

# Ensamblés

Técnicas de Aprendizaje de Máquina

**Oscar Bustos**

August 24, 2025



Pontificia Universidad  
**JAVERIANA**  
Bogotá

VIGILADA MINEDUCACIÓN

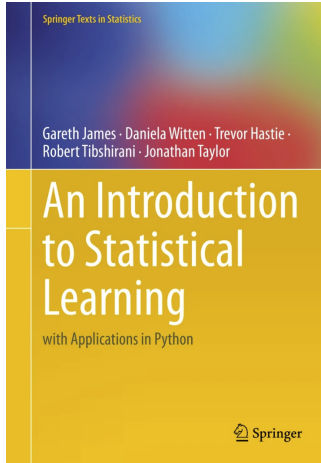
# Table of Contents

## 1 Intuición Matemática

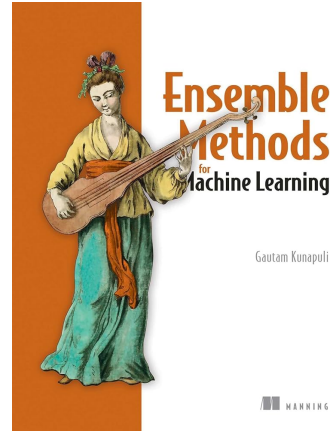
- ▶ Intuición Matemática
- ▶ Parallel Homogeneous Ensembles
- ▶ Parallel Heterogeneous Ensembles
- ▶ Sequential Gradient Boosting Ensembles
- ▶ Comparación
- ▶ Librerías

# Capítulos sobre Ensamblés

## 1 Intuición Matemática



Capítulo 9

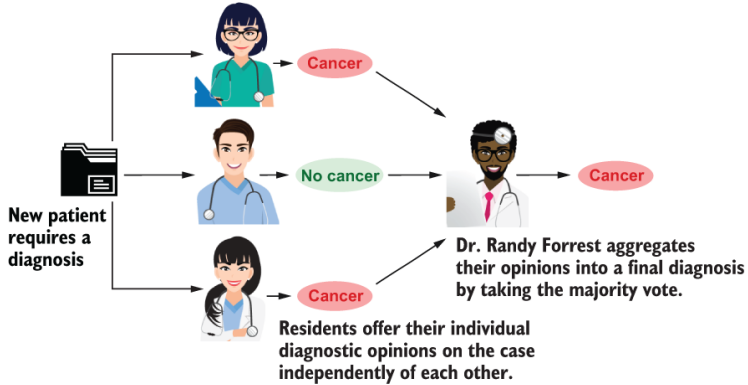


Capítulos 1,2,3,5

# Wisdom of Crowds

## 1 Intuición Matemática

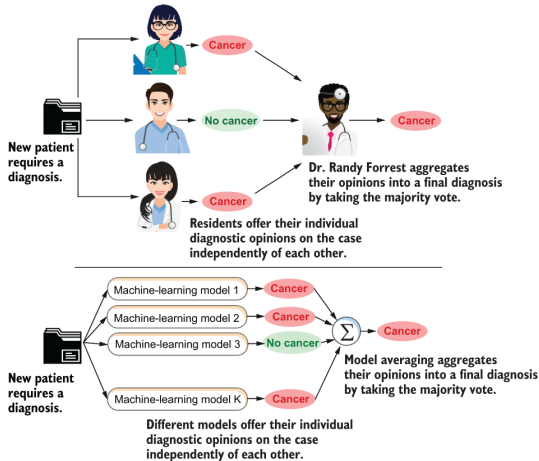
¿Cómo diagnostico si un paciente tiene cáncer o no?



# Wisdom of Crowds

## 1 Intuición Matemática

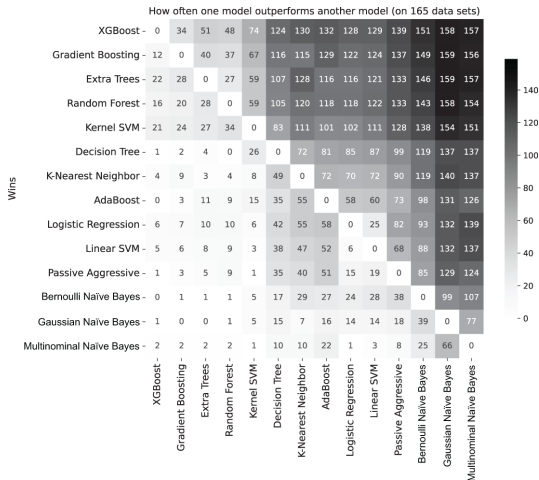
¿Cómo diagnostico si un paciente tiene cáncer o no usando IA?



# Evidencia empírica en los ensambles

## 1 Intuición Matemática

¿Qué algoritmo de aprendizaje de máquina debo utilizar? Olson et al. 2018



# ¿Qué son los Ensambles?

## 1 Intuición Matemática

- Combinación de múltiples modelos más débiles para crear un modelo más fuerte.
- Idea principal: "La sabiduría de la multitud".
- Reducen el sesgo y la varianza.
- Mejoran la precisión y la robustez del modelo.

# Resumen de Categorías de Ensamblés

## 1 Intuición Matemática

Existen tres principales categorías de ensambles en machine learning:

- **Parallel Homogeneous Ensembles:** Múltiples modelos del mismo tipo entrenados con diferentes muestras de datos. Mejora la estabilidad y reduce la varianza.
- **Parallel Heterogeneous Ensembles:** Combinación de diferentes tipos de modelos para aprovechar sus fortalezas individuales.
- **Sequential Gradient Boosting Ensembles:** Modelos entrenados secuencialmente, donde cada uno corrige los errores del anterior, logrando una mejor optimización de la función de pérdida.



# Table of Contents

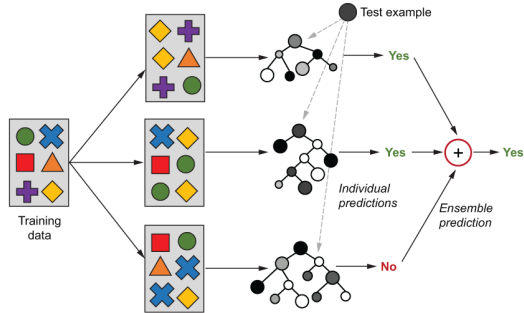
## 2 Parallel Homogeneous Ensembles

- ▶ Intuición Matemática
- ▶ **Parallel Homogeneous Ensembles**
- ▶ Parallel Heterogeneous Ensembles
- ▶ Sequential Gradient Boosting Ensembles
- ▶ Comparación
- ▶ Librerías

# Arquitectura de Ensamble Homogéneo

## 2 Parallel Homogeneous Ensembles

Cuando todos los clasificadores son del mismo tipo



**Bootstrap sampling**  
generates diverse subsets  
for training base learners.

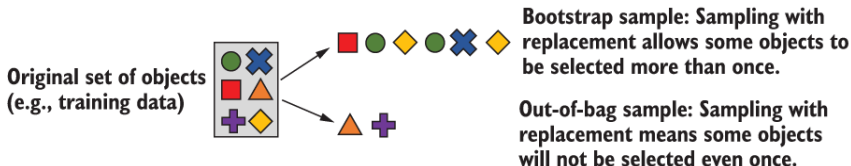
Diverse **base learners**  
are trained on sampled  
subsets of the data.

Final prediction of  
the ensemble is reached  
by **model aggregation**.

# Bootstrapping

## 2 Parallel Homogeneous Ensembles

El bootstrapping es un método estadístico que permite generar **múltiples subconjuntos de datos** a partir del conjunto de entrenamiento original. Se basa en el muestreo con reemplazo, lo que significa que algunos ejemplos pueden aparecer más de una vez en un subconjunto, mientras que otros pueden no aparecer en absoluto.



# Bagging - Entrenamiento

## 2 Parallel Homogeneous Ensembles

```
1 import numpy as np
2 from sklearn.tree import DecisionTreeClassifier
3
4 rng = np.random.RandomState(seed=4190)
5 def bagging_fit(X, y, n_estimators, max_depth=5, max_samples=200):
6     n_examples = len(y)
7     estimators = [DecisionTreeClassifier(max_depth=max_depth)
8                   for _ in range(n_estimators)]
9
10    for tree in estimators:
11        bag = np.random.choice(n_examples, max_samples,
12                               replace=True)
13        tree.fit(X[bag, :], y[bag])
14
15    return estimators
```

# Bagging - Inferencia

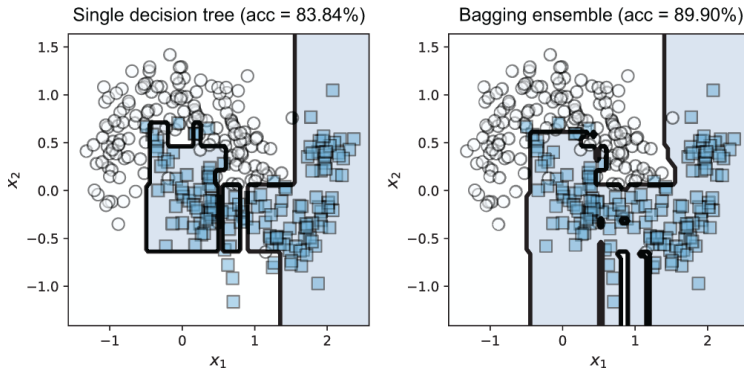
## 2 Parallel Homogeneous Ensembles

```
1 from scipy.stats import mode
2
3 def bagging_predict(X, estimators):
4     all_predictions = np.array([tree.predict(X)
5                                 for tree in estimators])
6     ypred, _ = mode(all_predictions, axis=0,
7                     keepdims=False)
8     return np.squeeze(ypred)
```

# Comparación Bagging en Árboles de Decisión

## 2 Parallel Homogeneous Ensembles

Un único árbol de decisión (izquierda) **sobreajusta** el conjunto de entrenamiento y puede ser sensible a los valores atípicos. Un conjunto de bagging (derecha) **suaviza** los efectos de sobreajuste y las clasificaciones erróneas de varios de esos estimadores base.



# Bagging en Sklearn

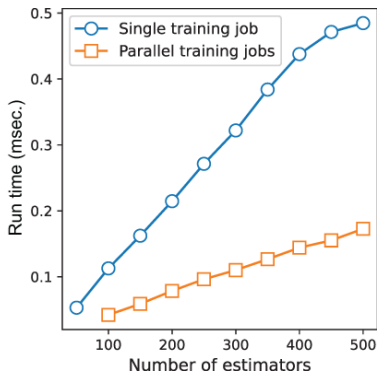
## 2 Parallel Homogeneous Ensembles

```
1 from sklearn.tree import DecisionTreeClassifier
2 from sklearn.ensemble import BaggingClassifier
3
4 base_estimator = DecisionTreeClassifier(max_depth=10)
5 bag_ens = BaggingClassifier(base_estimator=base_estimator,
6                             n_estimators=500,
7                             max_samples=100,
8                             oob_score=True,
9                             random_state=rng)
10 bag_ens.fit(Xtrn, ytrn)
11 ypred = bag_ens.predict(Xtst)
```

# Optimización en Bagging

## 2 Parallel Homogeneous Ensembles

```
1  
2 BaggingClassifier(base_estimator=DecisionTreeClassifier(),  
3                   n_estimators=100, max_samples=100,  
4                   oob_score=True, n_jobs=-1)
```

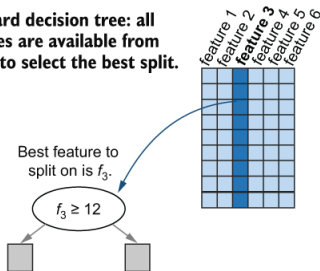




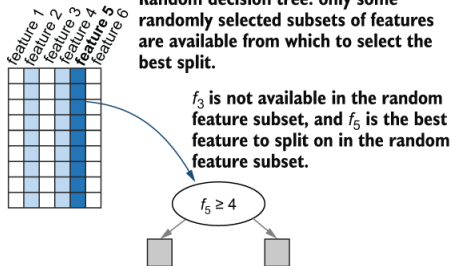
# Características Aleatorias

## 2 Parallel Homogeneous Ensembles

**Standard decision tree: all features are available from which to select the best split.**



**Random decision tree: only some randomly selected subsets of features are available from which to select the best split.**



# Random Forest: Ventajas

## 2 Parallel Homogeneous Ensembles

- Alta precisión y robustez.
- Maneja bien datos de alta dimensión.
- Reduce el sobreajuste.
- Proporciona importancia de las características.

# Random Forest: Hiperparámetros Clave

## 2 Parallel Homogeneous Ensembles

- `n_estimators`: Número de árboles en el bosque.
- `max_features`: Número de características a considerar en cada división.
- `max_depth`: Profundidad máxima de los árboles.
- `min_samples_split`: Número mínimo de muestras requeridas para dividir un nodo interno.

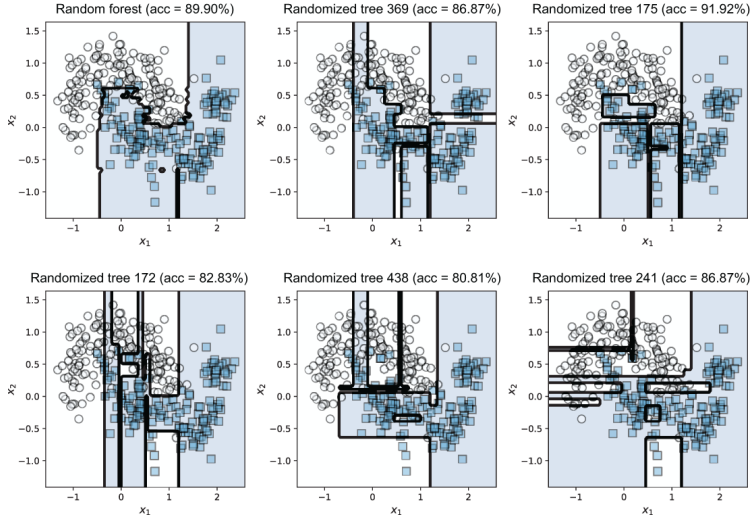
# Random Forest en Sklearn

## 2 Parallel Homogeneous Ensembles

```
1 from sklearn.ensemble import RandomForestClassifier
2
3 rf_ens = RandomForestClassifier(n_estimators=500,
4                               max_depth=10,
5                               oob_score=True,
6                               n_jobs=-1,
7                               random_state=rng)
8 rf_ens.fit(Xtrn, ytrn)
9 ypred = rf_ens.predict(Xtst)
```

# Random Forest en Sklearn

## 2 Parallel Homogeneous Ensembles



# Random Forest en Sklearn

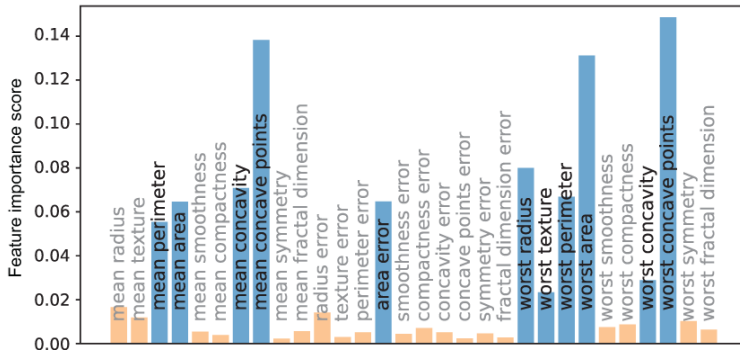
## 2 Parallel Homogeneous Ensembles

```
1 X_trn, X_tst, y_trn, y_tst = train_test_split(X, y, test_size=0.15)
2 n_features = X_trn.shape[1]
3
4 rf = RandomForestClassifier(max_leaf_nodes=24,
5                             n_estimators=50, n_jobs=-1)
6 rf.fit(X_trn, y_trn)
7 err = 1 - accuracy_score(y_tst, rf.predict(X_tst))
8
9 importance_threshold = 0.02
10 for i, (feature, importance) in enumerate(zip(dataset['feature_names',
11
12         rf.feature_importances_))):
13     if importance > importance_threshold:
14         print('{0} {1} (score={2:4.3f})'.format(i, feature, importance))
```

# Random Forest en Sklearn

## 2 Parallel Homogeneous Ensembles

Importancia de Características



# Table of Contents

## 3 Parallel Heterogeneous Ensembles

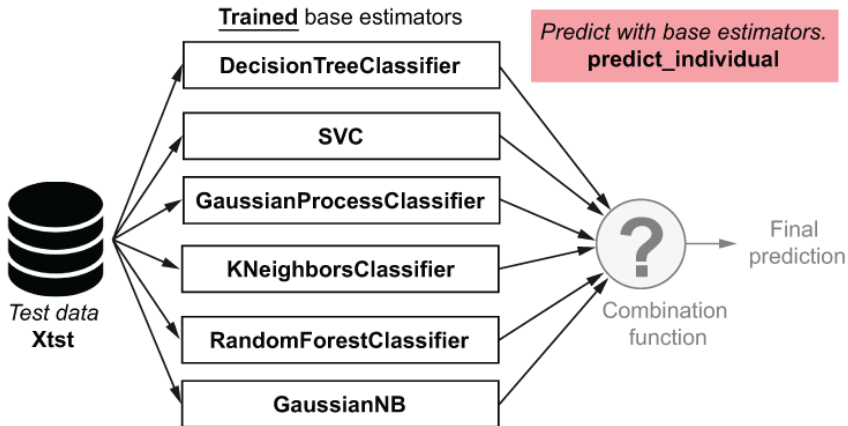
- ▶ Intuición Matemática
- ▶ Parallel Homogeneous Ensembles
- ▶ **Parallel Heterogeneous Ensembles**
- ▶ Sequential Gradient Boosting Ensembles
- ▶ Comparación
- ▶ Librerías



# Arquitectura de Ensamble Heterogéneo

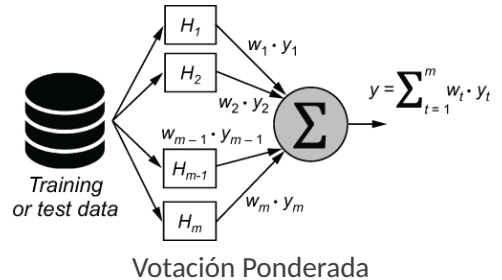
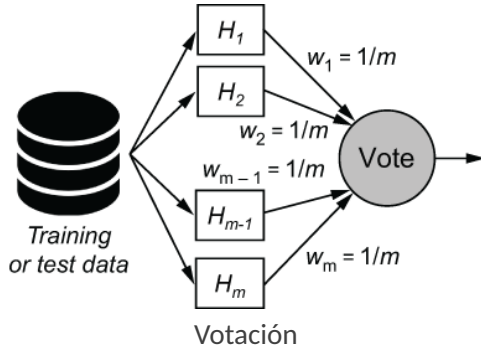
## 3 Parallel Heterogeneous Ensembles

Cuando todos los clasificadores son de tipos distintos



# Voting

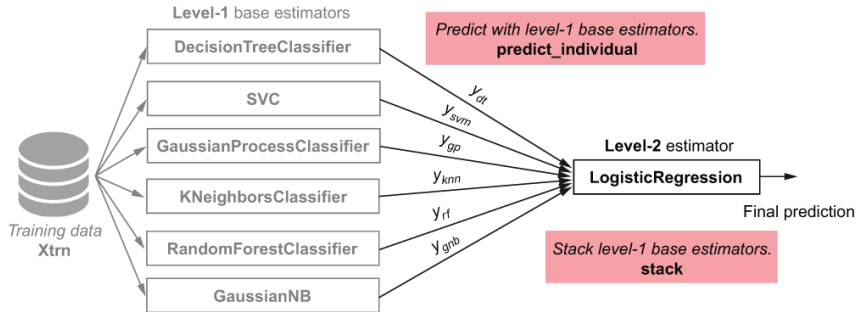
## 3 Parallel Heterogeneous Ensembles



# Stacking

## 3 Parallel Heterogeneous Ensembles

Entrenando un modelo que usa como entrada las salidas de otro modelo



# Table of Contents

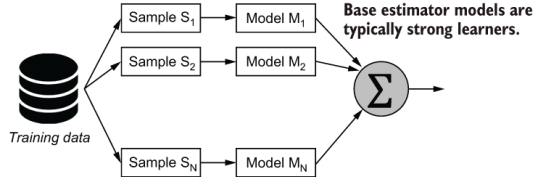
## 4 Sequential Gradient Boosting Ensembles

- ▶ Intuición Matemática
- ▶ Parallel Homogeneous Ensembles
- ▶ Parallel Heterogeneous Ensembles
- ▶ Sequential Gradient Boosting Ensembles
- ▶ Comparación
- ▶ Librerías

# Arquitectura de Gradient Boosting

## 4 Sequential Gradient Boosting Ensembles

### Parallel ensembles

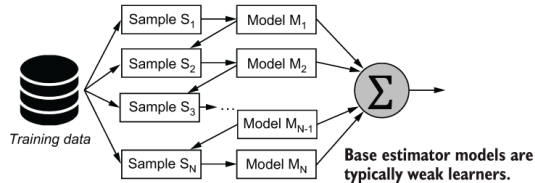


Generate subsets  
from original data.

Train multiple  
base classifiers.

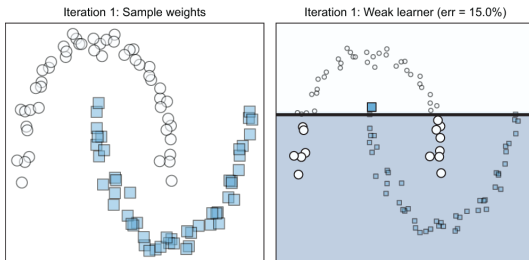
Combine/aggregate  
base classifiers.

### Sequential ensembles



# Gradient Boosted Trees (GBT): Intuición

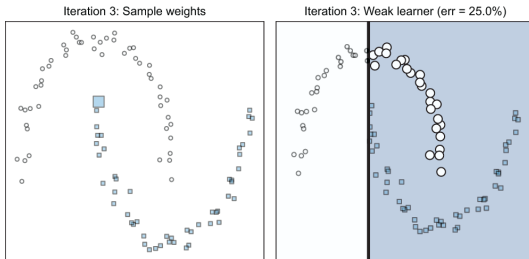
## 4 Sequential Gradient Boosting Ensembles



- Construye modelos de forma secuencial, donde cada nuevo modelo se enfoca en corregir los errores del modelo anterior.
- Utiliza el gradiente descendente para minimizar la función de pérdida.

# Gradient Boosted Trees (GBT): Intuición

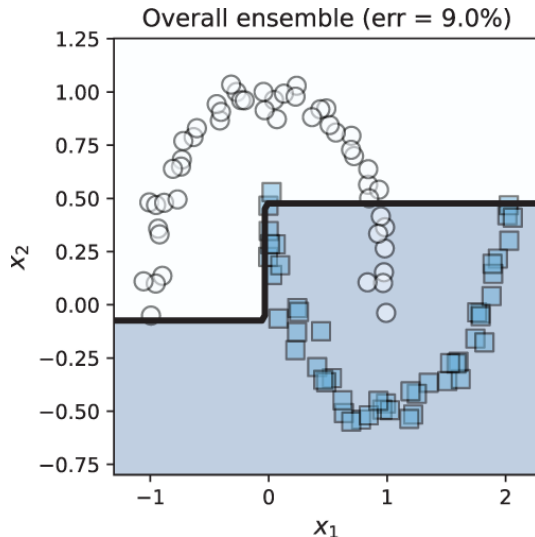
## 4 Sequential Gradient Boosting Ensembles



- Construye modelos de forma secuencial, donde cada nuevo modelo se enfoca en corregir los errores del modelo anterior.
- Utiliza el gradiente descendente para minimizar la función de pérdida.

# Gradient Boosted Trees (GBT): Intuición

## 4 Sequential Gradient Boosting Ensembles





# Gradient Boosted Trees (GBT): Hiperparámetros Clave

## 4 Sequential Gradient Boosting Ensembles

- `n_estimators`: Número de árboles en el modelo.
- `learning_rate`: Tasa de aprendizaje, controla el impacto de cada árbol.
- `max_depth`: Profundidad máxima de los árboles.
- `subsample`: Fracción de muestras utilizadas para entrenar cada árbol.
- `colsample_bytree`: Fracción de características utilizadas para entrenar cada árbol.

# Table of Contents

## 5 Comparación

- ▶ Intuición Matemática
- ▶ Parallel Homogeneous Ensembles
- ▶ Parallel Heterogeneous Ensembles
- ▶ Sequential Gradient Boosting Ensembles
- ▶ Comparación
- ▶ Librerías

# Random Forest vs Gradient Boosted Trees

## 5 Comparación

Característica	Random Forest	Gradient Boosted Trees
Enfoque	Bagging	Boosting
Paralelización	Sí	No (secuencial)
Precisión	Alta	Muy alta
Sobreaajuste	Menor riesgo	Mayor riesgo (requiere ajuste cuidadoso)
Velocidad	Rápido	Más lento (depende de la implementación)
Importancia de características	Sí	Sí



# Table of Contents

## 6 Librerías

- ▶ Intuición Matemática
- ▶ Parallel Homogeneous Ensembles
- ▶ Parallel Heterogeneous Ensembles
- ▶ Sequential Gradient Boosting Ensembles
- ▶ Comparación
- ▶ **Librerías**



# Ensamblen en Python

## 6 Librerías

- **Random Forest:** `sklearn.ensemble.RandomForestClassifier`,  
`sklearn.ensemble.RandomForestRegressor`
- **Gradient Boosting:** `sklearn.ensemble.GradientBoostingClassifier`,  
`sklearn.ensemble.GradientBoostingRegressor`
- **XGBoost:** `xgboost`
- **LightGBM:** `lightgbm` (Microsoft)
- **CatBoost:** `catboost`

# Ejemplo de Código: Random Forest (scikit-learn)

## 6 Librerías

```
1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn.model_selection import train_test_split
3 from sklearn.metrics import accuracy_score
4
5 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
    =0.2)
6
7 rf = RandomForestClassifier(n_estimators=100, max_depth=5)
8 rf.fit(X_train, y_train)
9
10 y_pred = rf.predict(X_test)
11 accuracy = accuracy_score(y_test, y_pred)
12 print("Accuracy:", accuracy)
```

# Ejemplo de Código: Gradient Boosting (XGBoost)

6 Librerías

```
1 import xgboost as xgb
2 from sklearn.model_selection import train_test_split
3 from sklearn.metrics import accuracy_score
4
5 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
    =0.2)
6
7 xgb_model = xgb.XGBClassifier(n_estimators=100, max_depth=5,
    learning_rate=0.1)
8 xgb_model.fit(X_train, y_train)
9
10 y_pred = xgb_model.predict(X_test)
11 accuracy = accuracy_score(y_test, y_pred)
12 print("Accuracy:", accuracy)
```



# References

7 Librerías



James, Gareth et al. (2023). *An introduction to statistical learning: With applications in python*. Springer Nature.



Kunapuli, Gautam (2023). *Ensemble methods for machine learning*. Simon and Schuster.



Olson, Randal S et al. (2018). "Data-driven advice for applying machine learning to bioinformatics problems". In: *Pacific symposium on biocomputing 2018: Proceedings of the pacific symposium*. World Scientific, pp. 192–203.