



Selección de características

Metodos para restringir el conjunto de atributos con el cuál se entrenará un modelo de Machine Learning con diferentes fines como

- ☐ Reducir dimensionalidad
- ☐ Evitar información redundante, irrelevante o que "engañe" al modelo
- ☐ Disminuir el tiempo de entrenamiento

Selección de características

```
graph TD; A[Selección de características] --> B[Supervisada]; A --> C[No supervisada]; B --> D[Tiene en cuenta las etiquetas del conjunto de datos]; C --> E[No las tiene en cuenta 😊]
```

Supervisada



Tiene en cuenta las etiquetas
del conjunto de datos

No supervisada



No las tiene en cuenta 😊

Supervisada

```
graph TD; A[Supervisada] --- B[Métodos "wrapper"]; A --- C[Métodos de filtro]; A --- D[Métodos intrínsecos];
```

Métodos
"wrapper"

Métodos de
filtro

Métodos
intrínsecos

Supervisada

Métodos
"Wrapper"

Se crean varios modelos con diferentes subconjuntos de atributos y se seleccionan aquellos que resultaron en un mejor desempeño

Selección forward

Selección backward

Classic

RFE (Recursive Feature elimination)

Selección forward

Inicia con un conjunto vacío de atributos y va añadiendo atributos nuevos de forma iterativa

En cada iteración prueba entrenar el modelo con cada uno de los atributos faltantes y se queda con el que se desempeñe mejor

esto ocurre hasta que se cumpla un criterio de mejoramiento en el desempeño del modelo

Sklearn: `SequentialFeatureSelector`

Selección forward

Sklearn

SequentialFeature Selector

```
SequentialFeatureSelector ( model,  
                             k_features = 5,  
                             tol = incremento mínimo del score para continuar  
                             direction = forward  
                             scoring = 'accuracy',  
                             cv = 5 )
```

Selección backward

Classic

RFE (Recursive Feature elimination)

Classic

Igual que el forward pero inicia con todos los atributos y va quitando uno a uno

RFE (Recursive Feature elimination)

→

RFE (Recursive Feature elimination)

Evita entrenar muchos modelos en cada paso de la búsqueda.

Se le da una medida de importancia a las variables. Esto las organiza en un ranking

Top 1	X_{i1}
2	X_{i2}
\vdots	\vdots
k	X_{ik}

En cada paso de la iteración, los atributos menos importantes se eliminan **antes** de reconstruir el modelo

RFE (Recursive Feature elimination)

- 1 Tune/train the model on the training set using all P predictors
- 2 Calculate model performance
- 3 Calculate variable importance or rankings
- 4 **for** each subset size S_i , $i = 1 \dots S$ **do**
- 5 Keep the S_i most important variables
- 6 [Optional] Pre-process the data
- 7 Tune/train the model on the training set using S_i predictors
- 8 Calculate model performance
- 9 [Optional] Recalculate the rankings for each predictor
- 10 **end**
- 11 Calculate the performance profile over the S_i
- 12 Determine the appropriate number of predictors (i.e. the S_i associated with the best performance)
- 13 Fit the final model based on the optimal S_i

Supervisada

Métodos de
filtro

Se utilizan medidas estadísticas para
evaluar la relación entre cada atributo y la
variable de salida.

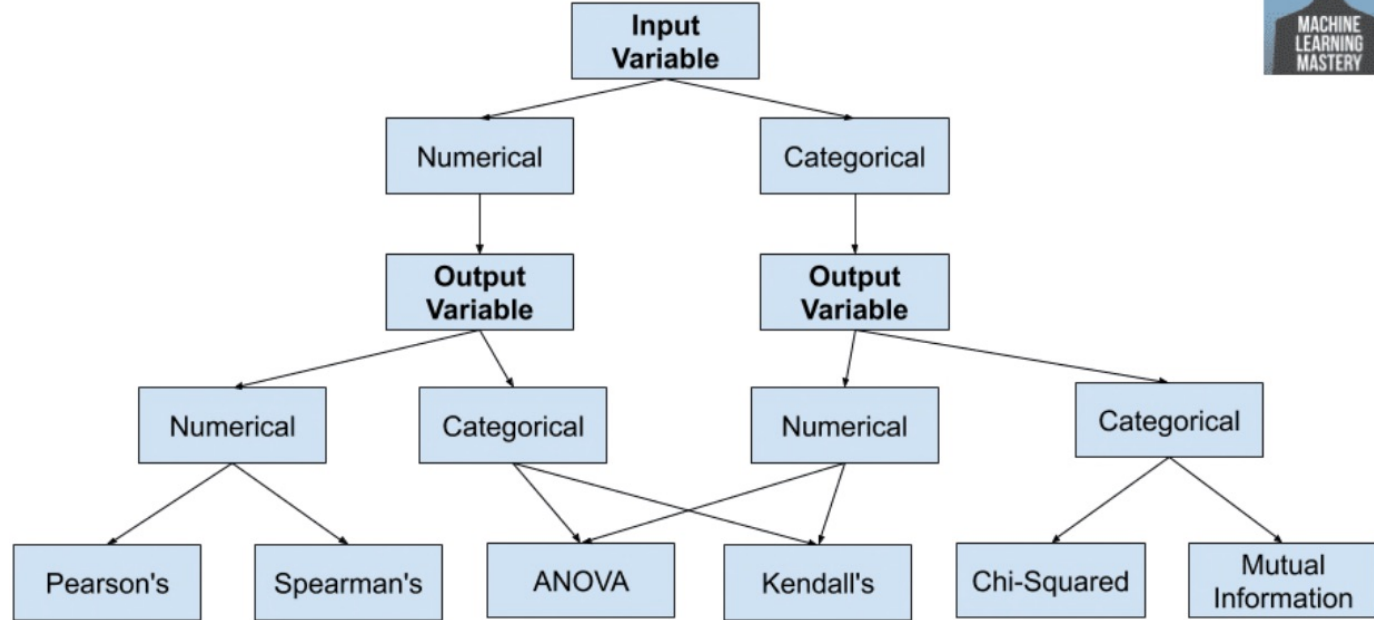
como la correlación

Se mide con un atributo a la vez

Sólo los mejores atributos de acuerdo a un criterio determinado
se utilizarán para entrenar el modelo

Con cuál medida estadística medir la relación de las variables con la variable objetivo?

How to Choose a Feature Selection Method



Pearson

$$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} = \frac{\text{cov}(x, y)}{\sqrt{\text{Var}(x)} \sqrt{\text{Var}(y)}}$$

$$\text{cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - E[X])(y_i - E[Y])$$

Spearman

$$\rho_{R(x), R(y)} = \frac{\text{cov}(R(x), R(y))}{\sigma_{R(x)} \sigma_{R(y)}}$$

R es un ranking entre las variables

Supervisada



```
graph TD; A[Supervisada] --- B[Métodos intrínsecos]
```

Métodos
intrínsecos

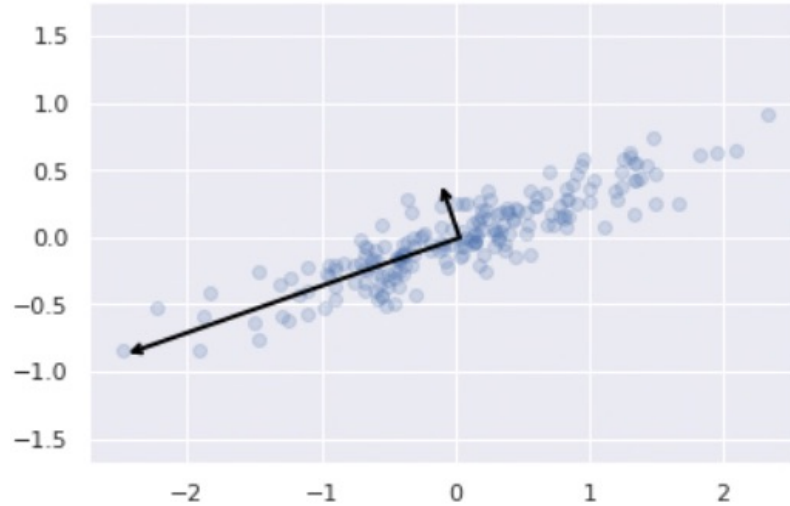
Son aquellos modelos de Machine Learning que durante el entrenamiento seleccionan atributos de manera "natural".

Ej: - Árboles de decisión
- LASSO

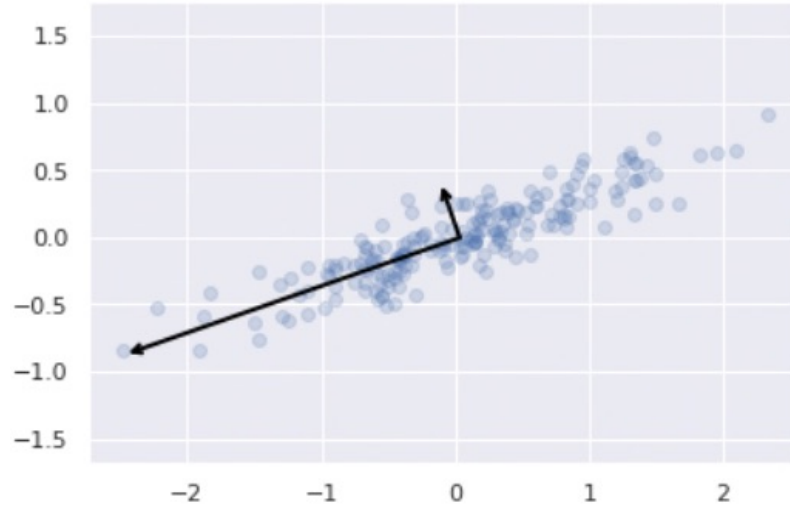
Principal Component Analysis



No supervised



Principal Component Analysis



Dirección de
máxima
varianza

Nuevas coordenadas

PCA permite :

Reducir la dimensionalidad de los datos

Cuando hay muchas variables correlacionadas, se reduce a un número bajo de variables independientes

Cómo se calculan las componentes?

- Se calcula el promedio de las filas

$$\bar{X} = \frac{1}{n} \sum_{j=1}^n x_j$$

X matriz de características

$$X = \begin{bmatrix} \text{---} x_1 \text{---} \\ \vdots \\ \text{---} x_m \text{---} \end{bmatrix}$$

Cómo se calculan las componentes?

- Se calcula el promedio de las filas

$$\bar{X} = \frac{1}{n} \sum_{j=1}^n x_j$$

- Se resta la media \bar{X} a cada fila

$$B = X - \bar{X} \quad [\text{pequeño abuso de notación}]$$

X matriz de características

$$X = \begin{bmatrix} \text{---} x_1 \text{---} \\ \vdots \\ \text{---} x_m \text{---} \end{bmatrix}$$

Cómo se calculan las componentes?

- Se calcula el promedio de las filas

$$\bar{X} = \frac{1}{n} \sum_{j=1}^n x_j$$

- Se resta la media \bar{X} a cada fila

$$B = X - \bar{X} \quad [\text{pequeño abuso de notación}]$$

- Consideremos la matriz de covarianza de las filas de B

$$C = B^T B$$

X matriz de características

$$X = \begin{bmatrix} \text{---} x_1 \text{---} \\ \vdots \\ \text{---} x_m \text{---} \end{bmatrix}$$

Cómo se calculan las componentes?

X matriz de características

- Se calcula el promedio de las filas

$$\bar{X} = \frac{1}{n} \sum_{j=1}^n x_j$$

$$X = \begin{bmatrix} \text{---} x_1 \text{---} \\ \vdots \\ \text{---} x_m \text{---} \end{bmatrix}$$

- Se resta la media \bar{X} a cada fila

$$B = X - \bar{X} \quad [\text{pequeño abuso de notación}]$$

- Consideremos la matriz de covarianza de las filas de B

$$C = B^T B$$

- Se calculan los vectores/valores propios.

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$$

$$C = V D V^{-1}$$

$$\underbrace{\begin{bmatrix} | & & | \\ e_1 & \dots & e_n \\ | & & | \end{bmatrix}}_V \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} \begin{bmatrix} V^{-1} \end{bmatrix}$$

El vector propio e_1 asociado al mayor valor propio indica la dirección de la primera componente, de la misma forma los demás.