

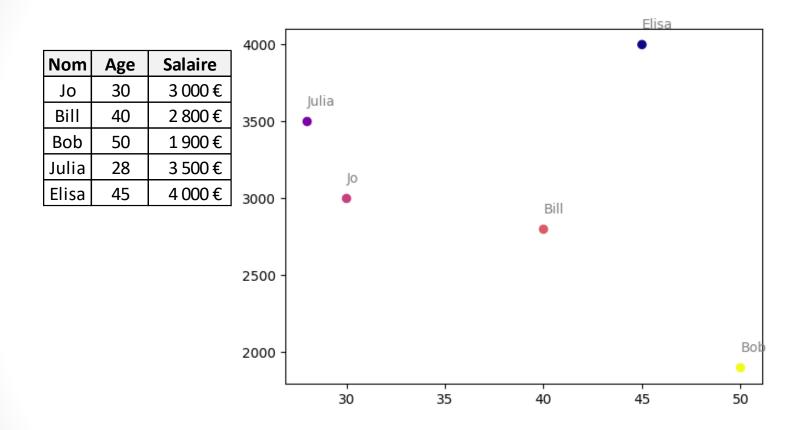
# Big Data & Traitement par l'IA II



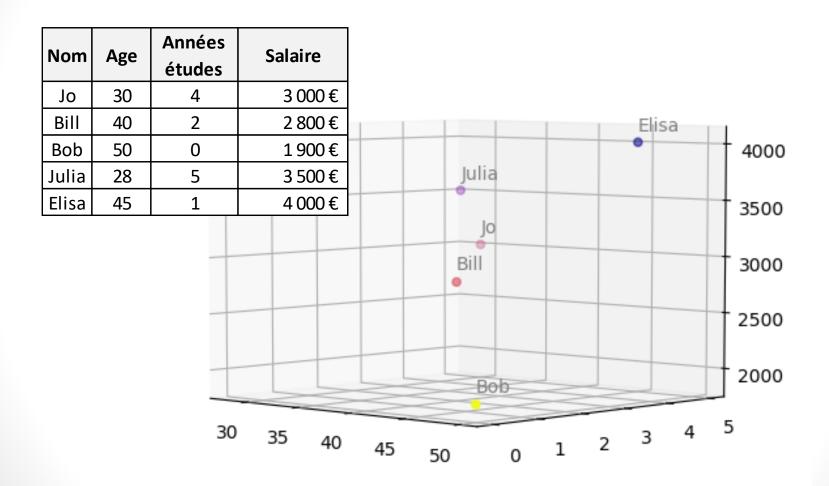
## Réduction de dimensionnalité PCA, tSNE

Clustering avec KMeans

#### Affichage des données 2D

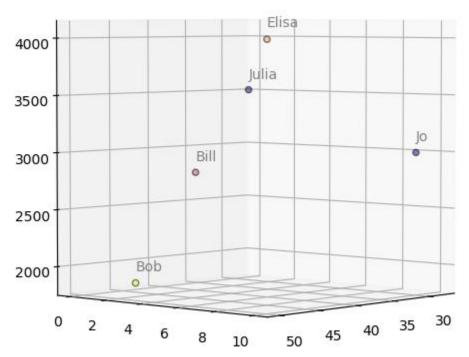


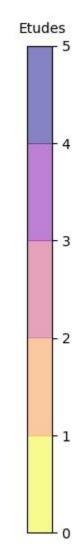
#### Affichage des données 3D



#### Gagner une dimension en utilisant un code couleur

Nom	Age	Années études	Personnes gérées	Salaire
Jo	30	4	10	3 000 €
Bill	40	2	2	2 800 €
Bob	50	0	3	1900€
Julia	28	5	0	3 500 €
Elisa	45	1	10	4 000 €





#### Et Après?

Nom	Age	Années études	Personnes gérées	Ancienneté	Salaire
Jo	30	4	10	2	3 000 €
Bill	40	2	2	5	2 800 €
Bob	50	0	3	8	1900€
Julia	28	5	0	4	3 500 €
Elisa	45	1	10	10	4 000 €

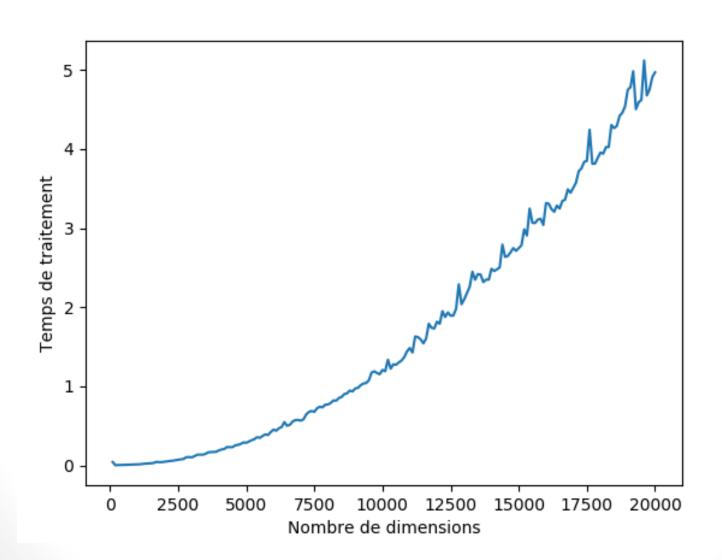


Solution : Construire des nouvelles variables qui permettent de réduire le nombre de dimensions tout en perdant un minimum d'information.

Avantage : On peut visualiser en 2D ou 3D les points proches, on peut aussi réduire le nombre de dimensions, et donc le temps de calcul, pour l'apprentissage.

Inconvénient : Les variables n'ont plus vraiment de sens.

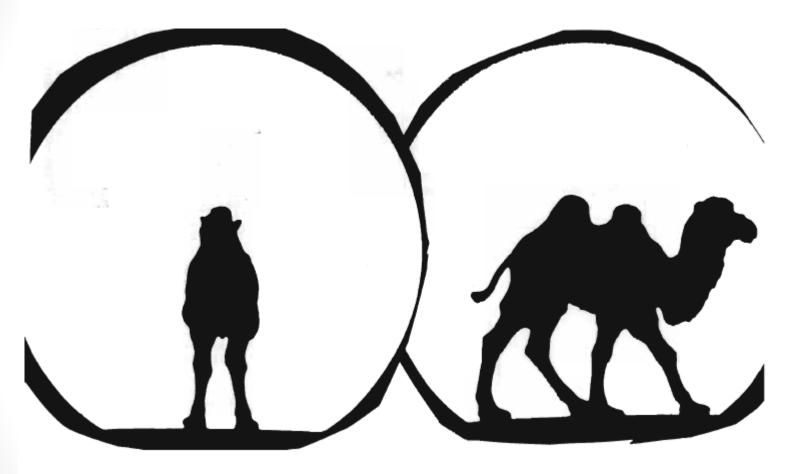
#### La « Malédiction des dimensions »





# Analyse en composantes principales (ou PCA)

#### Trouver la projection qui explique le mieux



Créé en 1901

## Le Dataset Iris







**Iris Versicolor** 

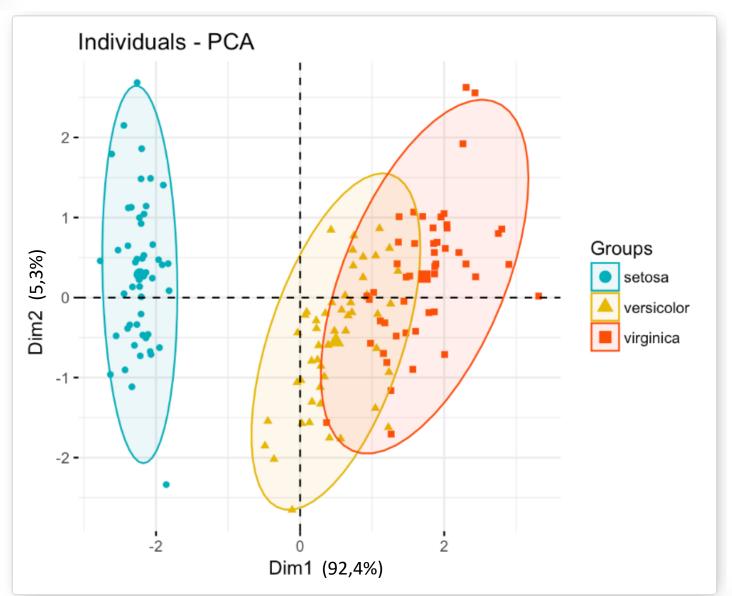
**Iris Setosa** 

Iris Virginica

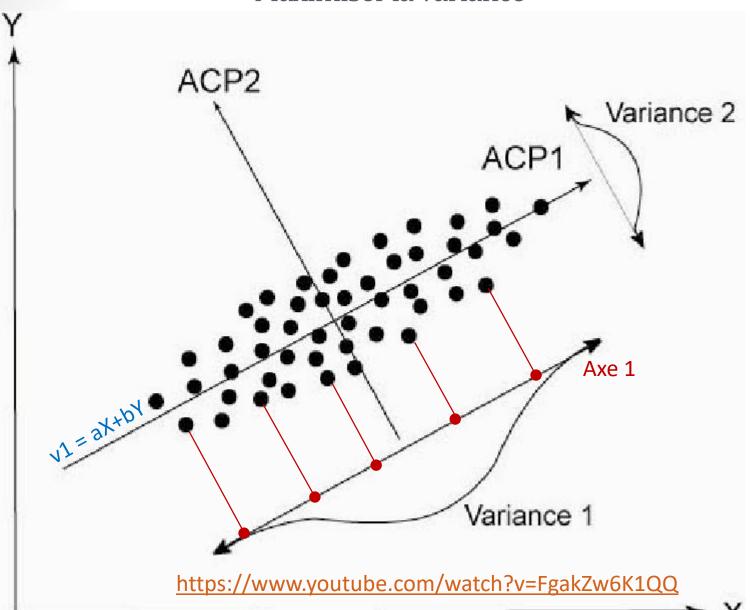
	sepal length	sepal width	petal length	petal width	target
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa

- 50 exemples pour chaque type d'Iris.
- Disponible dans Scikit\_Learn en utilisant la fonction load\_iris() du package sklearn.datasets.

## 4D vers 2D



#### Maximiser la variance

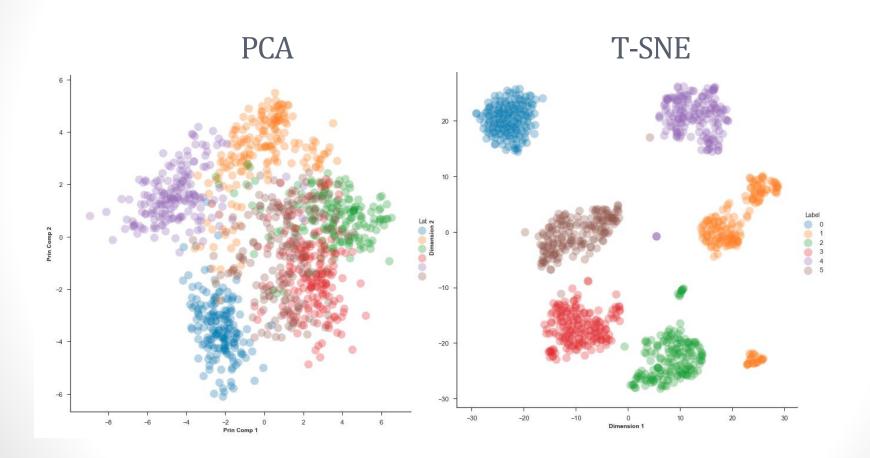




# T-SNE t-distributed stochastic neighbor embedding

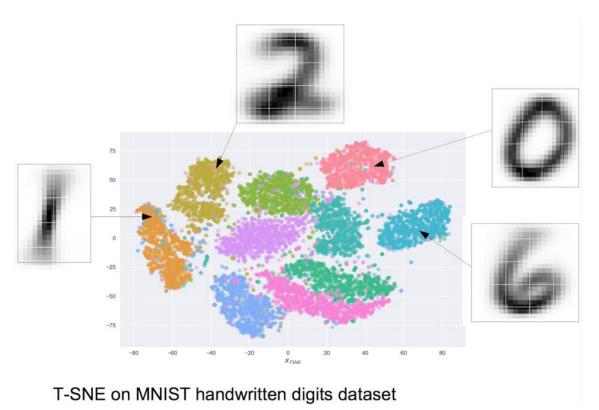
#### T-SNE

- Une méthode non-linéaire imaginée en 2008.
- Représentation d'un ensemble de points d'un espace à grande dimension dans un espace de deux ou trois dimensions.



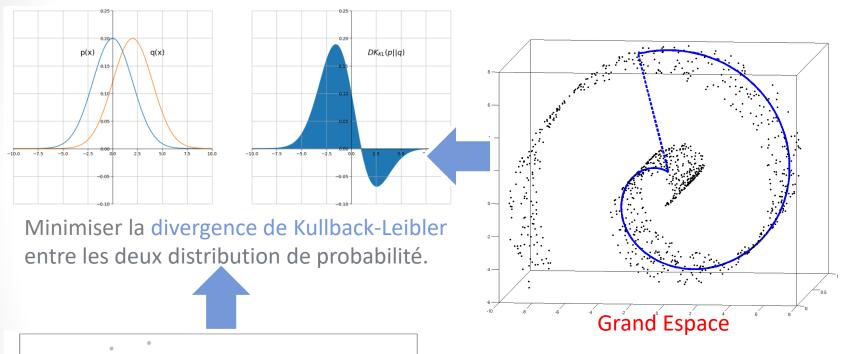
#### T-SNE

- Deux points qui sont proches (éloignés) dans l'espace d'origine devront être proches (éloignés) dans l'espace de faible dimension.
- ➤ Forte Complexité algorithmique o(n²)



#### Respect statistique des distances

Pour chaque couple de points, la distance est transformé en probabilité (similaire à la fonction softmax) ? Un variance est introduite la perplexité, qui va correspondre à la chute de la probabilité avec la distance.



Petit Espace

En utilisant la descente de gradient, avec un Learning Rate epsilon, on va faire converger les distributions des points du petit espace vers celle du grand espace.

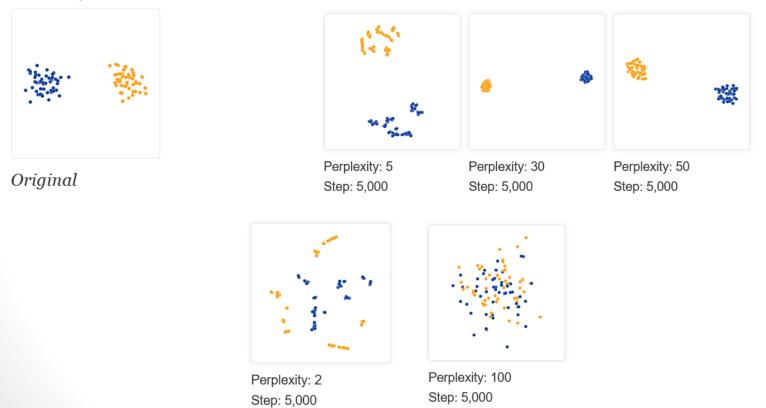
https://distill.pub/2016/misread-tsne/

#### Perplexité dans t-sne

Pour chaque point quelle est la probabilité pour un autre point d'être un voisin ? On utilise une loi normale dont l'écart type basé sur le nombre moyen de voisins perplexité

- => Une perplexité trop faible surestime les variations locales.
- => Une perplexité trop importante donne un résultat aléatoire (surtout si supérieur au nombre de points).

La valeur préconisée est entre 5 et 50.



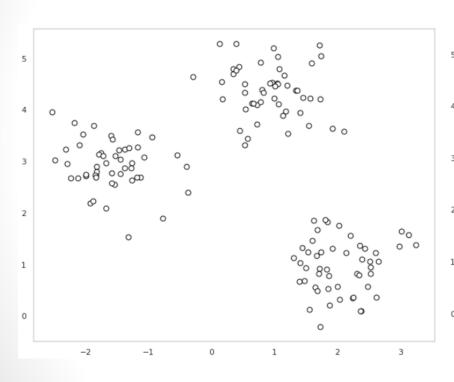


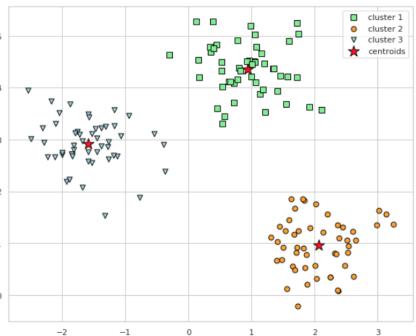
## Clustering

## Clustering

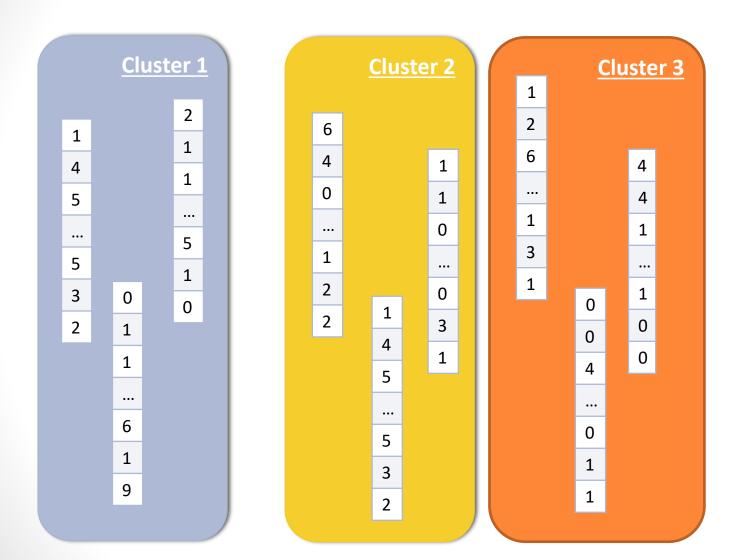
Regroupement les données (en N dimensions) en un nombre quelconque de clusters.

Peut souvent précéder une réduction de dimension.





# Regroupement de vecteurs



## **KMeans**

L'algorithme le plus populaire pour faire du clustering est le KMeans.

KMeans est un apprentissage non supervisé.

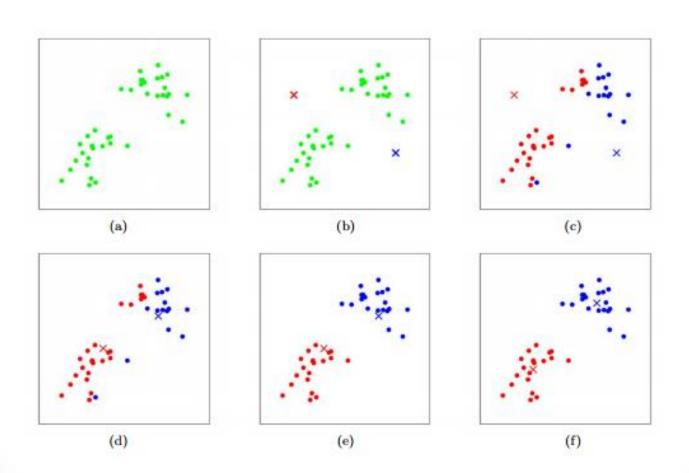
Pour faire K clusters, on prend K points (centroïds) au hasard et on itère sur 2 étapes :

- 1. On associe les points les plus proches pour chaque centroïd
- 2. On déplace les centroïds au centre des points qui lui étaient associés.
- ... et on continue jusqu'à ce que les centroïdes ne bougent plus.

Si à un moment, un centroïd se retrouve sans point, alors soit on l'élimine (le + courant) soit on le déplace aléatoirement.

# Kmeans/Illustration

2 clusters, 2 itérations





## Codage avec SKLearn

### sklearn.cluster.KMeans

class <u>sklearn.cluster.KMeans</u>(n\_clusters=8, init='k-means++', n\_init=10, max\_iter=300, tol=0.0001, verbose=0, random\_state=None, copy\_x=True, algorithm='auto')

<pre>fit(self, X[, y, sample_weight])</pre>	Compute k-means clustering.
<pre>fit_predict(self, X[, y, sample_weight])</pre>	Compute cluster centers and predict cluster index for each sample.
<pre>fit transform(self, X[, y, sample_weight])</pre>	Compute clustering and transform X to cluster-distance space.
<pre>get_params(self[, deep])</pre>	Get parameters for this estimator.
<pre>predict(self, X[, sample_weight])</pre>	Predict the closest cluster each sample in X belongs to.
<pre>score(self, X[, y, sample_weight])</pre>	Opposite of the value of X on the K-means objective.
<pre>set params(self, \*\*params)</pre>	Set the parameters of this estimator.
transform(self, X)	Transform X to a cluster-distance space.

## sklearn.decomposition.PCA

class <u>sklearn.decomposition.PCA</u>(n\_components=None, copy=True, whiten=False, svd\_solver='auto', tol=0.0, iterated\_power='auto', random\_state=None)

<pre>fit(self, X[, y])</pre>	Fit the model with X.
<pre>fit transform(self, X[, y])</pre>	Fit the model with X and apply the dimensionality reduction on X.
get covariance(self)	Compute data covariance with the generative model.
<pre>get_params(self[, deep])</pre>	Get parameters for this estimator.
get precision(self)	Compute data precision matrix with the generative model.
inverse transform(self V)	
inverse transform(self, X)	Transform data back to its original space.
score(self, X[, y])	Return the average log-likelihood of all samples.
	Return the average log-likelihood of all
score(self, X[, y])	Return the average log-likelihood of all samples.
<pre>score(self, X[, y]) score samples(self, X)</pre>	Return the average log-likelihood of all samples.  Return the log-likelihood of each sample.

### sklearn.manifold.TSNE

class <u>sklearn.manifold.TSNE</u>(n\_components=2, \*, perplexity=30.0, early\_exaggeration=12.0, learning\_rate='warn', n\_iter=1000, n\_iter\_without\_progress=300, min\_grad\_norm=1e-07, metric='euclidean', init='warn', verbose=0, random\_state=None, method='barnes\_hut', angle=0.5, n\_jobs=None, square\_distances='legacy')

<pre>fit(self, X[, y])</pre>	Fit X into an embedded space.
<pre>fit transform(self, X[, y])</pre>	Fit X into an embedded space and return that transformed output.
<pre>get_params(self[, deep])</pre>	Get parameters for this estimator.
<pre>set params(self, \*\*params)</pre>	Set the parameters of this estimator.