

# Electiva Bases de datos

## Ensamblés

Jimmy Mateo Guerrero<sup>1</sup>

<sup>1</sup>Departamento de Sistemas  
Facultad de Ingeniería

Ensamblés, 2019

1 Ensamblés

2 Bagging

3 Boosting

- Conjunto de modelos que se usan juntos como un meta modelo.
- Usar conocimiento de distintas fuentes al tomar decisiones.

- Comité de expertos:
  - muchos elementos
  - todos con alto conocimiento
  - todos sobre el mismo tema
  - votan

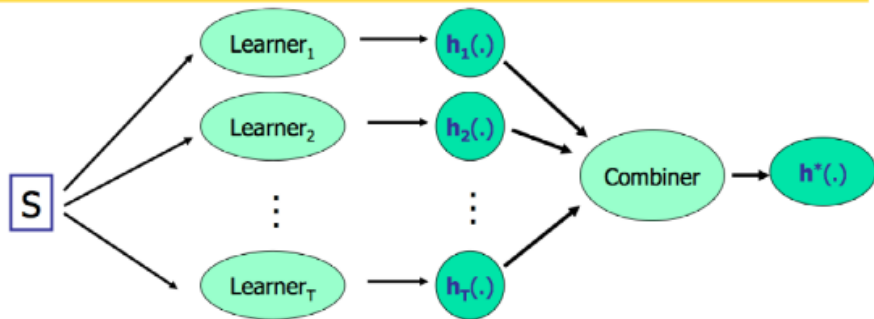
- Dos componentes base
  - Un método para seleccionar o construir los miembros
    - Distintos datasets x distintos modelos x distintas configuraciones
  - Un método para combinar las decisiones
    - Votación simple, votación ponderada, promedio, función específica, selectividad ...

- Ensamblajes Planos:
  - Muchos expertos, todos buenos.
  - Necesito que sean lo mejor posible individualmente (De lo contrario, usualmente no sirven.).
  - Pero necesito que opinen distinto en algunos casos (Si todos opinan siempre igual... me quedo con uno solo!).

# Ensembles

Original Data

Classifier



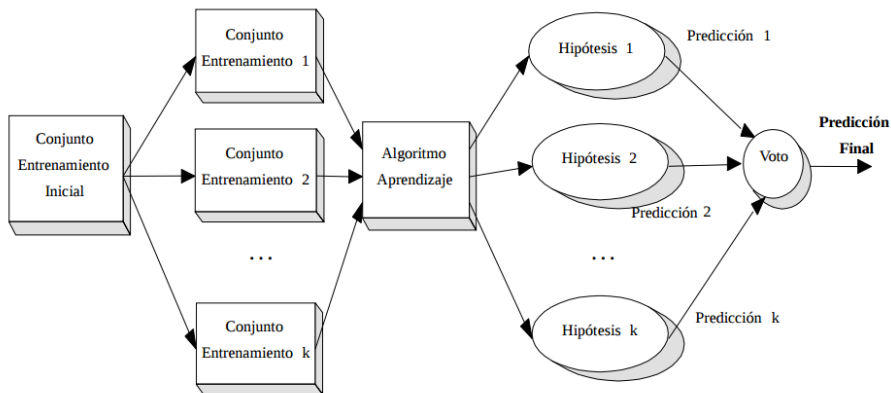
- Un aprendiz se dice inestable si el clasificador que produce sufre cambios importantes ante pequeñas variaciones en los datos de entrenamiento.
  - Inestables: árbol de decisiones, redes neuronales.
  - Estables: La regresión lineal, el vecino más cercano.
- Subsampling es mejor para los alumnos inestables.



- Existen dos tipos de algoritmos de votación:
  - aquellos que cambian adaptativamente la distribución del conjunto de entrenamiento basado en el desempeño de clasificadores anteriores (boosting).
  - y los que no (como en Bagging).

- Bootstrap aggregating (Breiman 96)
- Clasificadores de votos generados por diferentes muestras de bootstrap (réplicas)
- Se generan  $T$  muestras de bootstrap  $B_1, B_2, \dots, B_T$  y se construye un  $C_i$  clasificador a partir de cada muestra de bootstrap  $B_i$
- Un último clasificador  $C^*$  se construye a partir de  $C_1, C_2, \dots, C_T$  cuya salida es la clase predicha mas frecuente o votada por los clasificadores.

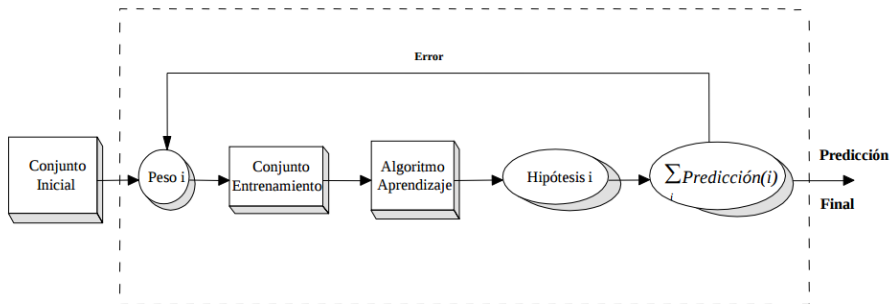
# Bagging



- Boosting (Schapire 90), AdaBoost M1 (Freund y Schapire 96)
- Genera los clasificadores secuencialmente, mientras que Bagging puede generarlos en paralelo.
- AdaBoost cambia los pesos de las instancias de entrenamiento proporcionadas como entrada para cada inductor en función de los clasificadores que se construyeron previamente.
- El objetivo es obligar al inductor a minimizar el error esperado sobre diferentes distribuciones de entrada.
- $C^*$  = votación ponderada. El peso de cada classifier depende de su rendimiento en el conjunto de entrenamiento utilizado para construir

- Boosting: buscar nuevos modelos para las instancias mal clasificadas por los anteriores.
- Fuerza al algoritmo a centrarse en los ejemplos mal clasificados por las hipótesis anteriores.
- Las instancias incorrectas son ponderadas por un factor inversamente proporcional al error en el conjunto de entrenamiento, es decir,  $1 / (2E_i)$ . Pequeños errores de entrenamiento, como 0.1 %, harán que los pesos crezcan en varios órdenes de magnitud.

# Boosting





Aurélien Géron

*Hands-On Machine Learning with Scikit-Learn and TensorFlow:  
Concepts, Tools, and Techniques to Build Intelligent Systems.*

Publisher: O'Reilly Media, Year: 2017



Jake VanderPlas.

Python Data Science Handbook: Essential Tools for Working with  
Data

*O'Reilly Media, Year: 2016*



Brett Lantz.

Machine learning with R

*Packt Publishing, Year: 2013*