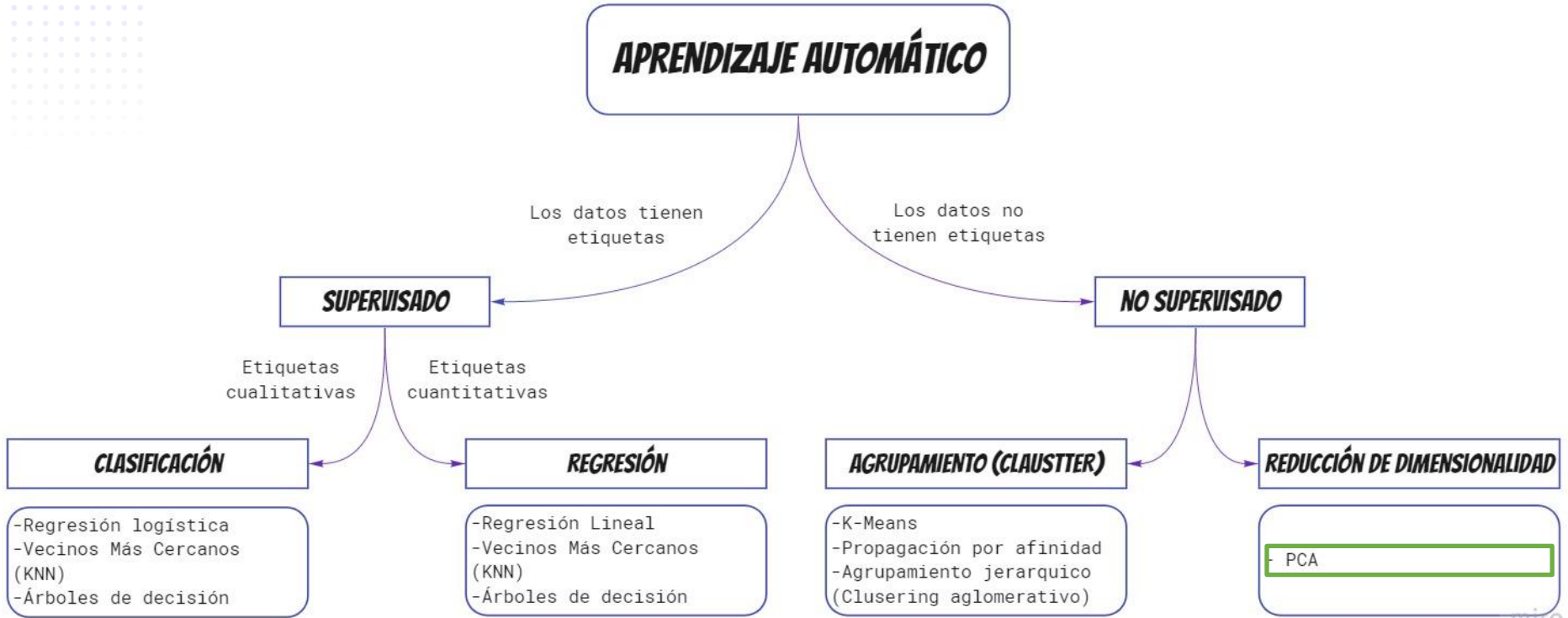


# Ciencia de Datos

- Módulo 4

*PCA*

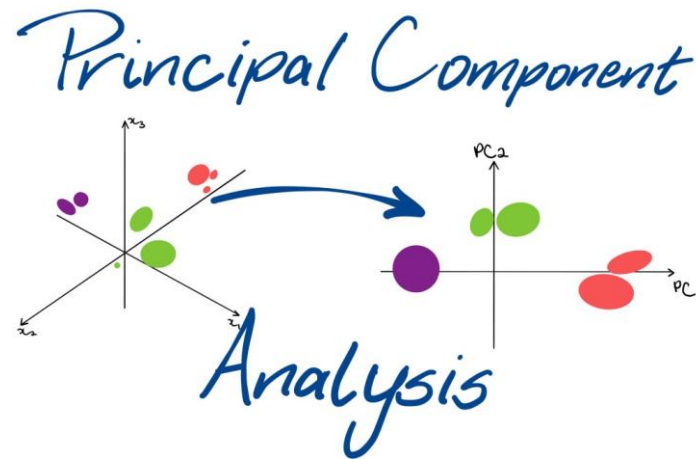




# PCA

Más conocido como PCA (por sus siglas en inglés Principal Components Analysis), es un algoritmo de aprendizaje no supervisado, comúnmente usado para lo que llamamos reducción de dimensionalidad.

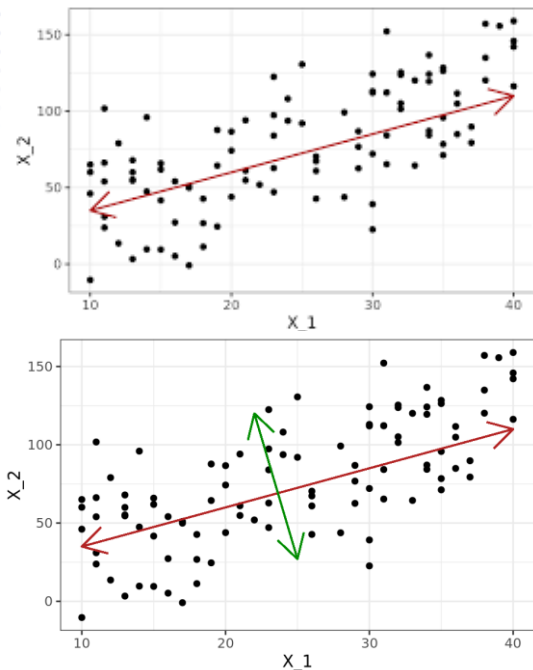
En muchos problemas de la vida real se tienen una gran cantidad de variables. En estos casos puede no ser conveniente incluir todas las variables, y decidir cuáles eliminar es un problema en si, este problema es el que trata de resolver PCA.



Fuente: <https://laptrinhx.com/understanding-principal-component-analysis-1459875749/>



# PCA



PCA busca eliminar variables afectando lo menos posible la variabilidad de los datos, esto es, mantener la mayor cantidad de casos existentes, para esto lo que se hace es aplicar **transformaciones** a las variables (descomponerlas), convirtiendo a las mismas en lo que llamamos **componentes**.

Cada componente aporta un grado de variabilidad al conjunto de datos, estas componentes están ordenadas de forma tal que la primera componente es aquella que mantiene la mayor variabilidad del conjunto, y la última la que menor variabilidad aporta, luego nos quedamos con las N primeras componentes.

Además estos componentes son independientes entre sí, siendo que puede suceder que algunas variables originales no lo sean (esto es muy importante en algunos modelos)

Fuente: [https://rpubs.com/Joaquin\\_AR/287787](https://rpubs.com/Joaquin_AR/287787)

# PCA



Fuente: <https://www.seoperu.com/los-mejores-procesadores-para-servidores/>

PCA se utiliza para:

- Comprimir los datos lo que es importante en términos de memoria y procesamiento. Además sirve para reducir el número de variables sin identificar cuales son las más explicativas, permite descartar información redundante o ruido lo que puede ser importante para la generalización
- Visualizar en 2 dimensiones de conjuntos de datos.

El mayor inconveniente con la utilización de PCA es que limita de un manera significativa la interpretabilidad de los datos



# PCA – Código



```
31 self.file = None
32 self.fingerprints = []
33 self.logspaces = True
34 self.debug = debug
35 self.logger = logging.getLogger(__name__)
36
37 if path:
38     self.file = os.path.join(path, 'fingerprints.log')
39     self.file.write('')
40     self.fingerprints.append('')
41
42 @classmethod
43 def from_settings(cls, settings):
44     debug = settings.get('debug', False)
45     return cls([os.path.join(settings['path'], 'fingerprints.log')], debug)
46
47 def request_seen(self, request):
48     fp = self.request_fingerprint(request)
49     if fp in self.fingerprints:
50         return True
51     self.fingerprints.add(fp)
52     if self.file:
53         self.file.write(fp + os.linesep)
54
55 def request_fingerprint(self, request):
56     return request_fingerprint(request)
```

En Python con la librería Scikit-Learn se utiliza:

- Método `.fit` para entrenar el modelo
- Método `.transform` para transformar los datos.

Es posible utilizar todo en una misma línea utilizando el método `.fit_transform`

La cantidad de componentes las podemos seleccionar de forma directa definiéndolo como **hiperparámetros** (decir cuántas componentes se quieren trabajar) o seleccionar la cantidad de componentes a mantener tal que se mantenga cierto porcentaje de variabilidad del conjunto de datos.

El máximo de componentes es el mínimo entre la cantidad de registros o la cantidad de variables.

Solamente se utiliza para **variables numéricas**. Para variables categóricas se utilizan otras técnicas (análisis de correspondencias y categorical PCA)

Este es un modelo NO supervisado por lo tanto no se desea predecir ninguna etiqueta. Solamente se **transforman** los datos que se están utilizando.

# PCA

---

## Vídeo

- <https://www.youtube.com/watch?v=UVHneBUBW0>

## Documentación Scikit-Learn

- <https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.PCA.html>

## Adicional

- [https://www.youtube.com/watch?v=HMOI\\_lkzW08](https://www.youtube.com/watch?v=HMOI_lkzW08)
- <https://empresas.blogthinkbig.com/python-para-todos-que-es-el-pca/>
- [https://www.jacobsoft.com.mx/es\\_mx/analisis-del-componente-principal/](https://www.jacobsoft.com.mx/es_mx/analisis-del-componente-principal/)

