

```
1. Importación y Análisis Inicial (Versión Mejorada)
import pandas as pd
```

[illegible]

```
X1954: p = 0.00000
X1389: p = 0.00000

Método BH01:
Genes significativos: 72
Top 10 genes más significativos:
X1003: p = 0.00000
X187: p = 0.00000
X2050: p = 0.00000
X1951: p = 0.00000
X1445: p = 0.00000
X244: p = 0.00000
X2046: p = 0.00000
X509: p = 0.00000
X1954: p = 0.00000
X1389: p = 0.00000

Método FDR_BH:
Genes significativos: 286
Top 10 genes más significativos:
X1003: p = 0.00000
X187: p = 0.00000
X2050: p = 0.00000
X1951: p = 0.00000
X1445: p = 0.00000
X244: p = 0.00000
X2046: p = 0.00000
X509: p = 0.00000
X1954: p = 0.00000
X1389: p = 0.00000
```

Análisis:

Borferoni con un criterio más conservador identificó 72 genes con diferencias significativas ($p < 0.05$), incluyendo X1003, X187 y X2050 como los más relevantes $p=0$. Hizo pareado a Borferoni pero ligeramente menos estricto, también detectó 72 genes significativos, coincidiendo en los primeros 10 genes con Borferoni. FDR_BH más flexible encontró 286 genes diferenciales, manteniendo los mismos 10 genes top con $p=0$, pero incluyendo muchos más candidatos posibles. Todos los métodos encuentran que X1003, X187 y X2050 son los genes con diferencias más extremas entre clases.

3. Realizaremos un experimento parecido al anterior, pero ahora comparando las medias de las 4 clases de la base de datos. Para lograrlo, en vez de trabajar con el estadístico t , probemos el análisis de variancia ANOVA. Dicha prueba la realizaremos con la función `oneway` de `scipy.stats`, pero para usar la función, necesitaremos primero estratificarlos por clase.

```
In [36]: from scipy.stats import f_oneway

# Preparar datos para ANOVA
anova_results = []
for genes in X.columns:
    # Separar datos por clase
    groups = [X[ly == clc][genes] for clc in sorted(y.unique())]
```

```
anova_results.append({'gene': gene, 'p': f_stat, 'p_val': p_val})

# Convertir a DataFrame y ordenar
anova_df = pd.DataFrame(anova_results).sort_values('p')

# Corrección Benjamini-Hochberg
rejected = corrected_p <= alpha * multipletests(anova_df['p'], alpha=0.05, method='fdr_bh')
anova_df['significant'] = rejected
anova_df['corrected_p'] = corrected_p

# Resultados
print(f"número de genes significativos: {sum(rejected)}")
print(f"número de genes más significativos: {sum(anova_df['corrected_p'] < 0.01)}")
print(anova_df[['gene', 'p', 'corrected_p']].head(10))

ANOVA - Genes significativos: 1162

Top 10 genes más significativos
      gene      p      corrected_p
1954 X1955  84.364086  2.04755e-21
3188 X1389  83.817193  2.04755e-21
1002 X1003  77.795622  1.24554e-20
2049 X2050  69.210739  1.731184e-19
240 X246  68.411042  3.06212e-19
741 X742  65.512797  6.44652e-19
0 X1  59.118264  1.26950e-17
2161 X161  56.397623  2.98397e-17
3393 X1594  55.419914  9.35774e-17
1644 X1645  54.768403  6.75174e-17
```

Aunque se mantenga cierta coincidencia en genes tipo X1955, X1003, X2050, X1645 que aparecen en ambos análisis, el ANOVA revela genes adicionales relevantes, como X742 y X1, que las comparaciones por pares no capturan, sugiriendo que estos podrían estar expresados diferente en otros contrastes entre cistitas o mostrar patrones complejos que solo se podrían ver al considerar todas las clases simultáneamente.

4. Separaremos los datos en entrenamiento y prueba, construiremos y entrenaremos un modelo de SVM con un kernel lineal, con un kernel polinomial de orden 3, y con un kernel radial. Para evitar que el tiempo de procesamiento sea exagerado, seleccionaremos solamente algunos variables, partiendo de los resultados que obtenimos en los puntos anteriores. Pero lo ideal sería que la selección de características se basara solamente en experimentos realizados con los datos de entrenamiento. Pero, en este caso, obviaremos este detalle.

```
In [181]: from sklearn.model_selection import train_test_split
from sklearn.svm import SVC
from sklearn.preprocessing import StandardScaler

# Seleccionar los 50 genes mas significativos del ANOVA
selected_genes = anova_df.head(50) ["gene"].values
X_selected = X[selected_genes]

# Estandarizar datos
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X_selected)
```

```
# Dividir datos (70% entrenamiento, 30% prueba)
X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.3, random_state=42)

# Modelos SVM
models = [
    'linear': SVC(kernel='linear', C=1, random_state=42),
    'Polinomial (grado 3)': SVC(kernel='poly', degree=3, C=1, random_state=42),
    'RBF': SVC(kernel='rbf', C=1, gamma='scale', random_state=42)
]

# Entrenamiento
```

```
In [40]: from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
import seaborn as sns
import matplotlib.pyplot as plt

for name, model in models.items():
    print(f"=== Evaluación modelo {name} ===")

    # Predicciones
```

	2	1.00	1.00	1.00	9
	3	1.00	1.00	1.00	4
	4	1.00	1.00	1.00	8
accuracy				1.00	25
macro avg	1.00	1.00	1.00	25	
weighted avg	1.00	1.00	1.00	25	

Matriz de Confusión - Lineal

	Verdadero = 1	Verdadero = 2	Verdadero = 3
Predicho = 1	4	0	0
Predicho = 2	0	9	0
Predicho = 3	0	0	4

```

=== Evaluación modelo Polynomial (grado 3) ===
Exactitud: 0.9200

```

Matriz de Confusión - Polinomial (grado 3)

	Real = 0	Real = 1
Pred = 0	4	0
Pred = 1	0	9

Vero \ Poli	1	2	3	4
3	0	2	2	0
4	0	0	0	8

Matriz de Confusión - RBF

	Actual 0	Actual 1	Actual 2	Actual 3
Predicted 0	4	0	0	0
Predicted 1	0	0	0	0

	2	3	4
2	9	0	0
3	0	4	0
4	0	0	8

Análisis:

Tanto el modelo lineal como el RBF alcanzaron una exactitud perfecta (100%), clasificando correctamente todos los casos (F1-score=1.00 en todas las clases), lo que sugiere una separabilidad óptima de los datos con estos kernels. El modelo polinomial (grado 3) mostró un rendimiento ligeramente inferior (98%, exactitud=0.98, precisión=0.98, F1=0.97), indicando posibles limitaciones para ciertos patrones complejos o ruido. La asociación del RBF y el lineal en todos los métricas (precisión, recall y F1) sugiere que los datos son altamente separables en estos espacios de características.

posición como los mejores candidatos, aunque el kernel lineal no parece por su simplicidad computacional ofrecer el mismo rendimiento perfecto. El polinomial, aunque útil, requiere ajustes (e. g. modular el grado o parámetros de regularización) para mejorar su capacidad predictiva en la clase 3.

Mientras que los resultados sugieren que el lineal y RBF son óptimos, la perfección es estadísticamente sospechosa. El polinomial (32%) podría ser más realista, y su error en la clase 3 merece investigación. Yo personalmente de debo la exactitud del 100% a la selección d genomas y la fuga de datos que tenemos en el código que permite correr más rápido el sistema.

Referencias:

OpenAI. (2025). ChatGPT (versión GPT-4) [Modelo de lenguaje grande]. <https://chat.openai.com/>