

PREPROCESSING

PREPROCESSING

Encodage

$\begin{bmatrix} \textit{chien} \\ \textit{chat} \\ \textit{chien} \\ \textit{oiseau} \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0 \\ 2 \end{bmatrix}$

PREPROCESSING

Encodage

$\begin{bmatrix} \textit{chien} \\ \textit{chat} \\ \textit{chien} \\ \textit{oiseau} \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0 \\ 2 \end{bmatrix}$

Normalisation

$\begin{bmatrix} 2 \\ 10 \\ 4 \\ 6 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0.25 \\ 0.5 \end{bmatrix}$

PREPROCESSING

Encodage

$\begin{bmatrix} \textit{chien} \\ \textit{chat} \\ \textit{chien} \\ \textit{oiseau} \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0 \\ 2 \end{bmatrix}$

Normalisation

$\begin{bmatrix} 2 \\ 10 \\ 4 \\ 6 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0.25 \\ 0.5 \end{bmatrix}$

Imputation

$\begin{bmatrix} 2 \\ 1 \\ 3 \\ \textit{'nan'} \end{bmatrix} \rightarrow \begin{bmatrix} 2 \\ 1 \\ 3 \\ 2 \end{bmatrix}$

PREPROCESSING

Encodage

$\begin{bmatrix} \text{chien} \\ \text{chat} \\ \text{chien} \\ \text{oiseau} \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0 \\ 2 \end{bmatrix}$

Normalisation

$\begin{bmatrix} 2 \\ 10 \\ 4 \\ 6 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0.25 \\ 0.5 \end{bmatrix}$

Imputation

$\begin{bmatrix} 2 \\ 1 \\ 3 \\ \text{'nan'} \end{bmatrix} \rightarrow \begin{bmatrix} 2 \\ 1 \\ 3 \\ 2 \end{bmatrix}$

Sélection

$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 1 \\ 3 \\ 4 \end{bmatrix}$

PREPROCESSING

Encodage

$\begin{bmatrix} \text{chien} \\ \text{chat} \\ \text{chien} \\ \text{oiseau} \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0 \\ 2 \end{bmatrix}$

Normalisation

$\begin{bmatrix} 2 \\ 10 \\ 4 \\ 6 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0.25 \\ 0.5 \end{bmatrix}$

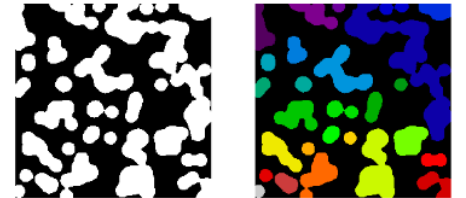
Imputation

$\begin{bmatrix} 2 \\ 1 \\ 3 \\ \text{'nan'} \end{bmatrix} \rightarrow \begin{bmatrix} 2 \\ 1 \\ 3 \\ 2 \end{bmatrix}$

Sélection

$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 1 \\ 3 \\ 4 \end{bmatrix}$

Extraction



PREPROCESSING

Encodage

$\begin{bmatrix} \text{chien} \\ \text{chat} \\ \text{chien} \\ \text{oiseau} \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0 \\ 2 \end{bmatrix}$

Normalisation

$\begin{bmatrix} 2 \\ 10 \\ 4 \\ 6 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0.25 \\ 0.5 \end{bmatrix}$

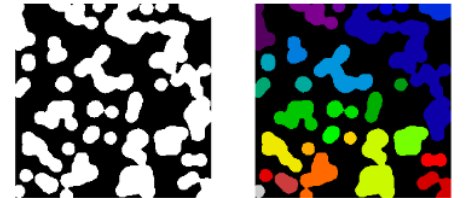
Imputation

$\begin{bmatrix} 2 \\ 1 \\ 3 \\ \text{'nan'} \end{bmatrix} \rightarrow \begin{bmatrix} 2 \\ 1 \\ 3 \\ 2 \end{bmatrix}$

Sélection

$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 1 \\ 3 \\ 4 \end{bmatrix}$

Extraction



`sklearn.preprocessing`

PREPROCESSING

Encodage

$\begin{bmatrix} \text{chien} \\ \text{chat} \\ \text{chien} \\ \text{oiseau} \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0 \\ 2 \end{bmatrix}$

$\begin{bmatrix} 2 \\ 10 \\ 4 \\ 6 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0.25 \\ 0.5 \end{bmatrix}$

`sklearn.preprocessing`

Normalisation

Imputation

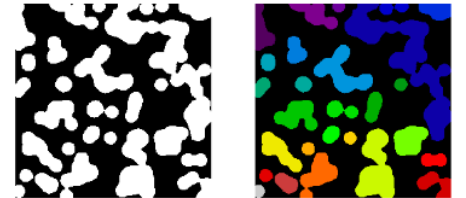
$\begin{bmatrix} 2 \\ 1 \\ 3 \\ \text{'nan'} \end{bmatrix} \rightarrow \begin{bmatrix} 2 \\ 1 \\ 3 \\ 2 \end{bmatrix}$

`sklearn.impute`

Sélection

$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 1 \\ 3 \\ 4 \end{bmatrix}$

Extraction



PREPROCESSING

Encodage

$\begin{bmatrix} \text{chien} \\ \text{chat} \\ \text{chien} \\ \text{oiseau} \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0 \\ 2 \end{bmatrix}$

$\begin{bmatrix} 2 \\ 10 \\ 4 \\ 6 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 1 \\ 0.25 \\ 0.5 \end{bmatrix}$

`sklearn.preprocessing`

Normalisation

Imputation

$\begin{bmatrix} 2 \\ 1 \\ 3 \\ \text{'nan'} \end{bmatrix} \rightarrow \begin{bmatrix} 2 \\ 1 \\ 3 \\ 2 \end{bmatrix}$

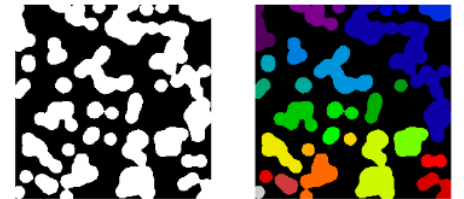
`sklearn.impute`

Sélection

$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ 1 \\ 3 \\ 4 \end{bmatrix}$

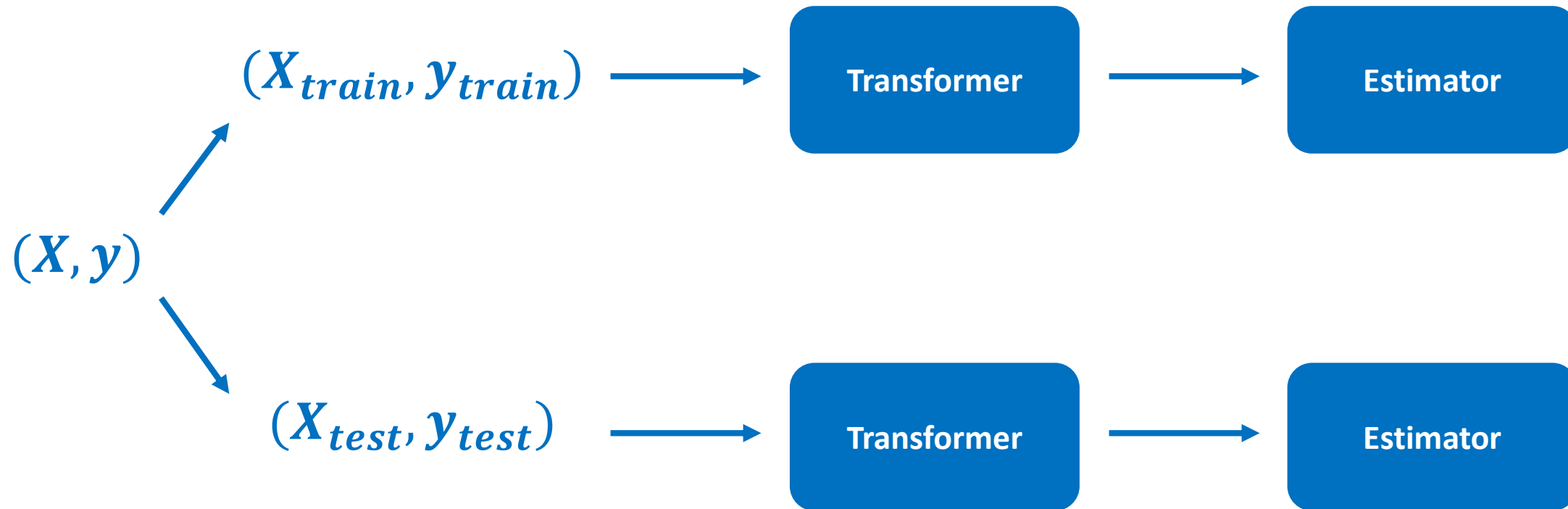
`sklearn
feature_selection`

Extraction



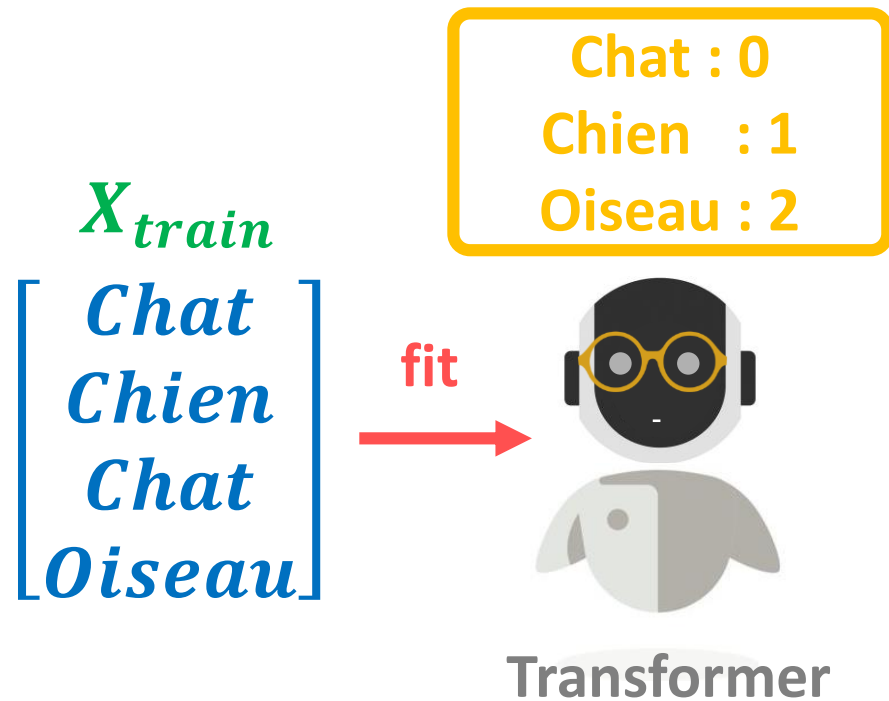
`sklearn
feature_extraction`

TRANSFORMER ET ESTIMATOR



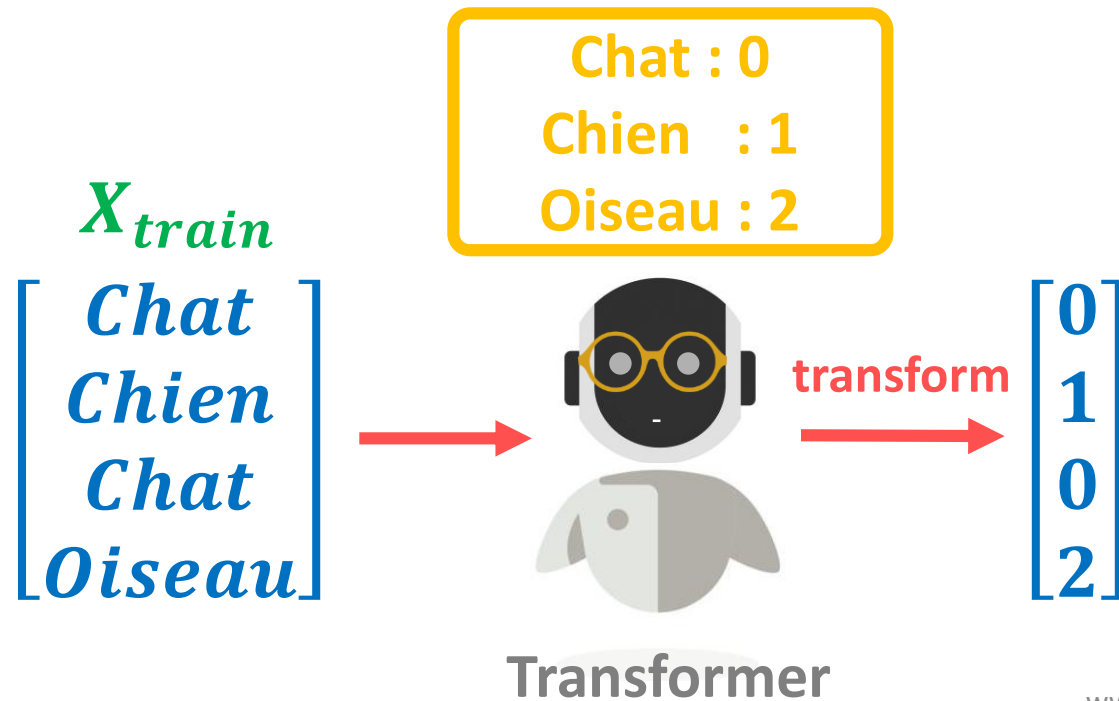
TRANSFORMER

fit(X_{train}) : développe une fonction de **transformation** à partir de X_{train}



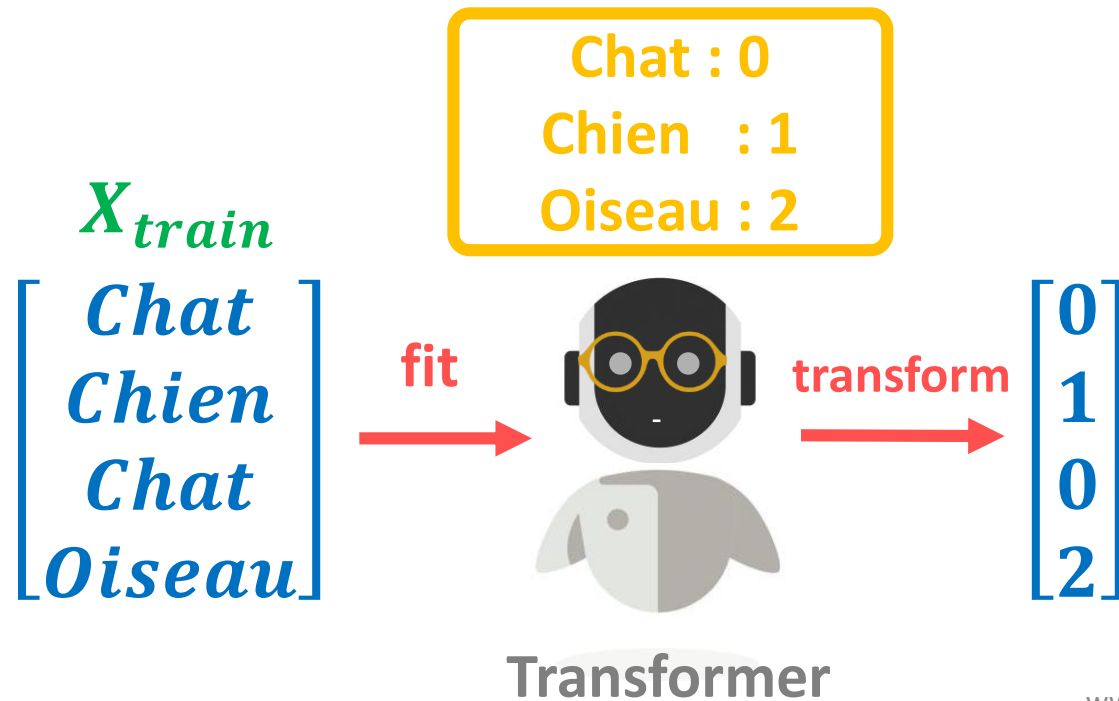
TRANSFORMER

transform(X) : applique la **transformation** sur les données X_{train} , X_{test} et toutes autres données futures.



TRANSFORMER

fit_transform(X_{train}) : développe la fonction de **transformation** puis l'utilise pour transformer X_{train}



ENCODAGE ORDINAL

Associe chaque catégorie ou classe d'une variable à une valeur **décimale unique**.

ENCODAGE ORDINAL

Associe chaque catégorie ou classe d'une variable à une valeur **décimale unique**.

<i>Chat</i>	-----	0
<i>Chat</i>	-----	0
<i>Chien</i>	-----	1
<i>Oiseau</i>	-----	2
<i>Chien</i>	-----	1

ENCODAGE ORDINAL

Associe chaque catégorie ou classe d'une variable à une valeur **décimale unique**.

<i>Chat</i>	-----	0
<i>Chat</i>	-----	0
<i>Chien</i>	-----	1
<i>Oiseau</i>	-----	2
<i>Chien</i>	-----	1

→ `LabelEncoder()`

→ `OrdinalEncoder()`

ENCODAGE ORDINAL

<i>Chat</i>	-----	<i>0</i>
<i>Chat</i>	-----	<i>0</i>
<i>Chien</i>	-----	<i>1</i>
<i>Oiseau</i>	-----	<i>2</i>
<i>Chien</i>	-----	<i>1</i>

$0 < 1 < 2$
 $\Rightarrow \textit{Chat} < \textit{Chien} < \textit{Oiseau}$

??? Ca n'a pas de sens...

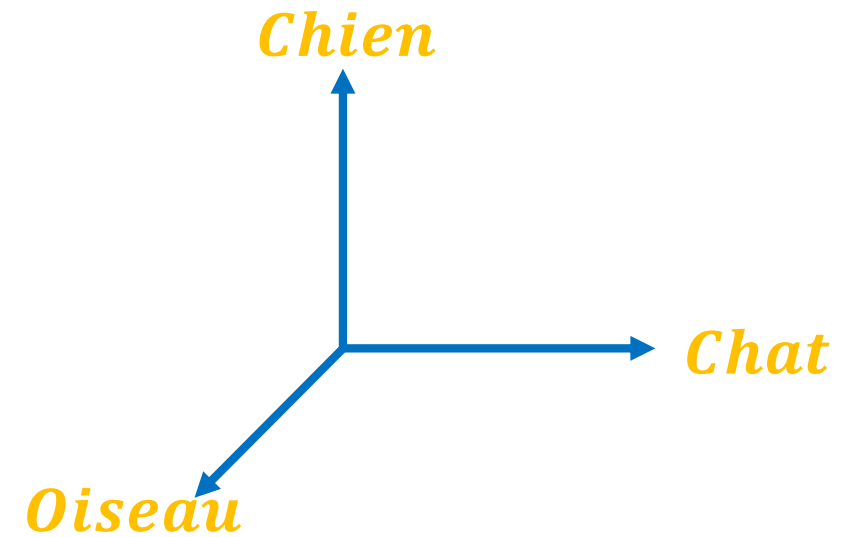
ENCODAGE ONE HOT

		Chat	Chien	Oiseau
Chat	-----	1	0	0
Chat	-----	1	0	0
Chien	-----	0	1	0
Oiseau	-----	0	0	1
Chien	-----	0	1	0

- LabelBinarizer()
- MultiLabelBinarizer()
- OneHotEncoder()

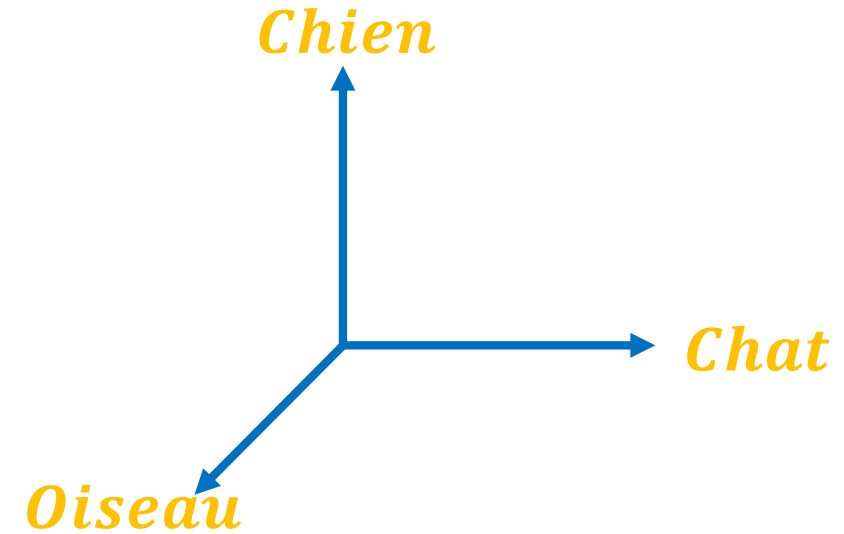
ENCODAGE ONE HOT

<i>Chat</i>	-----	<i>Chat</i>	<i>Chien</i>	<i>Oiseau</i>
<i>Chat</i>	-----	1	0	0
<i>Chat</i>	-----	1	0	0
<i>Chien</i>	-----	0	1	0
<i>Oiseau</i>	-----	0	0	1
<i>Chien</i>	-----	0	1	0



ENCODAGE ONE HOT

<i>Chat</i>	-----	<i>Chat</i>	<i>Chien</i>	<i>Oiseau</i>
<i>Chat</i>	-----	1	0	0
<i>Chien</i>	-----	1	0	0
<i>Oiseau</i>	-----	0	1	0
<i>Chien</i>	-----	0	0	1
	-----	0	1	0

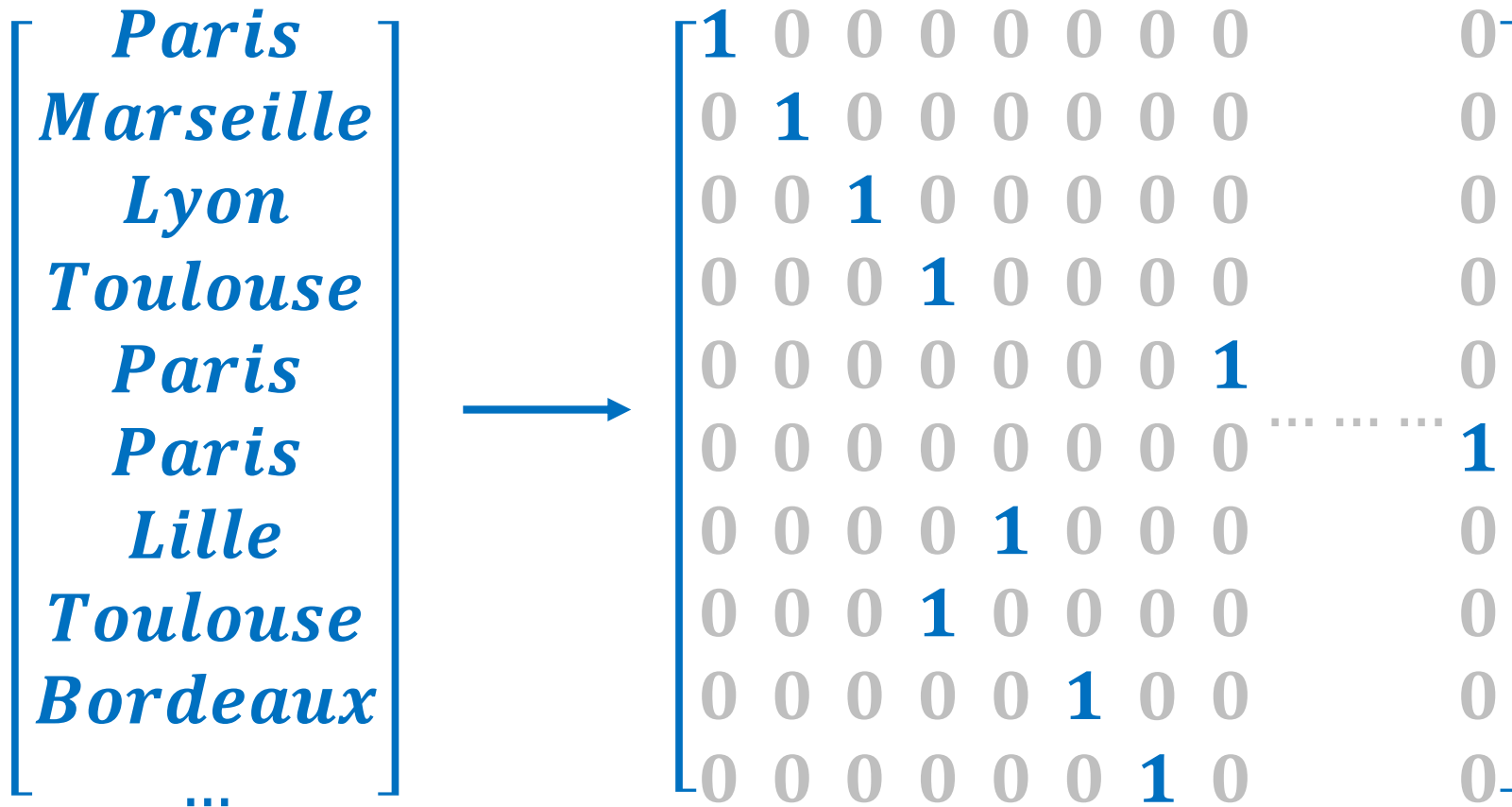


- `LabelBinarizer()`
- `MultiLabelBinarizer()`
- `OneHotEncoder()`

ENCODAGE ONE HOT

<i>Paris</i>	1	0	0	0	0	0	0	0	0
<i>Marseille</i>	0	1	0	0	0	0	0	0	0
<i>Lyon</i>	0	0	1	0	0	0	0	0	0
<i>Toulouse</i>	0	0	0	1	0	0	0	0	0
<i>Paris</i>	0	0	0	0	0	0	0	1	0
<i>Paris</i>	0	0	0	0	0	0	0	0	1
<i>Lille</i>	0	0	0	0	1	0	0	0	0
<i>Toulouse</i>	0	0	0	1	0	0	0	0	0
<i>Bordeaux</i>	0	0	0	0	0	1	0	0	0
...	0	0	0	0	0	0	1	0	0

ENCODAGE ONE HOT



Sparse Matrix

$$\begin{bmatrix} a & 0 & 0 & 0 & 0 \\ 0 & 0 & b & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & c \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Sparse Matrix

$$\begin{bmatrix} a & 0 & 0 & 0 & 0 \\ 0 & 0 & b & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & c \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

values = []

rows = []

cols = []

Sparse Matrix

$$\begin{bmatrix} a & 0 & 0 & 0 & 0 \\ 0 & 0 & b & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & c \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

values = $[a, b, c]$

rows = $[\quad]$

cols = $[\quad]$

Sparse Matrix

$$\begin{matrix} & 0 \\ 0 & \boxed{a} & 0 & 0 & 0 & 0 \\ & 0 & 0 & b & 0 & 0 \\ & 0 & 0 & 0 & 0 & 0 \\ & 0 & 0 & 0 & 0 & c \\ & 0 & 0 & 0 & 0 & 0 \end{matrix}$$

values = $[a, b, c]$

rows = $[0]$

cols = $[0]$

Sparse Matrix

$$\begin{matrix} & & 2 \\ 1 & \begin{bmatrix} a & 0 & 0 & 0 & 0 \\ 0 & 0 & b & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & c \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \end{matrix}$$

values = $[a, b, c]$

rows = $[0, 1]$

cols = $[0, 2]$

Sparse Matrix

$$\begin{matrix} & & & & 4 \\ \begin{matrix} 3 \\ \left[\begin{array}{ccccc} a & 0 & 0 & 0 & 0 \\ 0 & 0 & b & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & c \\ 0 & 0 & 0 & 0 & 0 \end{array} \right] \end{matrix} \end{matrix}$$

values = $[a, b, c]$

rows = $[0, 1, 3]$

cols = $[0, 2, 4]$

ENCODING

```
preprocessing.Binarizer([threshold, copy])
preprocessing.FunctionTransformer([func, ...])
preprocessing.KBinsDiscretizer([n_bins, ...])
preprocessing.KernelCenterer()
preprocessing.LabelBinarizer([neg_label, ...])
preprocessing.LabelEncoder
preprocessing.MultiLabelBinarizer([classes, ...])
preprocessing.MaxAbsScaler([copy])
preprocessing.MinMaxScaler([feature_range, copy])
preprocessing.Normalizer([norm, copy])
preprocessing.OneHotEncoder([categories, ...])
preprocessing.OrdinalEncoder([categories, dtype])
preprocessing.PolynomialFeatures([degree, ...])
preprocessing.PowerTransformer([method, ...])
preprocessing.QuantileTransformer([...])
preprocessing.RobustScaler([with_centering, ...])
preprocessing.StandardScaler([copy, ...])
```

ENCODING

```
preprocessing.Binarizer([threshold, copy])
preprocessing.FunctionTransformer([func, ...])
preprocessing.KBinsDiscretizer([n_bins, ...])
preprocessing.KernelCenterer()
preprocessing.LabelBinarizer([neg_label, ...])
preprocessing.LabelEncoder
preprocessing.MultiLabelBinarizer([classes, ...])
preprocessing.MaxAbsScaler([copy])
preprocessing.MinMaxScaler([feature_range, copy])
preprocessing.Normalizer([norm, copy])
preprocessing.OneHotEncoder([categories, ...])
preprocessing.OrdinalEncoder([categories, dtype])
preprocessing.PolynomialFeatures([degree, ...])
preprocessing.PowerTransformer([method, ...])
preprocessing.QuantileTransformer([...])
preprocessing.RobustScaler([with_centering, ...])
preprocessing.StandardScaler([copy, ...])
```

www.machinelearningia.com

Encodage
Ordinal

Encodage
One-Hot

y

LabelEncoder

LabelBinarizer
(et MultiLabelBinarizer)

X

OrdinalEncoder

OneHotEncoder

LABEL ENCODER

Encode chaque classe de la *variable y*
en une **valeur numérique** ($0, n_classe-1$)

y

<i>Chat</i>	<i>0</i>
<i>Chien</i>	<i>1</i>
<i>Chat</i>	<i>0</i>
<i>Oiseau</i>	<i>2</i>

Note : la méthode *inverse_transform*
permet de décoder les données

ORDINAL ENCODER

Encode les catégories des *variables X*
en **valeurs numériques** ($0, n_{\text{classe}}-1$)

X_1	X_2	
Chat	Poils	0 1
Chien	Poils	1 1
Chat	Poils	0 1
Oiseau	Plumes	2 0

*Note : C'est l'équivalent de **LabelEncoder**,
mais pour les **features X***

LABEL BINARIZER

Encode chaque classe de la *variable y*
en **One-Hot**

y

$$\begin{bmatrix} \text{Chat} \\ \text{Chien} \\ \text{Chat} \\ \text{Oiseau} \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Note : la méthode *inverse_transform*
permet de décoder les données

ONEHOT ENCODER

Encode les catégories des *variables X*
en **One-Hot**

X_1	X_2	
<i>Chat</i>	<i>Poils</i>	<i>Chat</i>
<i>Chien</i>	<i>Poils</i>	<i>Chien</i>
<i>Chat</i>	<i>Poils</i>	<i>Oiseau</i>
<i>Oiseau</i>	<i>Plumes</i>	<i>Poils</i>
		<i>Plumes</i>

→

<i>Chat</i>	<i>Chien</i>	<i>Oiseau</i>	<i>Poils</i>	<i>Plumes</i>
1	0	0	1	0
0	1	0	1	0
1	0	0	1	0
0	0	1	0	1

Note : C'est l'équivalent de *LabelBinarizer*,
mais pour les *features X*

ENCODING

Encodage
Ordinal

Encodage
One-Hot

<i>y</i>	LabelEncoder	LabelBinarizer (et MultiLabelBinarizer)
<i>X</i>	OrdinalEncoder	OneHotEncoder

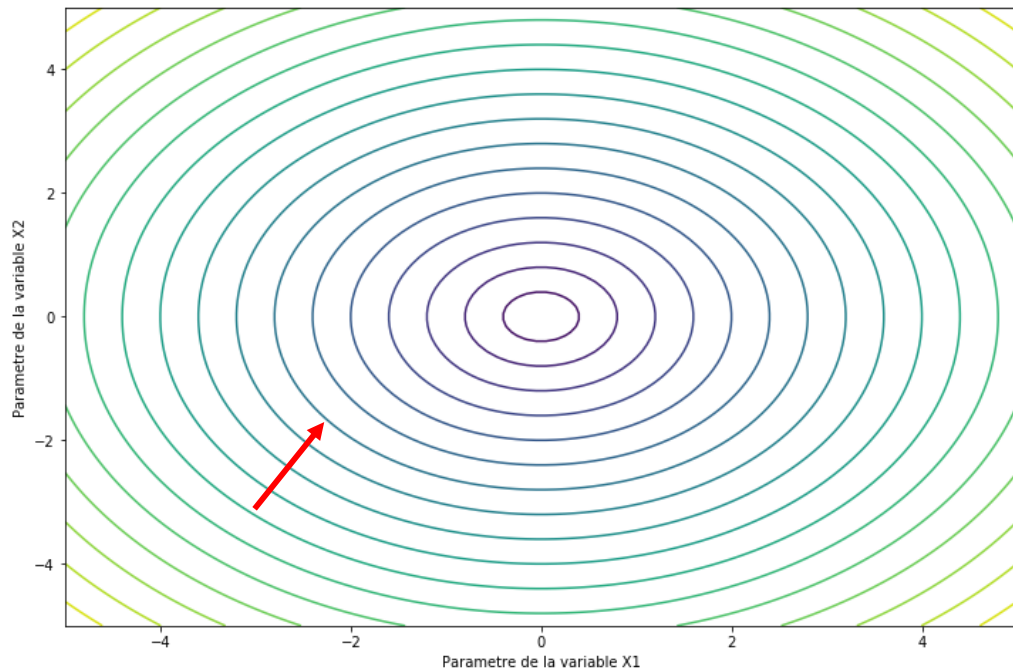
Notes :

- ***Inverse_transform** pour vos prédictions*
- ***sparse=True** pour OneHotEncoder*
- *Attention : Sklearn v22.0 gère mal les catégories inconnues.*

NORMALISATION

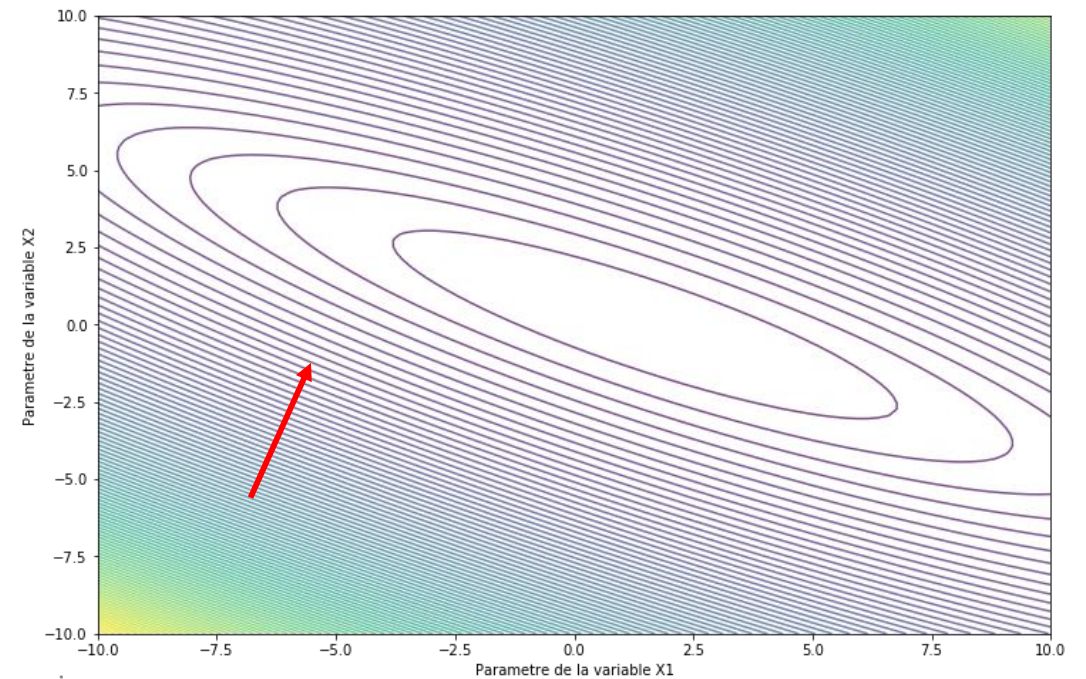
Avec Normalisation

Fonction Cout



Sans Normalisation

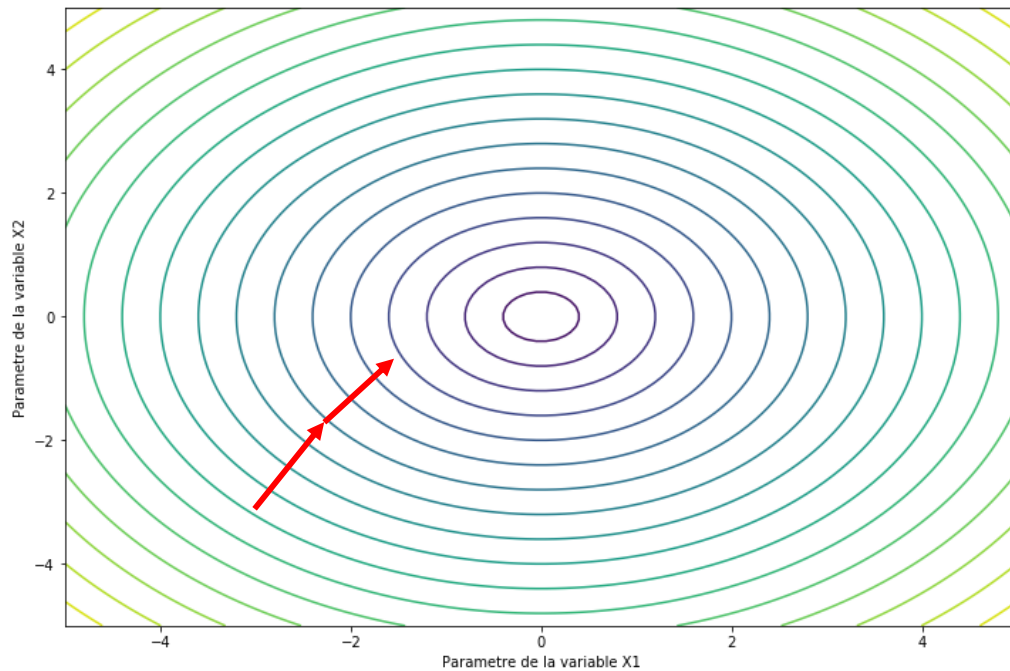
Fonction Cout



NORMALISATION

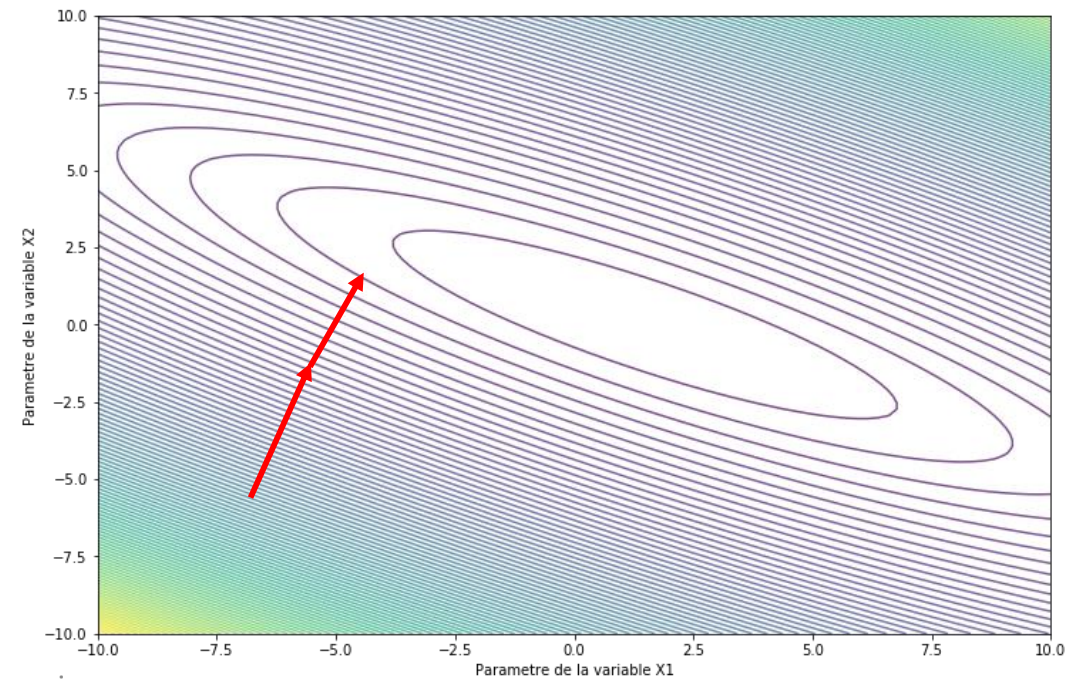
Avec Normalisation

Fonction Cout



Sans Normalisation

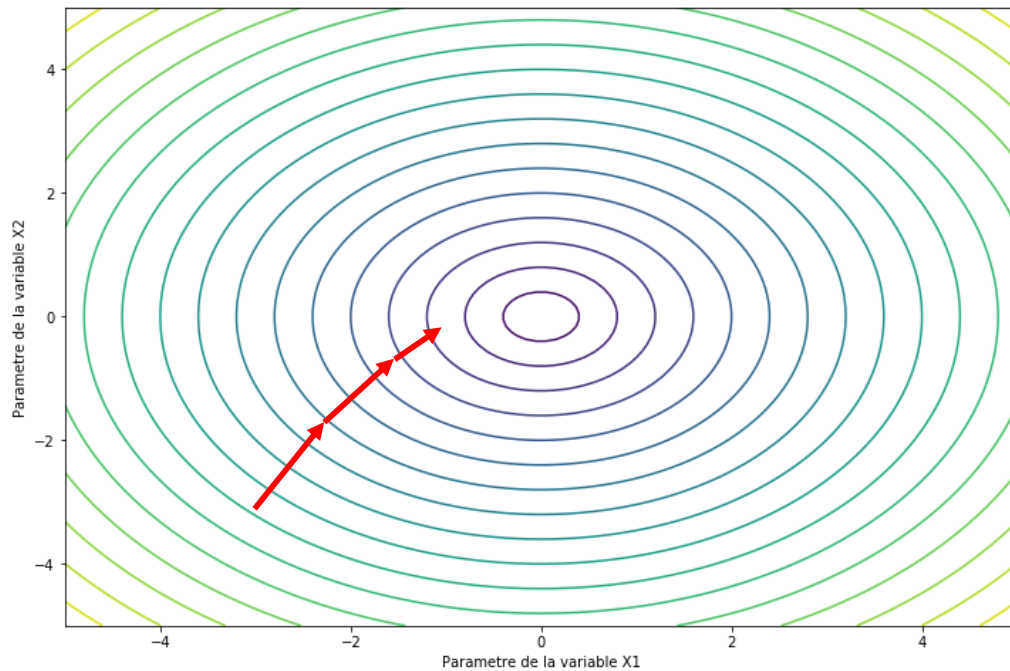
Fonction Cout



NORMALISATION

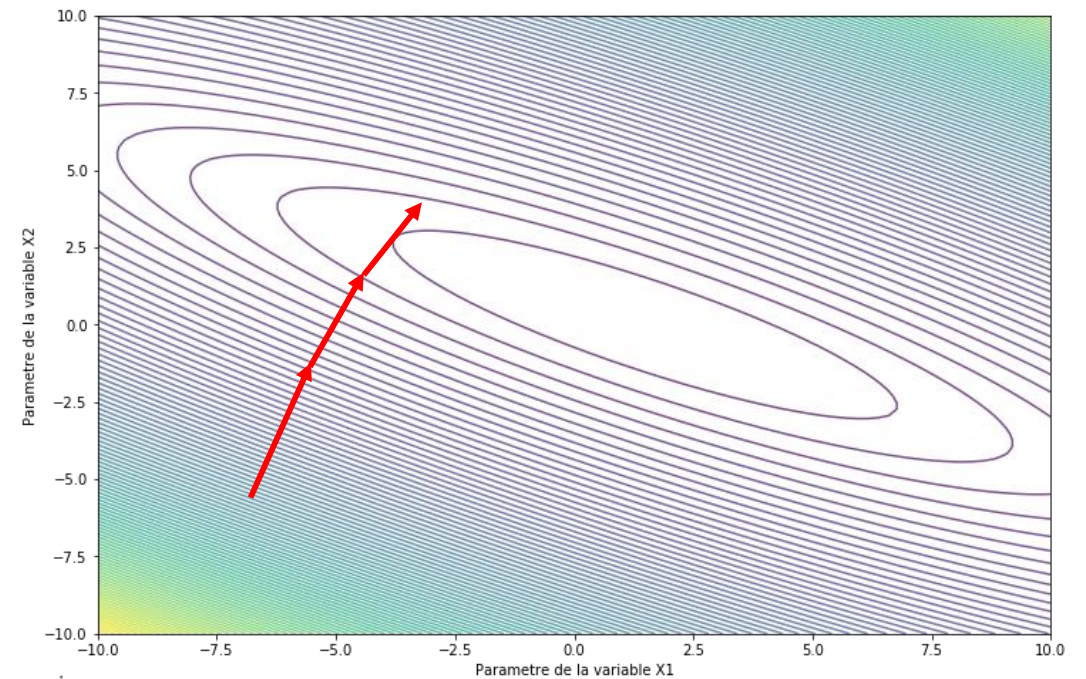
Avec Normalisation

Fonction Cout



Sans Normalisation

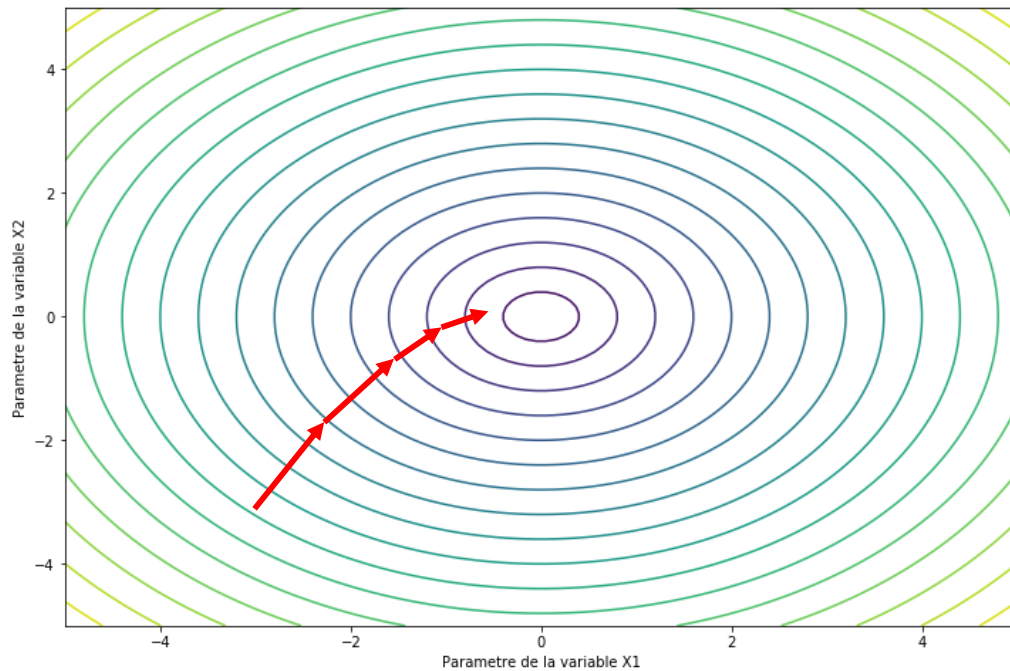
Fonction Cout



NORMALISATION

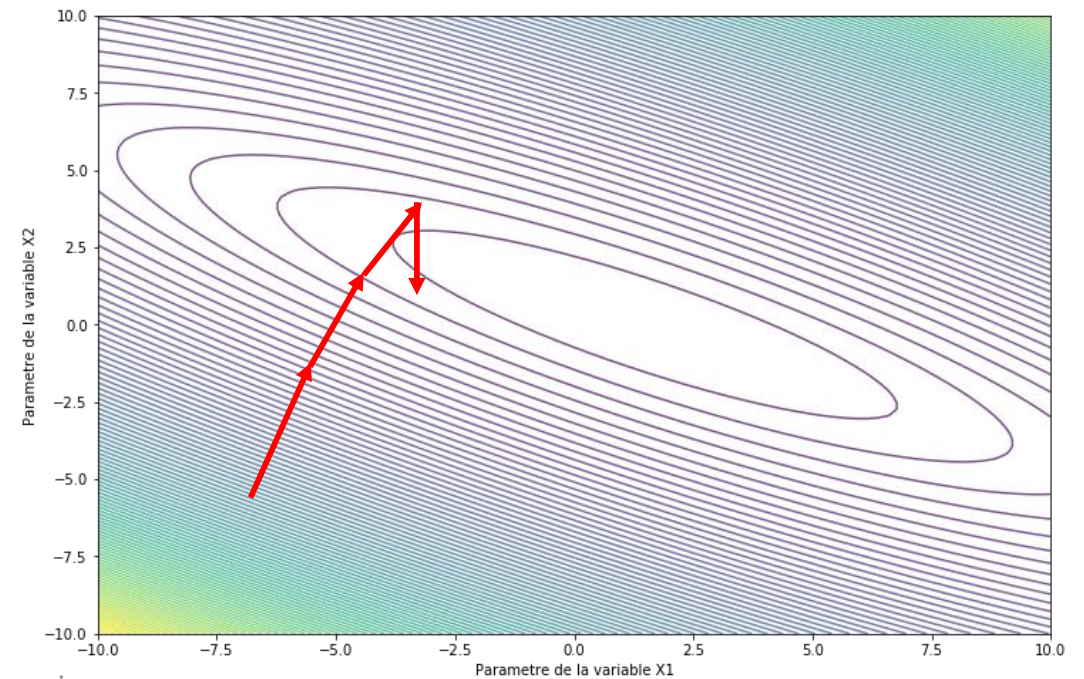
Avec Normalisation

Fonction Cout



Sans Normalisation

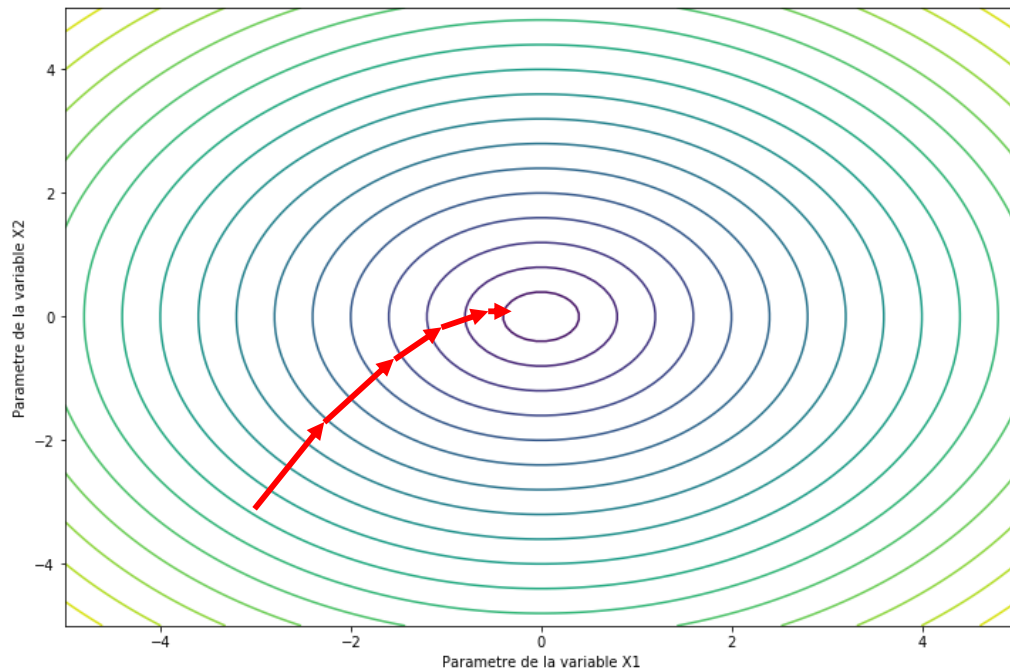
Fonction Cout



NORMALISATION

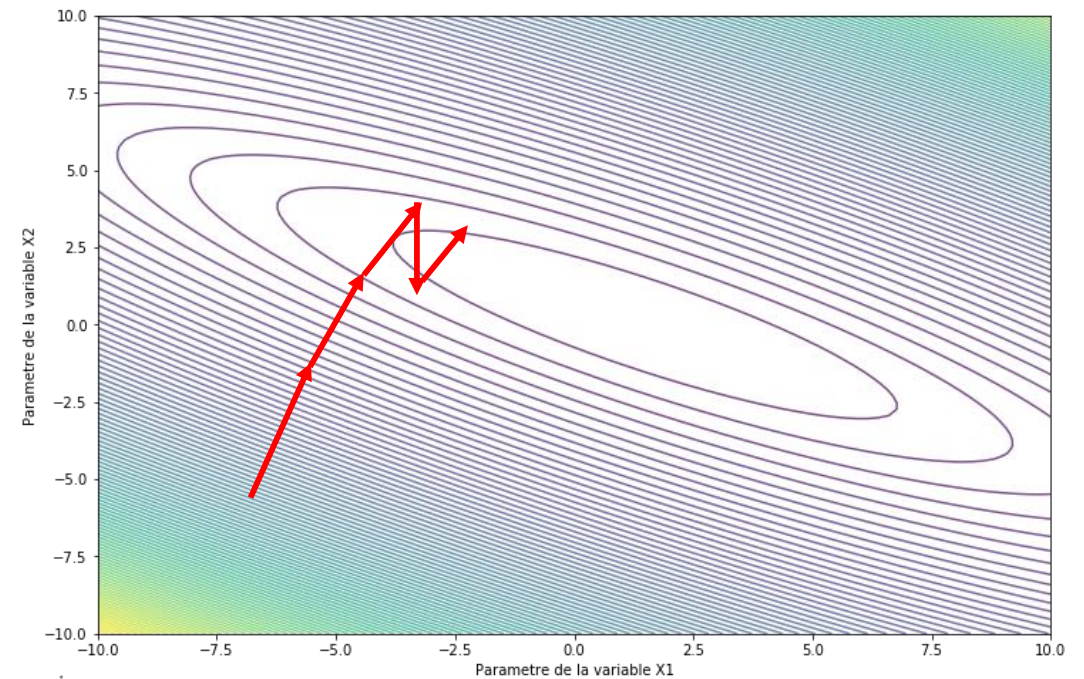
Avec Normalisation

Fonction Cout



Sans Normalisation

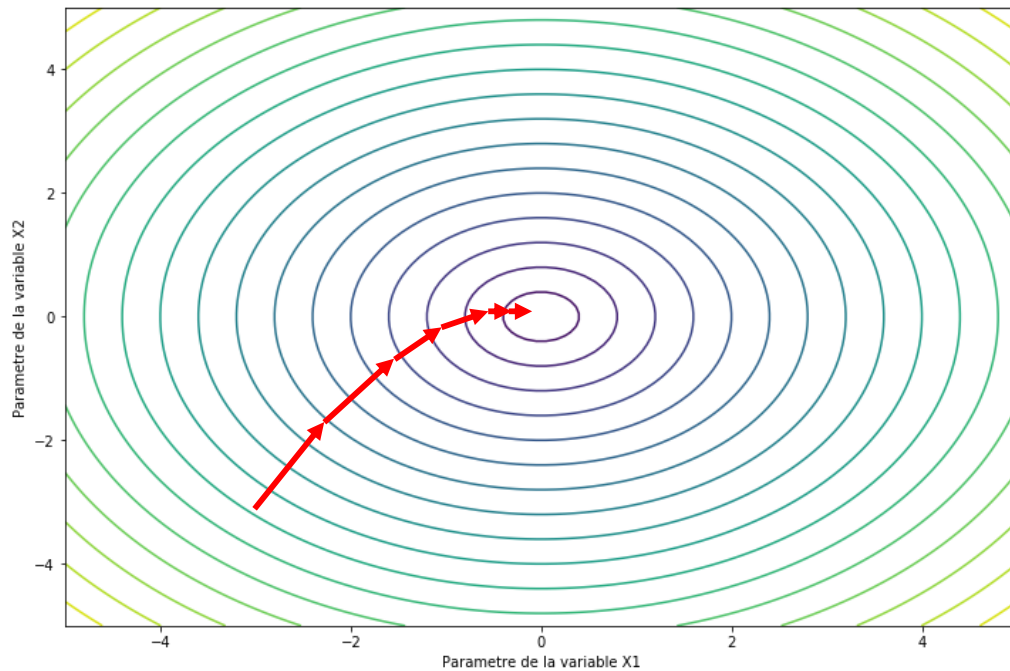
Fonction Cout



NORMALISATION

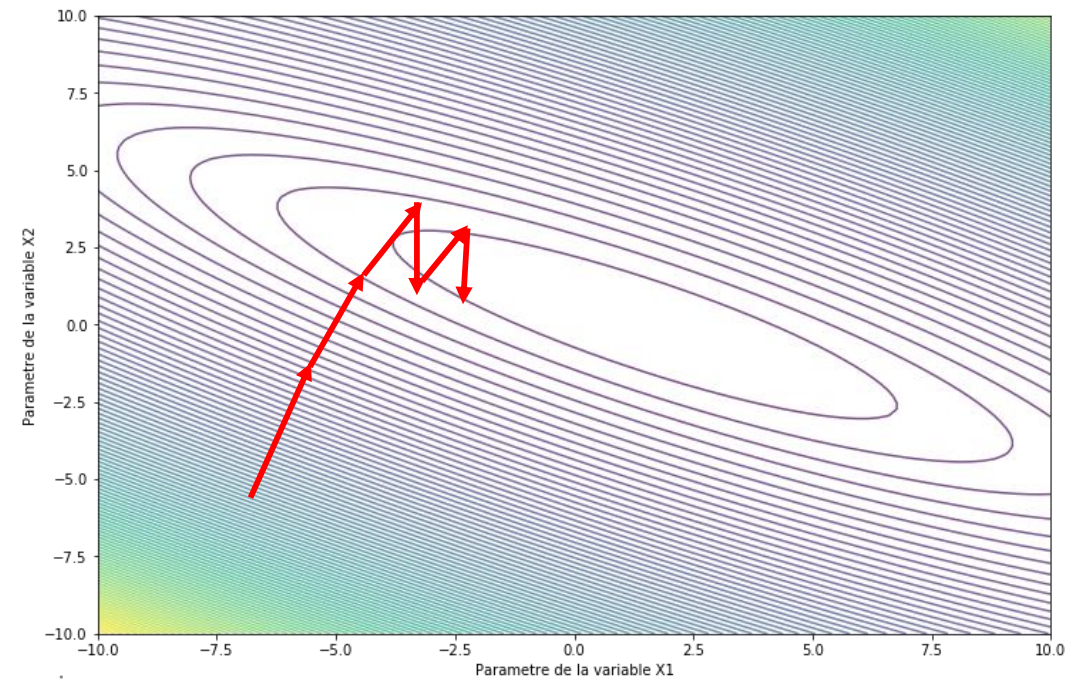
Avec Normalisation

Fonction Cout



Sans Normalisation

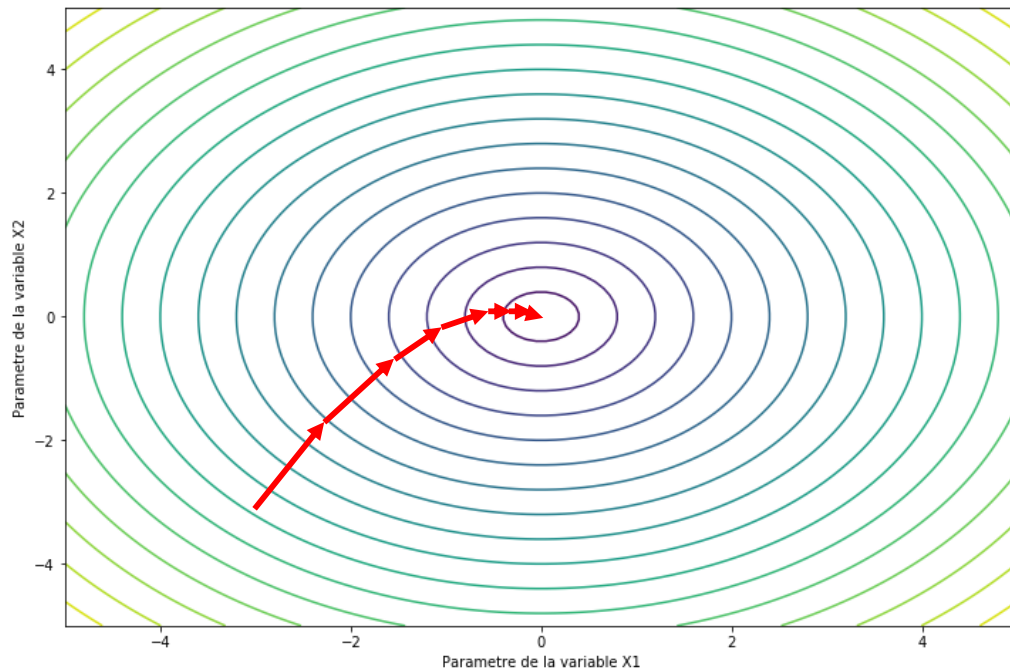
Fonction Cout



NORMALISATION

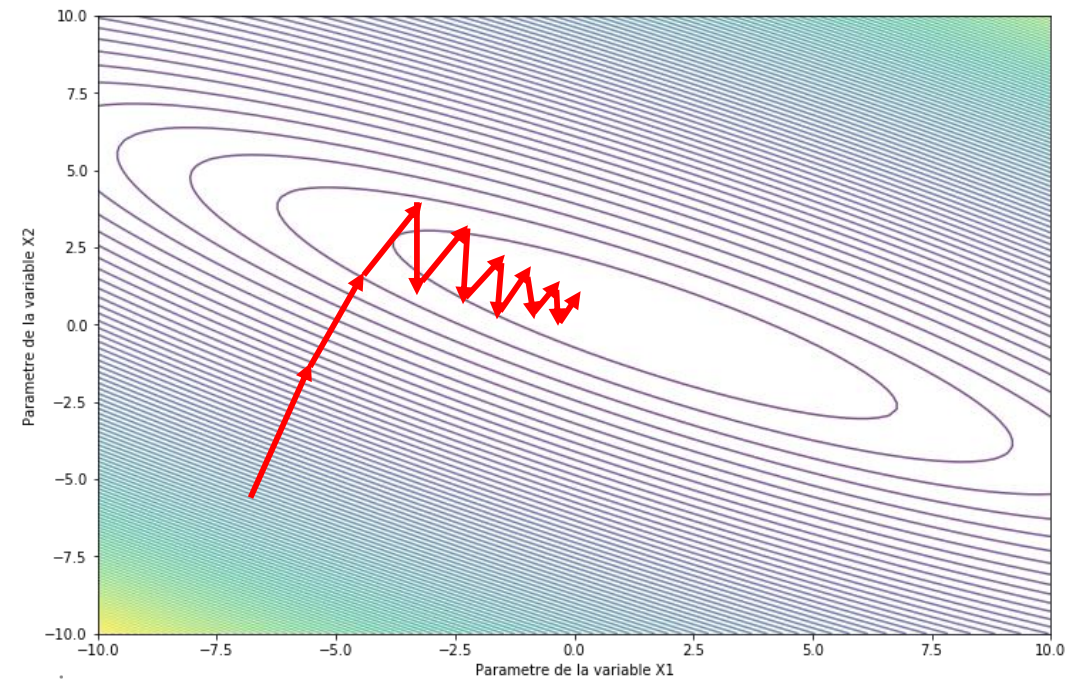
Avec Normalisation

Fonction Cout



Sans Normalisation

Fonction Cout



MINMAXSCALER

Transforme **chaque variable** X de telle sorte à être comprise entre **0 et 1**.

$$X_{scaled} = \frac{X - X_{min}}{X_{max} - X_{min}}$$

Exemple:

$$\begin{bmatrix} 70 \\ 80 \\ 120 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 0.2 \\ 1 \end{bmatrix}$$

Note : Pour remettre vos prédictions à leur véritable échelle : ***inverse_transform***

MINMAXSCALER

fit_transform(X_{train})

$$X_{train_scaled} = \frac{X_{train} - X_{train_min}}{X_{train_max} - X_{train_min}}$$

transform(X_{test})

$$X_{test_scaled} = \frac{X_{test} - X_{train_min}}{X_{train_max} - X_{train_min}}$$

STANDARD SCALER

Standardise chaque variable X : La Moyenne est **nulle** et l'écart type égale **1**

$$X_{scaled} = \frac{X - \mu_X}{\sigma_X}$$

Exemple:

$$\begin{bmatrix} 70 \\ 80 \\ 120 \end{bmatrix} \rightarrow \begin{bmatrix} -0.92 \\ -0.46 \\ 1.38 \end{bmatrix}$$

Note : Pour remettre vos prédictions à leur véritable échelle : ***inverse_transform***

ROBUST SCALER

Transforme **chaque variable** X en étant peu sensible aux **outliers**

$$X_{scaled} = \frac{X - \text{mediane}}{IQR}$$



Exemple:

$$\begin{bmatrix} 70 \\ 80 \\ 120 \end{bmatrix} \rightarrow \begin{bmatrix} -0.4 \\ 0 \\ 1.6 \end{bmatrix}$$

POLYFEATURES

Crée de nouvelles variables **polynômiales** à partir des variables existantes.

Exemple : 1 variable x \rightarrow Polynôme 2

$$\begin{array}{c} x \end{array} \rightarrow \begin{array}{ccc} 1 & x^1 & x^2 \end{array}$$
$$\begin{bmatrix} 1 \\ 2 \\ 0.5 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 4 \\ 1 & 0.5 & 0.25 \end{bmatrix}$$

La machine apprend alors un modèle :

$$f(x) = ax^2 + bx + c$$

POLYFEATURES

Crée de nouvelles variables **polynômiales** à partir des variables existantes.

Exemple : 2 variables $x_1, x_2 \rightarrow$ Polynôme 2

$$x_1, x_2 \rightarrow 1, x_1, x_2, x_1^2, x_1x_2, x_2^2$$

POLYFEATURES

Crée de nouvelles variables **polynômiales** à partir des variables existantes.

Exemple : 2 variables $x_1, x_2 \rightarrow$ Polynôme 2

$$x_1, x_2 \rightarrow 1, x_1, x_2, x_1^2, x_1x_2, x_2^2$$

Ces variables n'étant **pas sur la même échelle**, il faut les **normaliser** avant leur passage dans l'estimateur !

PREPROCESSING

Encode

```
preprocessing.Binarizer([threshold, copy])
preprocessing.FunctionTransformer([func, ...])
preprocessing.KBinsDiscretizer([n_bins, ...])
preprocessing.KernelCenterer()
preprocessing.LabelBinarizer([neg_label, ...])
preprocessing.LabelEncoder
preprocessing.MultiLabelBinarizer([classes, ...])
preprocessing.MaxAbsScaler([copy])
preprocessing.MinMaxScaler([feature_range, copy])
preprocessing.Normalizer([norm, copy])
preprocessing.OneHotEncoder([categories, ...])
preprocessing.OrdinalEncoder([categories, dtype])
preprocessing.PolynomialFeatures([degree, ...])
preprocessing.PowerTransformer([method, ...])
preprocessing.QuantileTransformer([...])
preprocessing.RobustScaler([with_centering, ...])
preprocessing.StandardScaler([copy, ...])
```

PREPROCESSING

Encode

`preprocessing.Binarizer([threshold, copy])`

`preprocessing.FunctionTransformer([func, ...])`

`preprocessing.KBinsDiscretizer([n_bins, ...])`

`preprocessing.KernelCenterer()`

`preprocessing.LabelBinarizer([neg_label, ...])`

`preprocessing.LabelEncoder`

`preprocessing.MultiLabelBinarizer([classes, ...])`

`preprocessing.MaxAbsScaler([copy])`

`preprocessing.MinMaxScaler([feature_range, copy])`

`preprocessing.Normalizer([norm, copy])`

`preprocessing.OneHotEncoder([categories, ...])`

`preprocessing.OrdinalEncoder([categories, dtype])`

`preprocessing.PolynomialFeatures([degree, ...])`

`preprocessing.PowerTransformer([method, ...])`

`preprocessing.QuantileTransformer([...])`

`preprocessing.RobustScaler([with_centering, ...])`

`preprocessing.StandardScaler([copy, ...])`

PREPROCESSING

`preprocessing.Binarizer([threshold, copy])`

`preprocessing.FunctionTransformer([func, ...])`

`preprocessing.KBinsDiscretizer([n_bins, ...])`

`preprocessing.KernelCenterer()`

`preprocessing.LabelBinarizer([neg_label, ...])`

`preprocessing.LabelEncoder`

`preprocessing.MultiLabelBinarizer([classes, ...])`

`preprocessing.MaxAbsScaler([copy])`

`preprocessing.MinMaxScaler([feature_range, copy])`

`preprocessing.Normalizer([norm, copy])`

`preprocessing.OneHotEncoder([categories, ...])`

`preprocessing.OrdinalEncoder([categories, dtype])`

`preprocessing.PolynomialFeatures([degree, ...])`

`preprocessing.PowerTransformer([method, ...])`

`preprocessing.QuantileTransformer([...])`

`preprocessing.RobustScaler([with_centering, ...])`

`preprocessing.StandardScaler([copy, ...])`

Encode

Normalisation

PREPROCESSING

`preprocessing.Binarizer([threshold, copy])`

`preprocessing.FunctionTransformer([func, ...])`

`preprocessing.KBinsDiscretizer([n_bins, ...])`

`preprocessing.KernelCenterer()`

`preprocessing.LabelBinarizer([neg_label, ...])`

`preprocessing.LabelEncoder`

`preprocessing.MultiLabelBinarizer([classes, ...])`

`preprocessing.MaxAbsScaler([copy])`

`preprocessing.MinMaxScaler([feature_range, copy])`

`preprocessing.Normalizer([norm, copy])`

`preprocessing.OneHotEncoder([categories, ...])`

`preprocessing.OrdinalEncoder([categories, dtype])`

`preprocessing.PolynomialFeatures([degree, ...])`

`preprocessing.PowerTransformer([method, ...])`

`preprocessing.QuantileTransformer([...])`

`preprocessing.RobustScaler([with_centering, ...])`

`preprocessing.StandardScaler([copy, ...])`

Encodage

Normalisation

PREPROCESSING

`preprocessing.Binarizer([threshold, copy])`

`preprocessing.FunctionTransformer([func, ...])`

`preprocessing.KBinsDiscretizer([n_bins, ...])`

`preprocessing.KernelCenterer()`

`preprocessing.LabelBinarizer([neg_label, ...])`

`preprocessing.LabelEncoder`

`preprocessing.MultiLabelBinarizer([classes, ...])`

`preprocessing.MaxAbsScaler([copy])`

`preprocessing.MinMaxScaler([feature_range, copy])`

`preprocessing.Normalizer([norm, copy])`

`preprocessing.OneHotEncoder([categories, ...])`

`preprocessing.OrdinalEncoder([categories, dtype])`

`preprocessing.PolynomialFeatures([degree, ...])`

`preprocessing.PowerTransformer([method, ...])`

`preprocessing.QuantileTransformer([...])`

`preprocessing.RobustScaler([with_centering, ...])`

`preprocessing.StandardScaler([copy, ...])`

Encodage

Normalisation

Création de polynômes

PREPROCESSING

`preprocessing.Binarizer([threshold, copy])`

`preprocessing.FunctionTransformer([func, ...])`

`preprocessing.KBinsDiscretizer([n_bins, ...])`

`preprocessing.KernelCenterer()`

`preprocessing.LabelBinarizer([neg_label, ...])`

`preprocessing.LabelEncoder`

`preprocessing.MultiLabelBinarizer([classes, ...])`

`preprocessing.MaxAbsScaler([copy])`

`preprocessing.MinMaxScaler([feature_range, copy])`

`preprocessing.Normalizer([norm, copy])`

`preprocessing.OneHotEncoder([categories, ...])`

`preprocessing.OrdinalEncoder([categories, dtype])`

`preprocessing.PolynomialFeatures([degree, ...])`

`preprocessing.PowerTransformer([method, ...])`

`preprocessing.QuantileTransformer([...])`

`preprocessing.RobustScaler([with_centering, ...])`

`preprocessing.StandardScaler([copy, ...])`

Encodage

Normalisation

Création de polynômes

Transformation non linéaire

PREPROCESSING

`preprocessing.Binarizer([threshold, copy])`

`preprocessing.FunctionTransformer([func, ...])`

`preprocessing.KBinsDiscretizer([n_bins, ...])`

`preprocessing.KernelCenterer()`

`preprocessing.LabelBinarizer([neg_label, ...])`

`preprocessing.LabelEncoder`

`preprocessing.MultiLabelBinarizer([classes, ...])`

`preprocessing.MaxAbsScaler([copy])`

`preprocessing.MinMaxScaler([feature_range, copy])`

`preprocessing.Normalizer([norm, copy])`

`preprocessing.OneHotEncoder([categories, ...])`

`preprocessing.OrdinalEncoder([categories, dtype])`

`preprocessing.PolynomialFeatures([degree, ...])`

`preprocessing.PowerTransformer([method, ...])`

`preprocessing.QuantileTransformer([...])`

`preprocessing.RobustScaler([with_centering, ...])`

`preprocessing.StandardScaler([copy, ...])`

Encodage

Normalisation

Création de polynômes

Transformation non linéaire

Discrétisation

PREPROCESSING

`preprocessing.Binarizer([threshold, copy])`

`preprocessing.FunctionTransformer([func, ...])`

`preprocessing.KBinsDiscretizer([n_bins, ...])`

`preprocessing.KernelCenterer()`

`preprocessing.LabelBinarizer([neg_label, ...])`

`preprocessing.LabelEncoder`

`preprocessing.MultiLabelBinarizer([classes, ...])`

`preprocessing.MaxAbsScaler([copy])`

`preprocessing.MinMaxScaler([feature_range, copy])`

`preprocessing.Normalizer([norm, copy])`

`preprocessing.OneHotEncoder([categories, ...])`

`preprocessing.OrdinalEncoder([categories, dtype])`

`preprocessing.PolynomialFeatures([degree, ...])`

`preprocessing.PowerTransformer([method, ...])`

`preprocessing.QuantileTransformer([...])`

`preprocessing.RobustScaler([with_centering, ...])`

`preprocessing.StandardScaler([copy, ...])`

Encodage

Normalisation

Création de polynômes

Transformation non linéaire

Discrétisation

Personnalisation

PREPROCESSING

```
preprocessing.Binarizer([threshold, copy])
preprocessing.FunctionTransformer([func, ...])
preprocessing.KBinsDiscretizer([n_bins, ...])
preprocessing.KernelCenterer()
preprocessing.LabelBinarizer([neg_label, ...])
preprocessing.LabelEncoder
preprocessing.MultiLabelBinarizer([classes, ...])
preprocessing.MaxAbsScaler([copy])
preprocessing.MinMaxScaler([feature_range, copy])
preprocessing.Normalizer([norm, copy])
preprocessing.OneHotEncoder([categories, ...])
preprocessing.OrdinalEncoder([categories, dtype])
preprocessing.PolynomialFeatures([degree, ...])
preprocessing.PowerTransformer([method, ...])
preprocessing.QuantileTransformer([...])
preprocessing.RobustScaler([with_centering, ...])
preprocessing.StandardScaler([copy, ...])
```

Encodage

Normalisation

Création de polynômes

Si vous débutez

PREPROCESSING TRANSFORMERS

`preprocessing.Binarizer([threshold, copy])`

`preprocessing.FunctionTransformer([func, ...])`

`preprocessing.KBinsDiscretizer([n_bins, ...])`

`preprocessing.KernelCenterer()`

`preprocessing.LabelBinarizer([neg_label, ...])`

`preprocessing.LabelEncoder`

`preprocessing.MultiLabelBinarizer([classes, ...])`

`preprocessing.MaxAbsScaler([copy])`

`preprocessing.MinMaxScaler([feature_range, copy])`

`preprocessing.Normalizer([norm, copy])`

`preprocessing.OneHotEncoder([categories, ...])`

`preprocessing.OrdinalEncoder([categories, dtype])`

`preprocessing.PolynomialFeatures([degree, ...])`

`preprocessing.PowerTransformer([method, ...])`

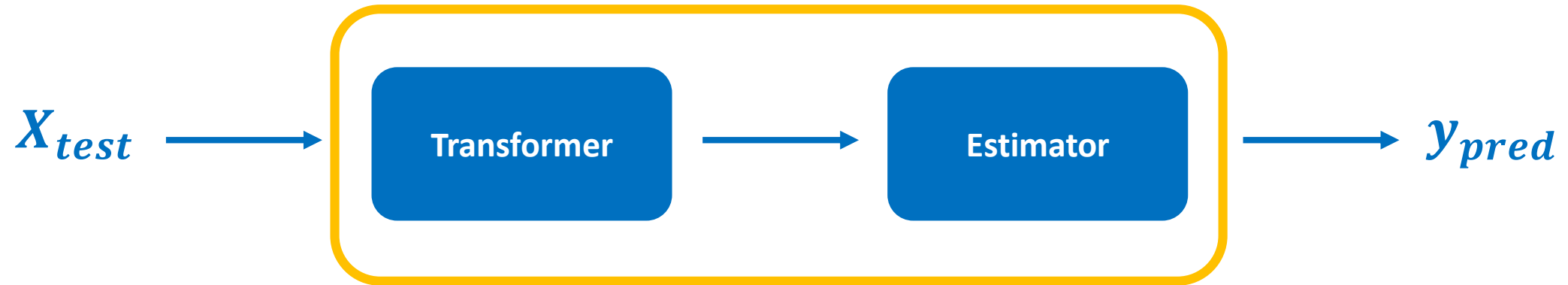
`preprocessing.QuantileTransformer([...])`

`preprocessing.RobustScaler([with_centering, ...])`

`preprocessing.StandardScaler([copy, ...])`

TRANSFORMER ET ESTIMATOR

PIPELINE



TRANSFORMER ESTIMATOR

X_{train}



Transformer

`.fit_transform(X_{train})`

Estimator

`.fit($X_{traintransformed}$, y_{train})`

X_{test}



Transformer

`.transform(X_{test})`

Estimator

`.predict($X_{testtransformed}$)`

TRANSFORMER ESTIMATOR

X_{train}



Transformer

`.fit_transform(X_{train})`

Estimator

`.fit($X_{traintransformed}$, y_{train})`

X_{test}



Transformer

`.transform(X_{test})`

Estimator

`.predict($X_{testtransformed}$)`

TRANSFORMER ESTIMATOR

X_{train}



Transformer

`.fit_transform(X_{train})`

Estimator

`.fit($X_{traintransformed}$, y_{train})`

X_{test}



Transformer

`.transform(X_{test})`

Estimator

`.predict($X_{testtransformed}$)`

TRANSFORMER ESTIMATOR

X_{train}



Transformer

`.fit_transform(X_{train})`

Estimator

`.fit($X_{traintransformed}$, y_{train})`

X_{test}



Transformer

`.transform(X_{test})`

Estimator

`.predict($X_{testtransformed}$)`

TRANSFORMER ESTIMATOR

X_{train}



Transformer

`.fit_transform(X_{train})`

Estimator

`.fit($X_{traintransformed}$, y_{train})`

X_{test}



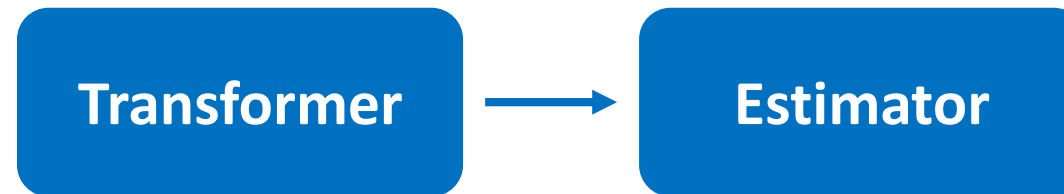
Transformer

`.transform(X_{test})`

Estimator

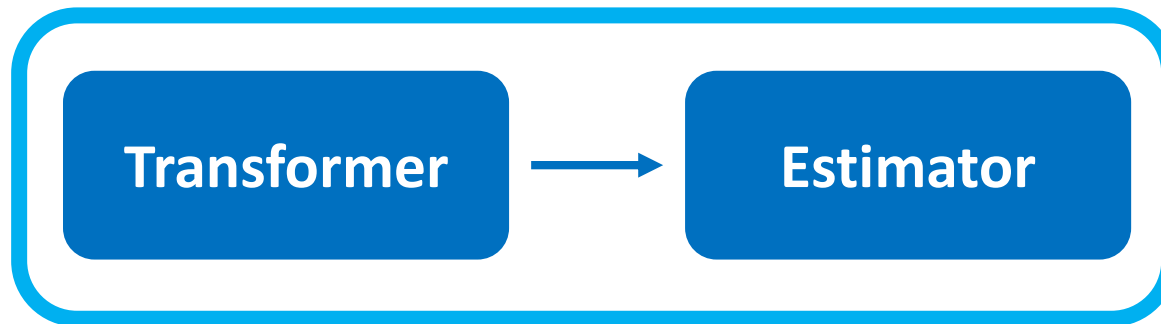
`.predict($X_{testtransformed}$)`

PIPELINE ET ESTIMATEUR COMPOSITE



PIPELINE ET ESTIMATEUR COMPOSITE

Pipeline



PIPELINE ET ESTIMATEUR COMPOSITE

Pipeline

Composite Estimator

PIPELINE ET ESTIMATEUR COMPOSITE

Pipeline

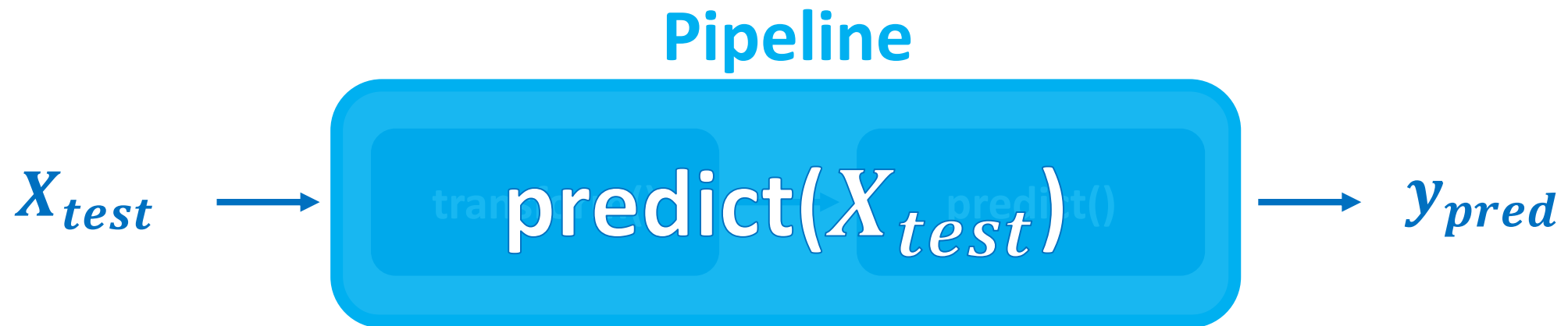


PIPELINE ET ESTIMATEUR COMPOSITE

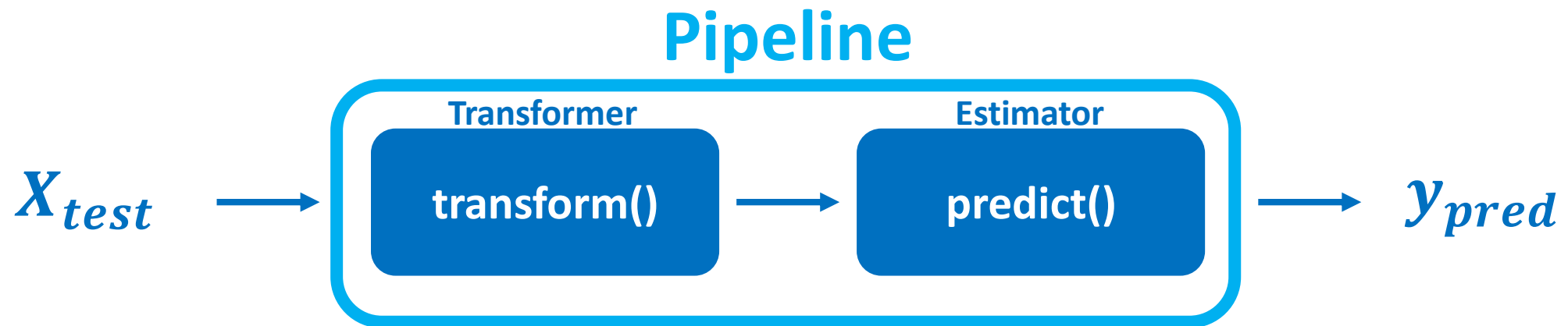
Pipeline



PIPELINE ET ESTIMATEUR COMPOSITE



PIPELINE ET ESTIMATEUR COMPOSITE



PIPELINE

Pour créer une **Pipeline**, 2 options :

- Classe Pipeline
- Fonction **make_pipeline**


```
model = make_pipeline(StandardScaler(),  
                      SGDClassifier())
```

```
model.fit( $X_{train}$ ,  $y_{train}$ )  
model.score( $X_{test}$ ,  $y_{test}$ )  
model.predict( $X_{test}$ )
```

PIPELINE GRIDSEARCHCV

On peut utiliser **GridSearchCV** pour trouver les meilleurs paramètres de la pipeline.

```
grid = GridSearchCV(pipeline, params, cv)
```



```
params = {  
    <Composant>__<paramètre> : [...]  
}
```

```
grid.fit( $X_{train}$ ,  $y_{train}$ )
```

```
grid.best_estimator
```

COLUMN TRANSFORMER

ColumnTransformer permet d'appliquer vos transformers sur les **colonnes** que vous sélectionnez.

```
make_column_transformer ((transformer,  
['colonne 1', 'colonne 2', ...] ))
```

COLUMN TRANSFORMER

ColumnTransformer permet d'appliquer vos transformers sur les **colonnes** que vous sélectionnez.

```
make_column_transformer ((transformer,  
['colonne 1', 'colonne 2', ...] ))
```

Application Classique :

Séparez les **catégories** et variables **numériques en 2 listes**, puis utilisez ColumnTransformer pour traiter chaque liste de variables séparément.

PIPELINE AVANCEE

Numerical_pipeline

SimpleImputer

StandardScaler

Categorical_pipeline

SimpleImputer

OneHotEncoder

Column_transformer

Numerical_pipeline

SimpleImputer

StandardScaler



Categorical_pipeline

SimpleImputer

OneHotEncoder

Pipeline

Column_transformer

Numerical_pipeline

SimpleImputer

StandardScaler



Categorical_pipeline

SimpleImputer

OneHotEncoder



SGDClassifier

COLUMN SELECTOR

Depuis la version 0.22 de sklearn, il existe **make_column_selector** qui permet de sélectionner certains **types** de colonnes.

3 arguments au choix :

dtype_include

dtype_exclude

pattern (avec des *Regular Expressions*)

Exemple de types :

np.number → variables numériques

object → variables catégorielles

FEATURE UNION

`make_union()` permet de créer des pipelines **parallèles**. Les résultats sont **concaténés** à la sortie du transformer.

