

Diseño avanzado de algoritmos

Clase 8: Ejemplo de aplicación divide y conquista

Dr. Ing. Fernando López Hernández

2018

Clase 8: Ejemplo de aplicación: divide y conquista

1

Árboles de decisión

Aprendizaje automático:

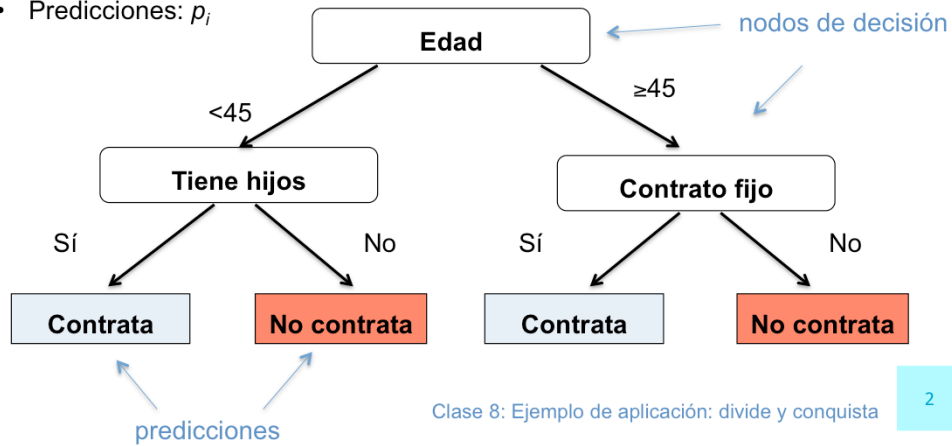
- Muestras: $\mathbf{s}_i = (\mathbf{x}_i, c_i)$ feature + clase

Clasificación:

- Observaciones: \mathbf{o}_i
- Predicciones: p_i

Ventajas:

- Rápidos
- Fáciles de interpretar
- Flexibilidad en las features



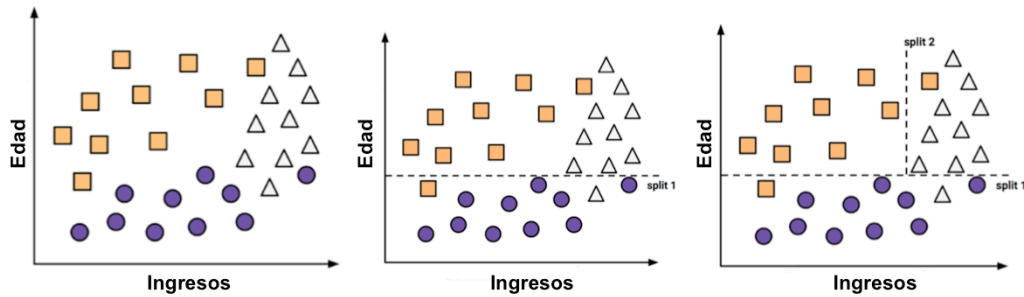
Construcción del árbol

Estrategias:

- Divide y conquista
- Búsqueda avara

Condición de parada:

- Pureza
- Falta de features
- Profundidad del árbol



Particiones paralelas a los ejes → SVM (Support Vector Machine)

Clase 8: Ejemplo de aplicación: divide y conquista

3

Elegir la mejor partición I

Pureza

$Purity(\mathbf{s}) = 1 \rightarrow$ Todas las muestras \mathbf{s} son de la misma clase c_j

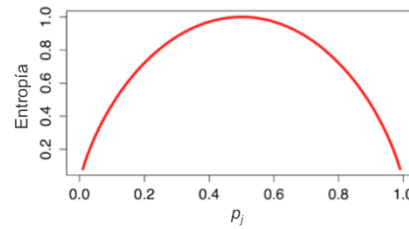
$Purity(\mathbf{s}) = 0 \rightarrow$ Las muestras \mathbf{s} equidistribuidas entre las m clases $c_1, \dots, c_j, \dots, c_m$

Métricas de impureza

1. Entropía

$$E(\mathbf{s}) = \sum_{j=1}^m -p_j \log_2 p_j$$

$$p_j = \frac{c_j \in |\mathbf{s}|}{|\mathbf{s}|}$$



2. Impureza Gini

$$I_G(\mathbf{s}) = \sum_{j=1}^m p_j (1 - p_j) = \sum_{j=1}^m (p_j - p_j^2) = \sum_{j=1}^m p_j - \sum_{j=1}^m p_j^2 = 1 - \sum_{j=1}^m p_j^2$$

Clase 8: Ejemplo de aplicación: divide y conquista

4

Elegir la mejor partición II

Ganancia de información

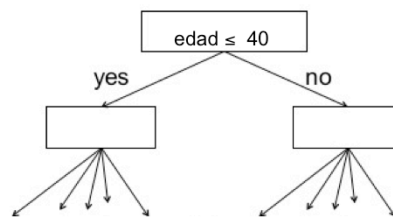
$$InfoGain(\mathbf{s}) = Impurity(\mathbf{s}) - \sum_{j=1}^m p_j \cdot Impurity(\mathbf{s}_j)$$

Partición:

- Elegir la feature que maximice la ganancia de información
- Criterio de parada

Features:

- Categóricas
- Numerales



Clase 8: Ejemplo de aplicación: divide y conquista

Construcción del árbol de decisión

```
decision_tree(s, f)
  IF empty(f) OR impurity(s)=0.0 THEN
    return support(strue, sfalse)
  strue, sfalse, fbest = best_partition(s, f)
  gain = impurity(s) - |strue|/|s| · impurity(strue) + |sfalse|/|s| · impurity(sfalse)
  IF (gain>0) THEN
    f = f - fbest
    branchtrue = decision_tree(strue, f)
    branchfalse = decision_tree(sfalse, f)
    return (fbest, branchtrue, branchfalse)
  ELSE
    return support(|strue|, |sfalse|)
```

Clase 8: Ejemplo de aplicación: divide y conquista

6

Algoritmo de clasificación

```
classify(o, T)
  IF is_prediction(T) THEN
    return majority_class(T)
  IF (NOT T.f in o)
    supporttrue = classify(o, T.branchtrue)
    supportfalse = classify(o, T.branchfalse)
    return majority_class(supporttrue, supportfalse)
  IF [is_numeric(o.f) AND o.f ≤ T.f] OR [is_categorical(o.f) AND o.f = T.f] THEN
    return classify(o, T.branchtrue)
  ELSE
    return classify(o, T.branchfalse)
```

ID3

- Permite decisiones no binarias
- Las predicciones son binarias
- No gestiona features numéricas
- Mide impureza con entropía

CART

- Las decisiones son binarias
- Las predicciones pueden no ser binarias
- Permite features numéricas
- Mide impureza con impureza Gini

C5.0

- Permite matrices de coste
- Permite podar

Clase 8: Ejemplo de aplicación: divide y conquista

8