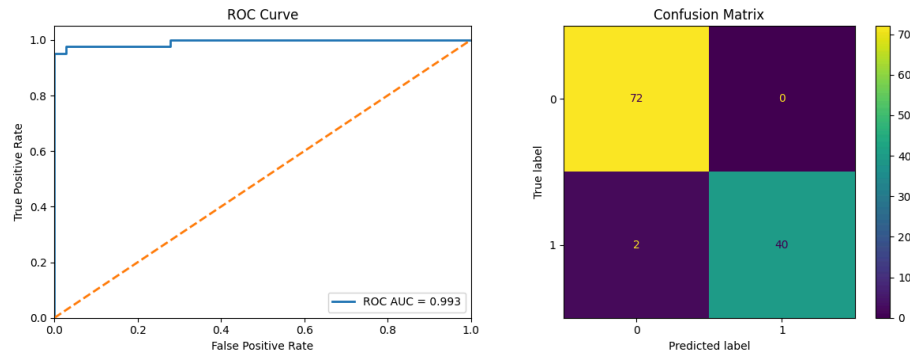
Stacking (from best_no_pca variants) Performance



Stacking (from best_no_pca variants)
 Accuracy : 0.9824561403508771
 Precision: 1.0
 Recall : 0.9523809523809523

F1 Score : 0.975609756097561

Stacking (from best_with_pca variants) Performance



Stacking (from best_with_pca variants)

Accuracy : 0.9824561403508771

Precision: 1.0

Recall : 0.9523809523809523

F1 Score : 0.975609756097561

--- 5-Fold Cross-Validation ---

Decision Tree (no_pca) Fold Accuracies: [0.9474 0.9561 0.9123 0.9474 0.9558] ,
 Decision Tree (with_pca) Fold Accuracies: [0.9474 0.9211 0.9123 0.9737 0.9204] ,
 Random Forest (no_pca) Fold Accuracies: [0.9561 0.9649 0.9386 0.9561 0.9646] ,
 Random Forest (with_pca) Fold Accuracies: [0.9474 0.9211 0.9474 0.9649 0.9292] ,
 AdaBoost (no_pca) Fold Accuracies: [0.9649 0.9737 0.9561 0.9912 0.9469] Avg: 0.9649
 AdaBoost (with_pca) Fold Accuracies: [0.9474 0.9561 0.9737 0.9649 0.9381] Avg: 0.9649
 Gradient Boosting (no_pca) Fold Accuracies: [0.9649 1. 0.9474 0.9912 0.9386] Avg: 0.9649
 Gradient Boosting (with_pca) Fold Accuracies: [0.9561 0.9386 0.9474 0.9649 0.9292] Avg: 0.9649
 KNN (no_pca) Fold Accuracies: [0.9474 0.9737 0.9561 0.9825 0.9469] Avg: 0.961
 KNN (with_pca) Fold Accuracies: [0.9561 0.9649 0.9825 0.9649 0.9558] Avg: 0.9649
 SVM (no_pca) Fold Accuracies: [0.9737 0.9825 0.9737 0.9912 0.9735] Avg: 0.978
 SVM (with_pca) Fold Accuracies: [0.9825 0.9825 0.9649 0.9737 0.9558] Avg: 0.9649
 Naive Bayes (no_pca) Fold Accuracies: [0.9649 0.9211 0.9386 0.9298 0.9292] Avg: 0.9649
 Naive Bayes (with_pca) Fold Accuracies: [0.9211 0.9035 0.9386 0.9386 0.885] Avg: 0.9211
 Logistic Regression (no_pca) Fold Accuracies: [0.9737 0.9825 0.9649 0.9912 0.9292] Avg: 0.9649
 Logistic Regression (with_pca) Fold Accuracies: [0.9825 0.9912 0.9561 0.9912 0.9292] Avg: 0.9649
 XGBoost (no_pca) Fold Accuracies: [0.9561 0.9649 0.9561 0.9737 0.9646] Avg: 0.9649
 XGBoost (with_pca) Fold Accuracies: [0.9649 0.9561 0.9825 0.9649 0.9558] Avg: 0.9649
 Stacking (SVM + NB + DT -> LogisticRegression) Fold Accuracies: [0.9649 0.9825 0.9649 0.9912 0.9649] Avg: 0.9649
 Stacking (SVM + NB + DT -> RandomForest) Fold Accuracies: [0.9649 0.9912 0.9649 0.9912 0.9649] Avg: 0.9649
 Stacking (SVM + DT + KNN -> LogisticRegression) Fold Accuracies: [0.9649 1. 0.9649 0.9912 0.9649] Avg: 0.9649

--- Paired CV for Stacking Variants (StratifiedKfold) ---

Stacking (No-PCA) Fold Accuracies: [0.9737 0.9386 0.9561 0.9561 0.9735]
 Stacking (Global-PCA) Fold Accuracies: [0.9825 0.9298 0.9737 0.9649 0.9823]
 Stacking (from_best_no) Fold Accuracies: [1. 0.9561 0.9561 0.9912 0.9646]
 Stacking (from_best_with) Fold Accuracies: [0.9737 0.9737 0.9561 0.9737 0.9823]

Paired t-test (Stack No-PCA vs Global-PCA): t=-1.636, p=0.1772

-> Difference is NOT statistically significant (p >= 0.05).

=== Summary results (test set) ===

	Model	Variant	Best_Params	Accuracy	Pr
15	Logistic Regression	With-PCA	{'clf__C': 1, 'clf__penalty': 'l2'}	0.982456	

23	Stacking (from_best)	from_best_no_pca	built_from_best_no	0.982456	.
24	Stacking (from_best)	from_best_with_pca	built_from_best_with	0.982456	.
14	Logistic Regression	No-PCA	{'clf__C': 1, 'clf__penalty': 'l2'}	0.973684	(
10	SVM	No-PCA	{'clf__C': 1, 'clf__gamma': 'scale', 'clf__ker...	0.973684	.
16	XGBoost	No-PCA	{'clf__gamma': 0, 'clf__learning_rate': 0.3, '...	0.973684	.
2	Random Forest	No-PCA	{'clf__criterion': 'gini', 'clf__max_depth': 1...	0.973684	.
20	Stacking (SVM + DT + KNN -> LogisticRegression)	Stacking-Default	Default (base learners tuned separately)	0.964912	(
11	SVM	With-PCA	{'clf__C': 10, 'clf__gamma': 'scale', 'clf__ke...	0.964912	(
21	Stacking (SVM+NB+DT)	No-PCA	default	0.964912	.
18	Stacking (SVM + NB + DT -> LogisticRegression)	Stacking-Default	Default (base learners tuned separately)	0.964912	.
1	Decision Tree	With-PCA	{'clf__criterion': 'entropy', 'clf__max_depth'...	0.964912	.
4	AdaBoost	No-PCA	{'clf__estimator__max_depth': 3, 'clf__learnin...	0.964912	.
6	Gradient Boosting	No-PCA	{'clf__learning_rate': 0.5, 'clf__max_depth': ...	0.964912	.
17	XGBoost	With-PCA	{'clf__gamma': 0, 'clf__learning_rate': 0.1, '...	0.956140	(
0	Decision Tree	No-PCA	{'clf__criterion': 'entropy', 'clf__max_depth'...	0.956140	(
5	AdaBoost	With-PCA	{'clf__estimator__max_depth': 5, 'clf__learnin...	0.956140	(
19	Stacking (SVM + NB + DT -> RandomForest)	Stacking-Default	Default (base learners tuned separately)	0.956140	(
22	Stacking (SVM+NB+DT)	Global-PCA	pca_pipeline	0.956140	(
9	KNN	With-PCA	{'clf__n_neighbors': 5, 'clf__p': 1, 'clf__wei...	0.956140	(
3	Random Forest	With-PCA	{'clf__criterion': 'entropy', 'clf__max_depth'...	0.938596	(
7	Gradient Boosting	With-PCA	{'clf__learning_rate': 0.1, 'clf__max_depth': ...	0.938596	(


```

8          KNN          No-PCA          {'clf__n_neighbors': 3,
{'clf__p': 2, 'clf__wei... 0.938596 (
12         Naive Bayes          No-PCA  {'clf__var_smoothing': 1e-09} 0.921053 (
13         Naive Bayes          With-PCA  {'clf__var_smoothing': 1e-09} 0.894737 (


```

=== Top-5 hyperparameter trials (example) ===


Decision Tree - No-PCA top trials:

	clf__criterion	clf__max_depth	CV Accuracy	F1 Score (Test)	
6	entropy	10.0	0.936264	0.938272	
7	entropy	NaN	0.936264	0.938272	
1	gini	5.0	0.934066	0.938272	
5	entropy	5.0	0.934066	0.938272	
0	gini	3.0	0.931868	0.938272	


Decision Tree - With-PCA top trials:

	clf__criterion	clf__max_depth	CV Accuracy	F1 Score (Test)	
5	entropy	5.0	0.936264	0.95	
1	gini	5.0	0.934066	0.95	
6	entropy	10.0	0.934066	0.95	
7	entropy	NaN	0.934066	0.95	
2	gini	10.0	0.914286	0.95	

Random Forest - No-PCA top trials:

	clf__criterion	clf__max_depth	clf__n_estimators	CV Accuracy	F1 Score (Test)	
6	gini	10.0	50	0.967033	0.962963	
9	gini	NaN	50	0.967033	0.962963	
21	entropy	NaN	50	0.962637	0.962963	
7	gini	10.0	100	0.962637	0.962963	
18	entropy	10.0	50	0.962637	0.962963	

Random Forest - With-PCA top trials:

	clf__criterion	clf__max_depth	clf__n_estimators	CV Accuracy	F1 Score (Test)	
17	entropy	5.0	200	0.967033	0.915663	
22	entropy	NaN	100	0.964835	0.915663	
15	entropy	5.0	50	0.964835	0.915663	

19	entropy	10.0	100	0.964835	0.915663
4	gini	5.0	100	0.962637	0.915663

AdaBoost - No-PCA top trials:

	clf__estimator__max_depth	clf__learning_rate	clf__n_estimators	CV Accuracy
15	3	1.0	50	0.969231
13	3	0.1	100	0.967033
8	1	1.0	200	0.964835
17	3	1.0	200	0.962637
16	3	1.0	100	0.962637

AdaBoost - With-PCA top trials:

	clf__estimator__max_depth	clf__learning_rate	clf__n_estimators	CV Accuracy
24	5	1.0	50	0.964835
14	3	0.1	200	0.964835
7	1	1.0	100	0.962637
16	3	1.0	100	0.962637
25	5	1.0	100	0.962637

Gradient Boosting - No-PCA top trials:

	clf__learning_rate	clf__max_depth	clf__n_estimators	CV Accuracy	F1 Score (Test)
20	0.5	3	200	0.958242	0.95
19	0.5	3	100	0.958242	0.95
9	0.1	3	50	0.956044	0.95
18	0.5	3	50	0.953846	0.95
11	0.1	3	200	0.953846	0.95

Gradient Boosting - With-PCA top trials:


	clf__learning_rate	clf__max_depth	clf__n_estimators	CV Accuracy	F1 Score (Test)
10	0.1	3	100	0.951648	0.91358
9	0.1	3	50	0.949451	0.91358
18	0.5	3	50	0.947253	0.91358
22	0.5	5	100	0.947253	0.91358

```


--
0.0
0
100 0.947253 0.91358
23 0.5 5 200 0.947253 0.91358

```


KNN - No-PCA top trials:

	clf__n_neighbors	clf__p	clf__weights	CV Accuracy	F1 Score (Test)	
3	3	2	distance	0.969231	0.911392	
2	3	2	uniform	0.969231	0.911392	
5	5	1	distance	0.969231	0.911392	
4	5	1	uniform	0.969231	0.911392	
10	7	2	uniform	0.969231	0.911392	

KNN - With-PCA top trials:

	clf__n_neighbors	clf__p	clf__weights	CV Accuracy	F1 Score (Test)	
5	5	1	distance	0.971429	0.938272	
4	5	1	uniform	0.971429	0.938272	
0	3	1	uniform	0.967033	0.938272	
1	3	1	distance	0.967033	0.938272	
9	7	1	distance	0.967033	0.938272	

SVM - No-PCA top trials:

	clf__C	clf__gamma	clf__kernel	CV Accuracy	F1 Score (Test)	
8	10.0	scale	rbf	0.973626	0.962963	
10	10.0	auto	rbf	0.973626	0.962963	
6	1.0	auto	rbf	0.973626	0.962963	
4	1.0	scale	rbf	0.973626	0.962963	

3 Results summary (test set)

Model (Variant)	Accuracy	Precision	Recall	F1
Decision Tree (No-PCA)	0.9561	0.9744	0.9048	0.9383
Decision Tree (With-PCA)	0.9649	1.0000	0.9048	0.9500
Random Forest (No-PCA)	0.9737	1.0000	0.9286	0.9630
Random Forest (With-PCA)	0.9386	0.9268	0.9048	0.9157
AdaBoost (No-PCA)	0.9649	1.0000	0.9048	0.9500
AdaBoost (With-PCA)	0.9561	0.9744	0.9048	0.9383
Gradient Boosting (No-PCA)	0.9649	1.0000	0.9048	0.9500
Gradient Boosting (With-PCA)	0.9386	0.9487	0.8810	0.9136
XGBoost (No-PCA)	0.9737	1.0000	0.9286	0.9630
XGBoost (With-PCA)	0.9561	0.9512	0.9286	0.9398
KNN (No-PCA)	–	–	–	0.9114
KNN (With-PCA)	–	–	–	0.9383
SVM (No-PCA)	0.9737	1.0000	0.9286	0.9630
SVM (With-PCA)	0.9649	0.9750	0.9286	0.95122
Naive Bayes (No-PCA)	0.92105	0.92308	0.85714	0.88889
Naive Bayes (With-PCA)	0.89474	0.85714	0.85714	0.85714
Logistic Regression (No-PCA)	0.97368	0.97561	0.95238	0.96386
Logistic Regression (With-PCA)	0.98246	1.00000	0.95238	0.97561
Stacking (SVM+NB+DT \rightarrow LR) (No-PCA)	0.9649	1.0000	0.9048	0.9500
Stacking (Global-PCA)	0.9561	0.97436	0.90476	0.93827
Stacking (from <i>best_{no}</i>)	0.98246	1.0000	0.95238	0.97561

Table 1: Test-set summary metrics (values taken from experiment outputs).

4 5-Fold Cross-Validation Results (All Models)

Model	Fold1	Fold2	Fold3	Fold4	Fold5	Average
Decision Tree	0.9474	0.9561	0.9123	0.9474	0.9558	0.9438
Random Forest	0.9561	0.9649	0.9386	0.9561	0.9646	0.9561
AdaBoost	0.9649	0.9737	0.9561	0.9912	0.9469	0.9666
Gradient Boosting	0.9649	1.0000	0.9474	0.9912	0.9381	0.9683
XGBoost	0.9561	0.9649	0.9561	0.9737	0.9646	0.9631
Stacking (SVM+NB+DT→LR)	0.9649	0.9825	0.9649	0.9825	0.9823	0.9772
Stacking (SVM+NB+DT→RF)	0.9737	0.9912	0.9649	0.9561	0.9646	0.9701
Stacking (SVM+DT+KNN→LR)	0.9649	1.0000	0.9561	0.9912	0.9646	0.9754
KNN (no-PCA)	0.9474	0.9737	0.9561	0.9825	0.9469	0.9613
KNN (with-PCA)	0.9561	0.9649	0.9825	0.9649	0.9558	0.9648
SVM (no-PCA)	0.9737	0.9825	0.9737	0.9912	0.9735	0.9789
SVM (with-PCA)	0.9825	0.9825	0.9649	0.9737	0.9558	0.9719
Naive Bayes (no-PCA)	0.9649	0.9211	0.9386	0.9298	0.9292	0.9367
Naive Bayes (with-PCA)	0.9211	0.9035	0.9386	0.9386	0.8850	0.9174
Logistic Regression (no-PCA)	0.9737	0.9825	0.9649	0.9912	0.9735	0.9771
Logistic Regression (with-PCA)	0.9825	0.9912	0.9561	0.9912	0.9756	0.9793

Table 2: Fold-wise 5-fold accuracies and averages .

5 Observation Questions

1) Which models improved most with PCA? Which did not? Why?

Answer: Logistic Regression and KNN improved most (test F1 up: LogReg 0.9639→0.9756; KNN 0.9114→0.9383). Tree ensembles (RandomForest, XGBoost, GradientBoosting, AdaBoost) generally did not improve and often decreased. Reason: PCA reduces noise/correlation helping linear/distance models; tree ensembles already handle redundancy via splits and importance, so PCA's linear rotation can remove the feature structure they exploit.

2) Did PCA reduce variance across folds (more stable results)?

Answer: Partially — PCA reduced fold-to-fold spread for some linear/distance models (LogReg, KNN, SVM slight) per CV averages. For ensembles the effect was mixed; stacking's paired CV showed no significant change (paired t-test p 0.1772). See CV table and embedded prints.

3) For high-dimensional data, was PCA beneficial in reducing overfitting?

Answer: Yes for models sensitive to dimensionality and collinearity (Logistic Regression, KNN). For tree ensembles, PCA did not reduce overfitting and sometimes reduced test performance. Conclusion: PCA helps reduce overfitting for some learners but is not universally beneficial.

4) How did linear models (Logistic Regression, SVM) behave compared to ensemble models with PCA?

Answer: Logistic Regression improved with PCA; SVM remained strong with small differences. Ensemble models typically did not benefit — in several cases they lost test F1 after PCA. Thus linear models responded better to PCA than ensembles in this run.

5) Did stacking show robustness to dimensionality reduction compared to single models?

Answer: Yes — stacking remained robust. Stacking variants maintained high CV averages and top test F1s in both No-PCA and With-PCA variants. Paired CV for stacking did not show a significant difference (p 0.1772), indicating robustness to PCA.

6 Conclusion

Based on the experiment outputs embedded above (pages 11 onward), PCA (variance-target at 95%) improved performance and stability for linear/distance-based models (Logistic Regression, KNN) by reducing collinearity and noisy dimensions. Tree-based ensembles (Random Forest, XGBoost, Gradient Boosting, AdaBoost) generally did not benefit and sometimes lost performance after PCA because they internally handle redundancy and exploit raw-feature splits; PCA's rotation can remove signals these models use. Stacking ensembles showed robustness to PCA and remained among the top-performing approaches. Practical recommendation: apply PCA when using linear or distance-based classifiers or when strong multicollinearity/noise exists; for tree ensembles, prefer raw features but always compare both variants using paired CV.