

Ciencia de datos

Para el sector público de salud

Módulo 4

Sesión 4: Modelos de clasificación



academia .opensaludlab.org

opensaludlab.org

Twitter / Instagram / LinkedIn

Elaboración de iniciativas de participantes

Presentaciones durante Junio 2021

Temática: libre, pero relacionada a la CD

Aplicar los conceptos aprendidos y las etapas de un proyecto de CD

Grupos de **2 a 5 personas (plazo 23.3.21)**

Twist: canal #CD Proyectos

Ok... vamos a lo que nos convoca

Machine Learning

El comienzo...



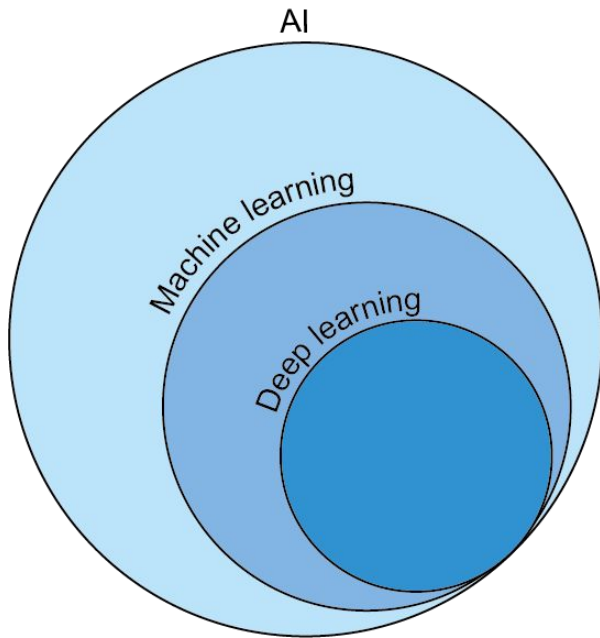
ML ¿Qué es y cuándo usarlo?

- Algoritmos de inteligencia artificial
- Perfectos para máxima predicción (modelos)
- En general son malos para explicar
- El mejor caso es cuando tengo al menos 1000 casos
- Si tengo millones de datos, mejor Deep Learning
- Tanto para clasificación como regresión

ML ó Aprendizaje automático

Es un subcampo de la inteligencia artificial donde los algoritmos “aprenden” patrones de los datos para realizar una tarea específica.

algoritmo != modelo



2. Then we pass data through these rules.



3. This gives us answers, based on our rules.

1. We create the algorithm by specifying all the rules ourselves.

1. We give data to an algorithm.



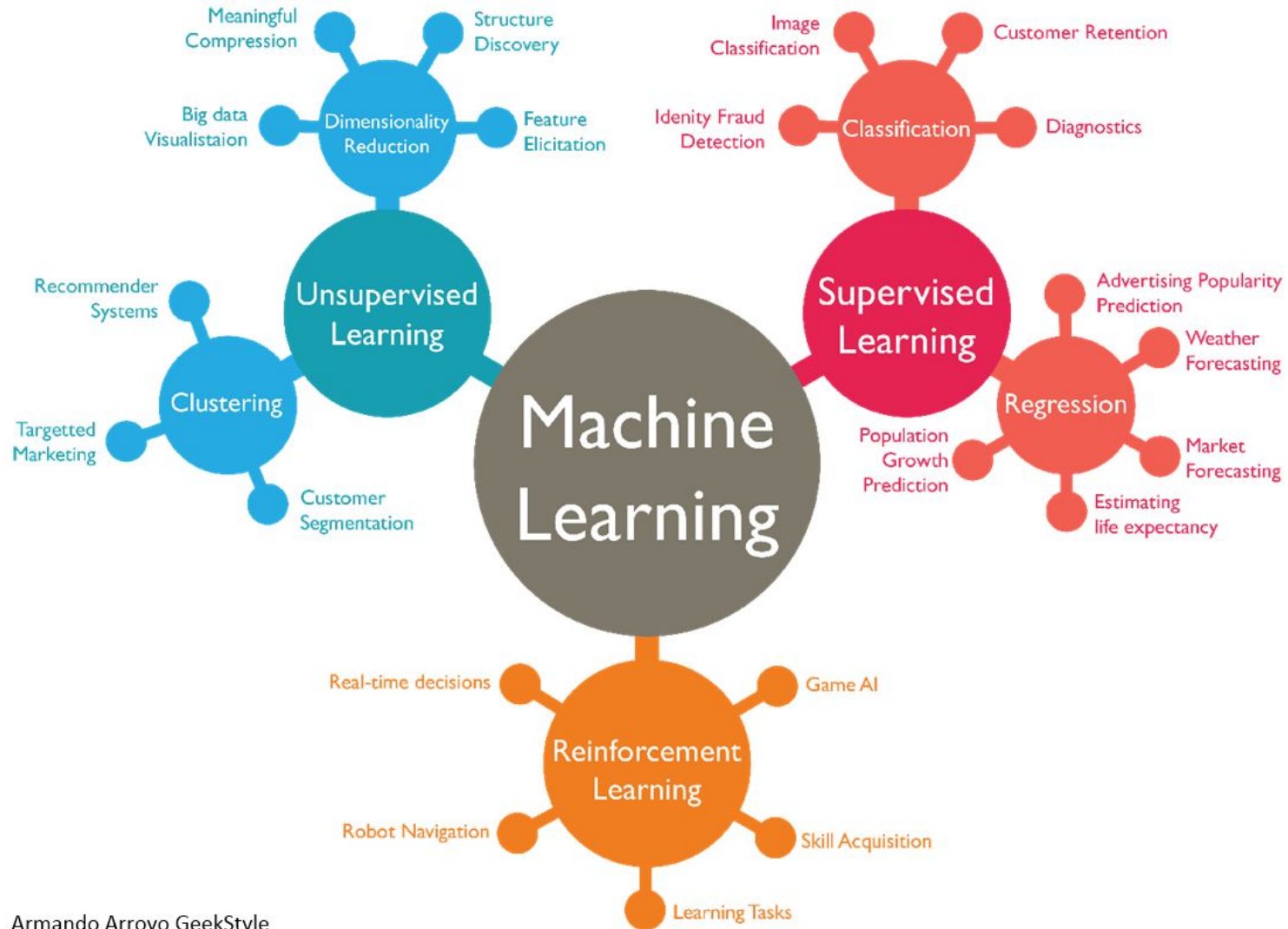
3. The algorithm learns the rules that map the answers to the data.

2. Either we also give it the answers, or it learns them for itself.



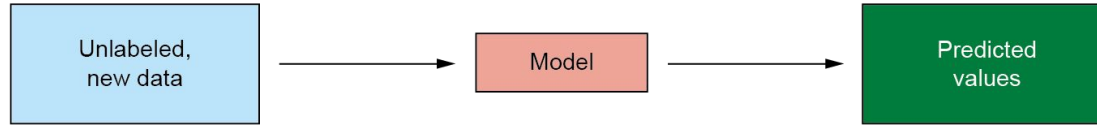
4. We pass new data through these rules...

5. ...and get answers for the new data.



1. We pass labeled data to a supervised algorithm.

2. The algorithm learns the relationships in the data and outputs a model.

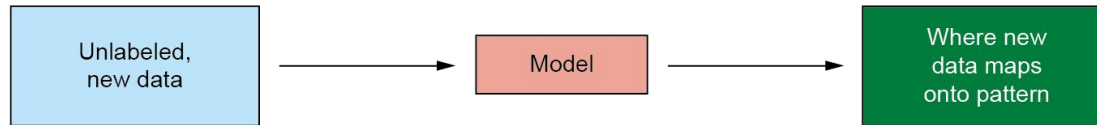
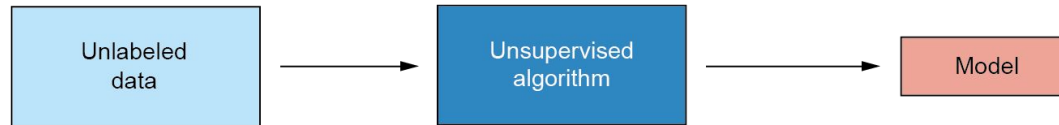


3. We pass unlabeled data into the model...

4. ...and get predicted values/labels for the new data.

1. We pass unlabeled data to an unsupervised algorithm.

2. The algorithm learns the patterns in the data and outputs a model.

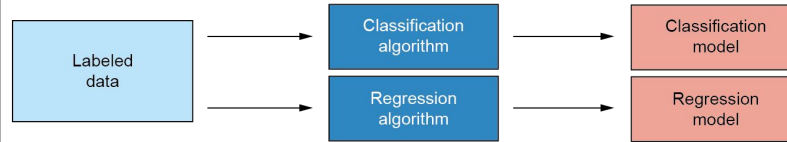


3. We pass new, unlabeled data into the model...

4. ...and get where the new data maps onto these patterns.

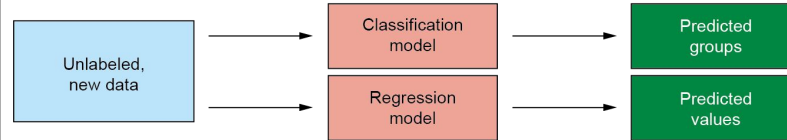
Supervised

1. Classification and regression algorithms are given labeled data.



2. They output classification and regression models, respectively.

3. We pass unlabeled, new data into the models.

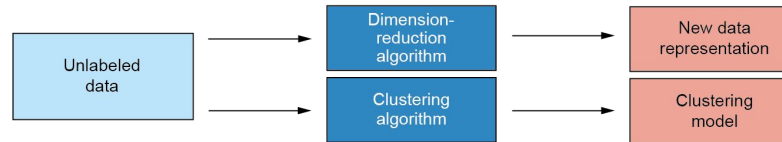


4. Classification models predict group membership.

5. Regression models predict continuous variables.

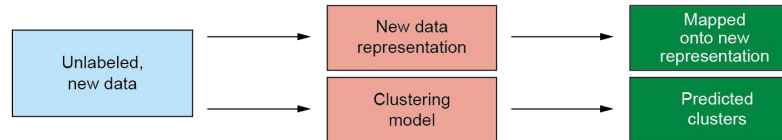
Unsupervised

1. Dimension reduction and clustering algorithms are given unlabeled data.



2. They output a lower-dimension representation of the data and clustering model, respectively.

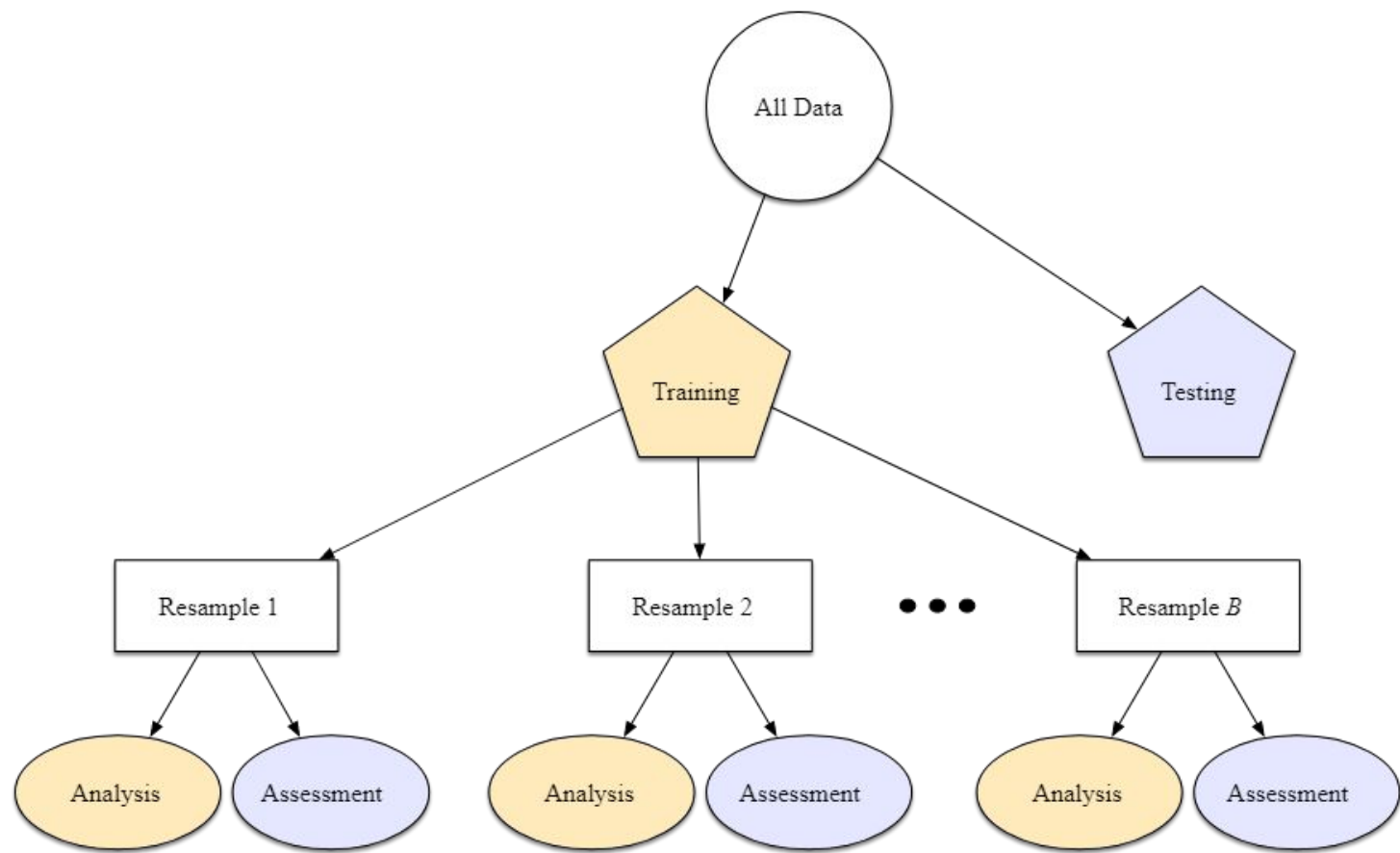
3. We pass unlabeled, new data into the models.



4. Dimension reduction maps new data onto the lower-dimensional representation.

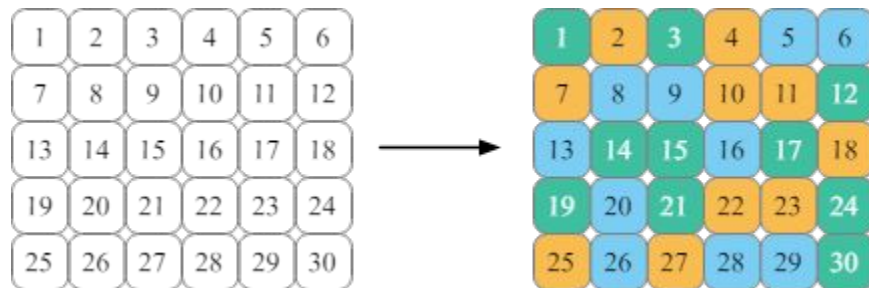
5. Clustering models predict cluster membership.

Resampleo



Cross-validation

La validación cruzada es un método de remuestreo bien establecido. Si bien hay una serie de variaciones, el método de validación cruzada más común es la validación cruzada en V . Los datos se dividen aleatoriamente en V conjuntos de aproximadamente el mismo tamaño (llamados Folds o "pliegues"). A modo de ilustración, a continuación se muestra $V = 3$ para un conjunto de datos de treinta puntos de ajuste de entrenamiento con asignaciones de pliegues (folds) aleatorias. El número dentro de los símbolos es el número de muestra:



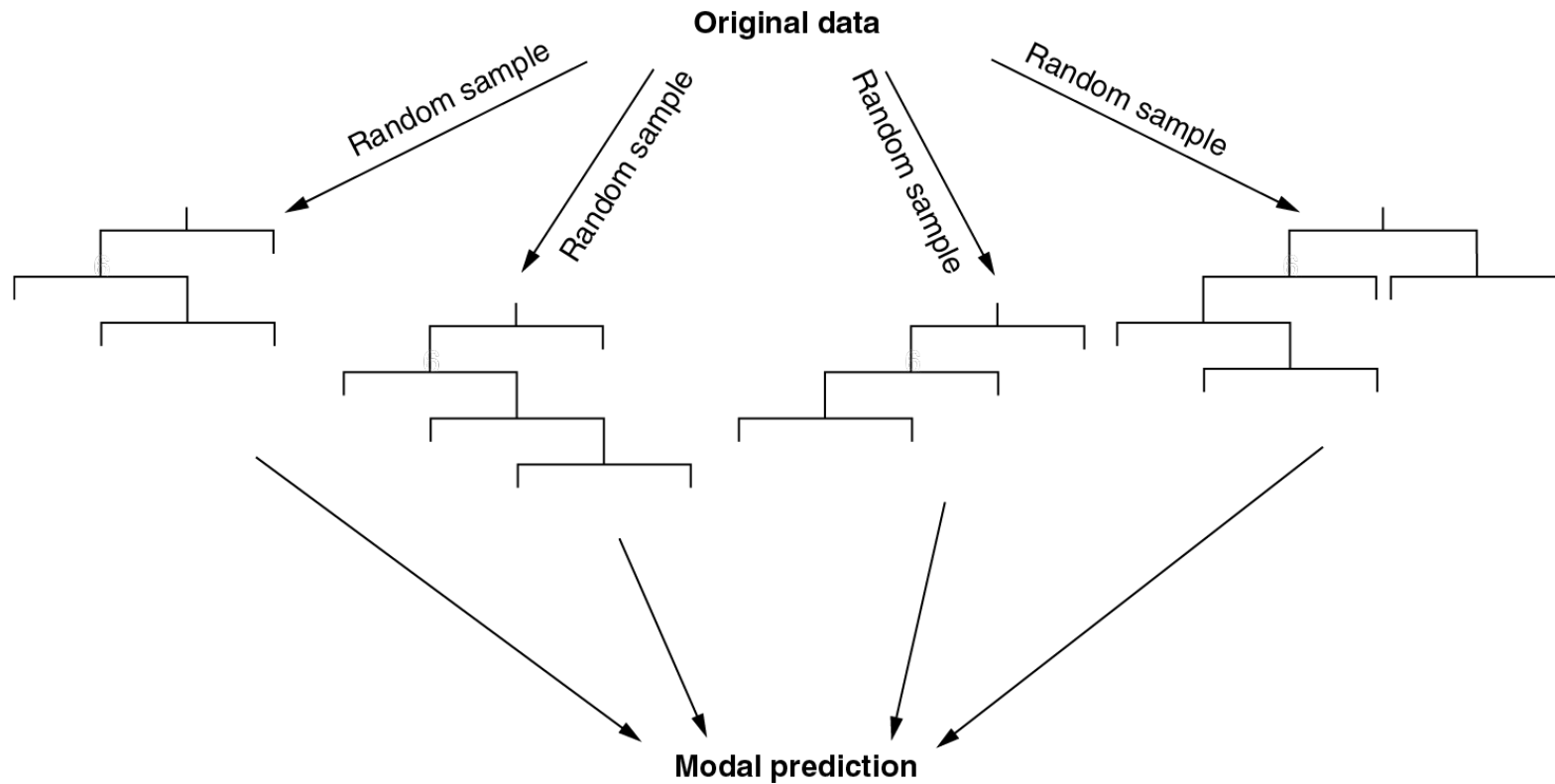
	Fold 1 Iteration	Fold 2 Iteration	Fold 3 Iteration
Model Fit Using	<div> <div>2</div><div>4</div><div>5</div><div>6</div> <div>7</div><div>8</div><div>9</div><div>10</div> <div>11</div><div>13</div><div>16</div><div>18</div> <div>20</div><div>22</div><div>23</div><div>25</div> <div>26</div><div>27</div><div>28</div><div>29</div> </div>	<div> <div>1</div><div>3</div><div>5</div><div>6</div> <div>8</div><div>9</div><div>12</div><div>13</div> <div>14</div><div>15</div><div>16</div><div>17</div> <div>19</div><div>20</div><div>21</div><div>24</div> <div>26</div><div>28</div><div>29</div><div>30</div> </div>	<div> <div>1</div><div>2</div><div>3</div><div>4</div> <div>7</div><div>10</div><div>11</div><div>12</div> <div>14</div><div>15</div><div>17</div><div>18</div> <div>19</div><div>21</div><div>22</div><div>23</div> <div>24</div><div>25</div><div>27</div><div>30</div> </div>
Estimate Performance Using	<div> <div>1</div><div>3</div> <div>12</div><div>14</div> <div>15</div><div>17</div> <div>19</div><div>21</div> <div>24</div><div>30</div> </div>	<div> <div>2</div><div>4</div> <div>7</div><div>10</div> <div>11</div><div>18</div> <div>22</div><div>23</div> <div>25</div><div>27</div> </div>	<div> <div>5</div><div>6</div> <div>8</div><div>9</div> <div>13</div><div>16</div> <div>20</div><div>26</div> <div>28</div><div>29</div> </div>

Bootstrapping

Un Bootstrap es una muestra que tiene el mismo tamaño que el conjunto de entrenamiento, pero se extrae con reemplazo. Esto significa que algunos puntos de datos del conjunto de entrenamiento se seleccionan varias veces para el conjunto de análisis.

	Bootstrap Iteration 1	Bootstrap Iteration 2	Bootstrap Iteration 3																																																																																										
Model Fit Using	<table><tr><td>1</td><td>1</td><td>4</td><td>7</td><td>8</td><td>8</td></tr><tr><td>10</td><td>13</td><td>13</td><td>13</td><td>14</td><td>15</td></tr><tr><td>16</td><td>16</td><td>16</td><td>17</td><td>19</td><td>19</td></tr><tr><td>21</td><td>22</td><td>23</td><td>23</td><td>24</td><td>23</td></tr><tr><td>25</td><td>25</td><td>25</td><td>27</td><td>28</td><td>29</td></tr></table>	1	1	4	7	8	8	10	13	13	13	14	15	16	16	16	17	19	19	21	22	23	23	24	23	25	25	25	27	28	29	<table><tr><td>2</td><td>2</td><td>3</td><td>3</td><td>3</td><td>4</td></tr><tr><td>4</td><td>4</td><td>6</td><td>6</td><td>7</td><td>10</td></tr><tr><td>11</td><td>12</td><td>12</td><td>14</td><td>14</td><td>15</td></tr><tr><td>17</td><td>17</td><td>18</td><td>21</td><td>22</td><td>22</td></tr><tr><td>23</td><td>23</td><td>28</td><td>27</td><td>28</td><td>30</td></tr></table>	2	2	3	3	3	4	4	4	6	6	7	10	11	12	12	14	14	15	17	17	18	21	22	22	23	23	28	27	28	30	<table><tr><td>2</td><td>2</td><td>3</td><td>3</td><td>4</td><td>5</td></tr><tr><td>5</td><td>5</td><td>6</td><td>7</td><td>10</td><td>11</td></tr><tr><td>12</td><td>15</td><td>16</td><td>18</td><td>18</td><td>19</td></tr><tr><td>19</td><td>20</td><td>20</td><td>20</td><td>21</td><td>21</td></tr><tr><td>21</td><td>21</td><td>22</td><td>22</td><td>29</td><td>30</td></tr></table>	2	2	3	3	4	5	5	5	6	7	10	11	12	15	16	18	18	19	19	20	20	20	21	21	21	21	22	22	29	30
1	1	4	7	8	8																																																																																								
10	13	13	13	14	15																																																																																								
16	16	16	17	19	19																																																																																								
21	22	23	23	24	23																																																																																								
25	25	25	27	28	29																																																																																								
2	2	3	3	3	4																																																																																								
4	4	6	6	7	10																																																																																								
11	12	12	14	14	15																																																																																								
17	17	18	21	22	22																																																																																								
23	23	28	27	28	30																																																																																								
2	2	3	3	4	5																																																																																								
5	5	6	7	10	11																																																																																								
12	15	16	18	18	19																																																																																								
19	20	20	20	21	21																																																																																								
21	21	22	22	29	30																																																																																								
Estimate Performance Using	<table><tr><td>2</td><td>3</td><td>5</td><td>6</td><td>9</td><td>11</td></tr><tr><td>12</td><td>18</td><td>20</td><td>24</td><td>26</td><td>28</td></tr><tr><td></td><td></td><td>30</td><td></td><td></td><td></td></tr></table>	2	3	5	6	9	11	12	18	20	24	26	28			30				<table><tr><td>1</td><td>5</td><td>8</td><td>9</td><td>13</td><td>16</td></tr><tr><td>19</td><td>20</td><td>24</td><td>26</td><td>29</td><td></td></tr></table>	1	5	8	9	13	16	19	20	24	26	29		<table><tr><td>1</td><td>8</td><td>9</td><td>13</td><td>14</td><td>17</td></tr><tr><td>23</td><td>24</td><td>25</td><td>26</td><td>27</td><td>28</td></tr></table>	1	8	9	13	14	17	23	24	25	26	27	28																																																
2	3	5	6	9	11																																																																																								
12	18	20	24	26	28																																																																																								
		30																																																																																											
1	5	8	9	13	16																																																																																								
19	20	24	26	29																																																																																									
1	8	9	13	14	17																																																																																								
23	24	25	26	27	28																																																																																								

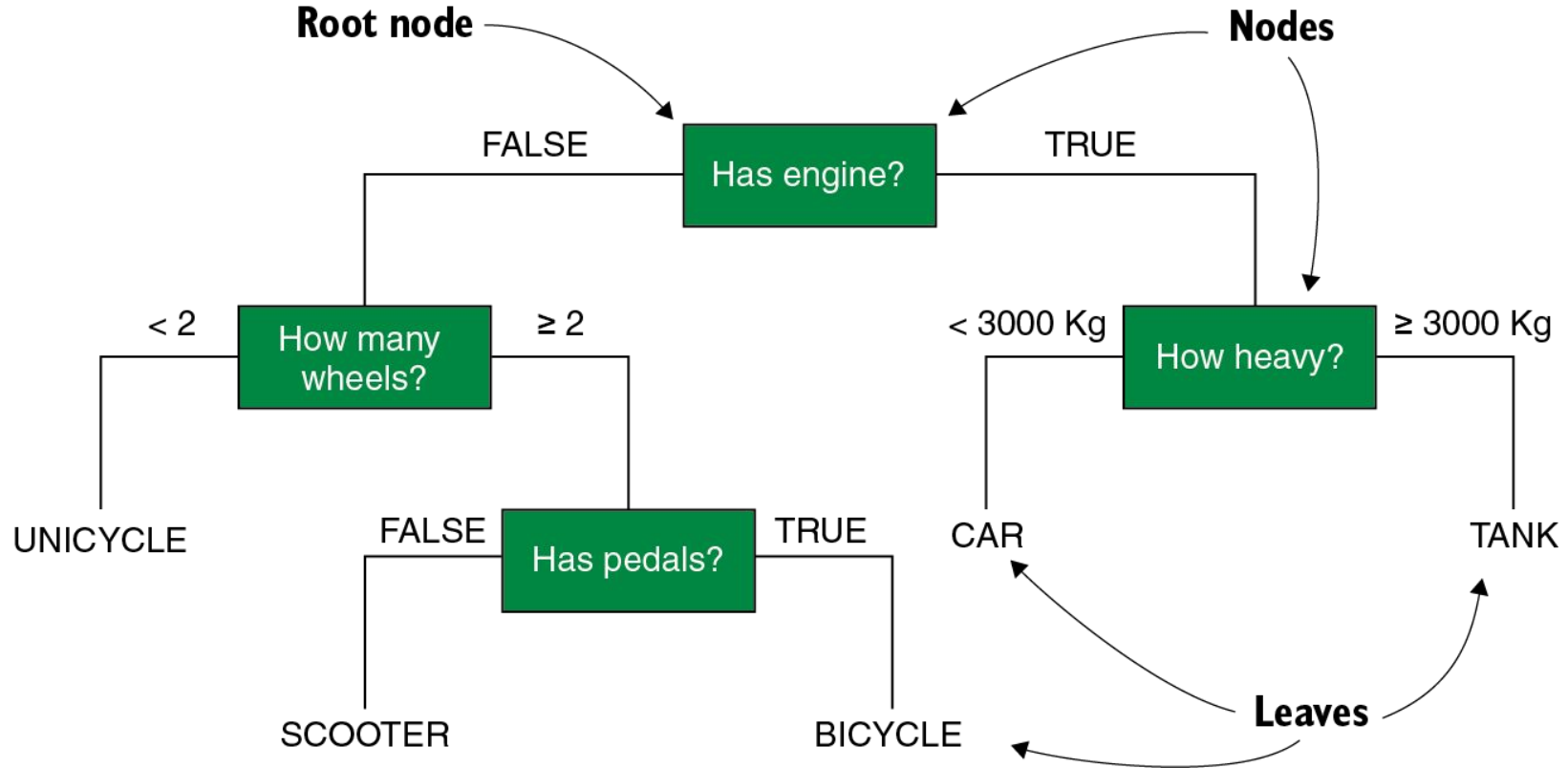
Bootstrapping



Machine Learning

Aprendizaje supervisado

Árboles de decisión



Árboles de decisión

Hiperparámetros

- **min_n:** mínimo de muestras por nodos
- **tree_depth:** pone límite a la profundidad máxima de un árbol. Es un método para detener el algoritmo y evitar *overfitting*
- **cost_complexity:** costo o penalización a los errores de los árboles más complejos. Es un parámetro de parada. Si adopta el enfoque de construir árboles realmente profundos, el valor predeterminado de 0.01 podría ser demasiado restrictivo.

Random Forest

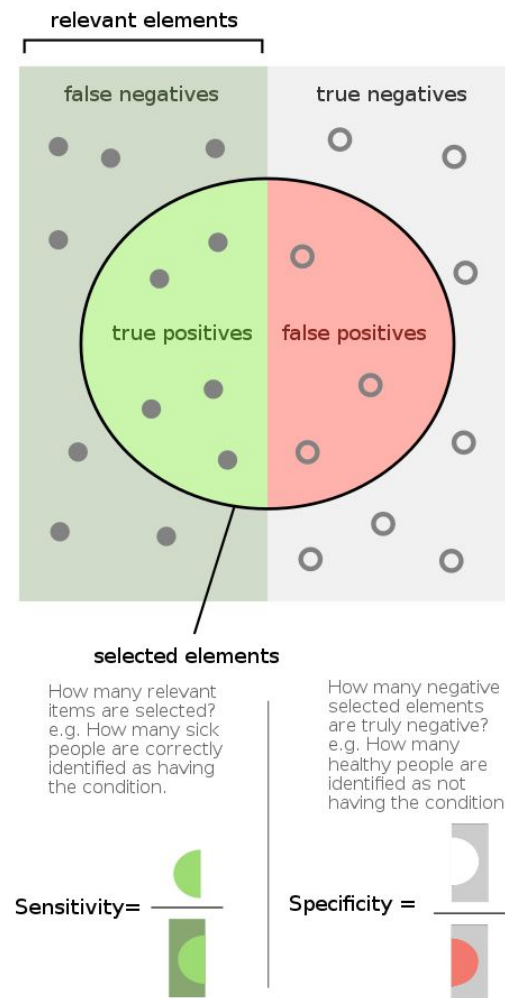
Hiperparámetros

- **mtry**: n° de predictores a muestrearse en cada split de árbol
- **min_n**: n° de observaciones necesarias para seguir dividiendo nodos

Métricas

Matriz de confusión

		Predicted Class		
		Positive	Negative	
Actual Class	Positive	True Positive (TP)	False Negative (FN) Type II Error	Sensitivity $\frac{TP}{(TP + FN)}$
	Negative	False Positive (FP) Type I Error	True Negative (TN)	Specificity $\frac{TN}{(TN + FP)}$
		Precision $\frac{TP}{(TP + FP)}$	Negative Predictive Value $\frac{TN}{(TN + FN)}$	Accuracy $\frac{TP + TN}{(TP + TN + FP + FN)}$



Métricas

RMSE (raíz del error cuadrático medio)

https://es.wikipedia.org/wiki/Ra%C3%ADz_del_error_cuadr%C3%A1tico_medio

Mide la cantidad de error que hay entre dos conjuntos de datos. En otras palabras, compara un valor predicho y un valor observado o conocido. A diferencia del error absoluto medio (MAE), utilizamos RMSE en una variedad de aplicaciones cuando comparamos dos conjuntos de datos.

MAE (error absoluto medio)

https://es.wikipedia.org/wiki/Error_absoluto_medio

Permite evaluar la diferencia entre dos variables continuas. Sirve para cuantificar la precisión de una técnica de predicción.

Métricas

Curva ROC (Receiver Operating Characteristic)

https://es.wikipedia.org/wiki/Curva_ROC

Una curva ROC es una representación gráfica de la sensibilidad frente a la especificidad para un sistema clasificador binario según se varía el umbral de discriminación.

Kappa https://es.wikipedia.org/wiki/Coeficiente_kappa_de_Cohen

Es una medida estadística que ajusta el efecto del azar en la proporción de la concordancia observada para elementos cualitativos (variables categóricas). En general se cree que es una medida más robusta que el simple cálculo del porcentaje de concordancia, ya que κ tiene en cuenta el acuerdo que ocurre por azar.

Veamos algo de código...