

Modelos de datos en Python

Juan Fernando Pérez

Departamento de Ingeniería Industrial

Universidad de los Andes

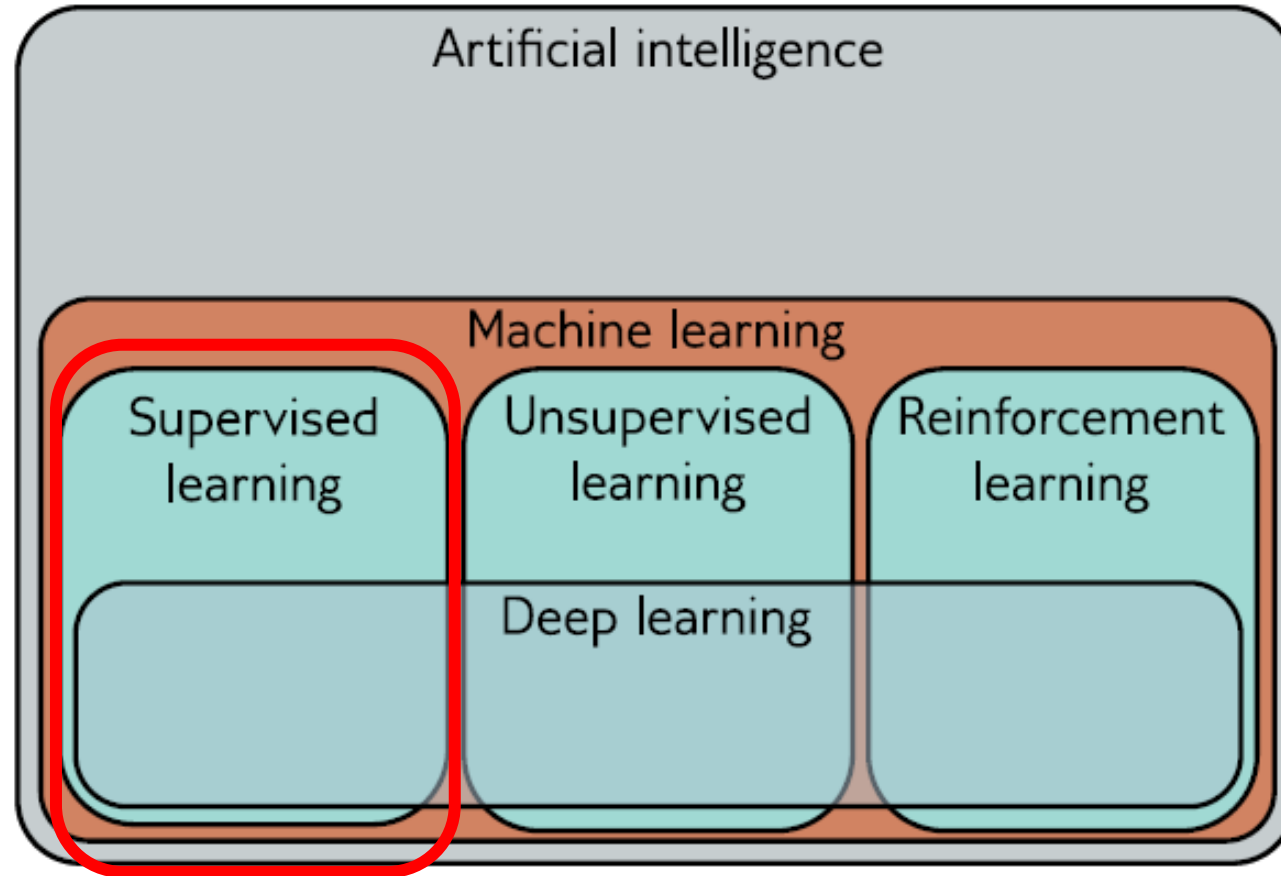
Agosto de 2024

jf.perez33@uniandes.edu.co

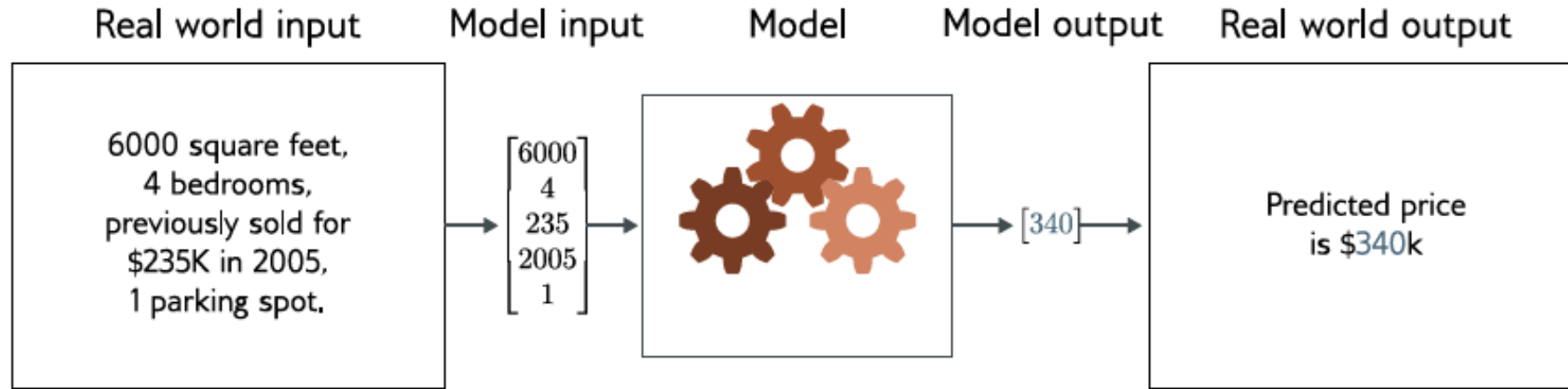
juanfperez.com

Aprendizaje supervisado

Problemas de Aprendizaje de máquina



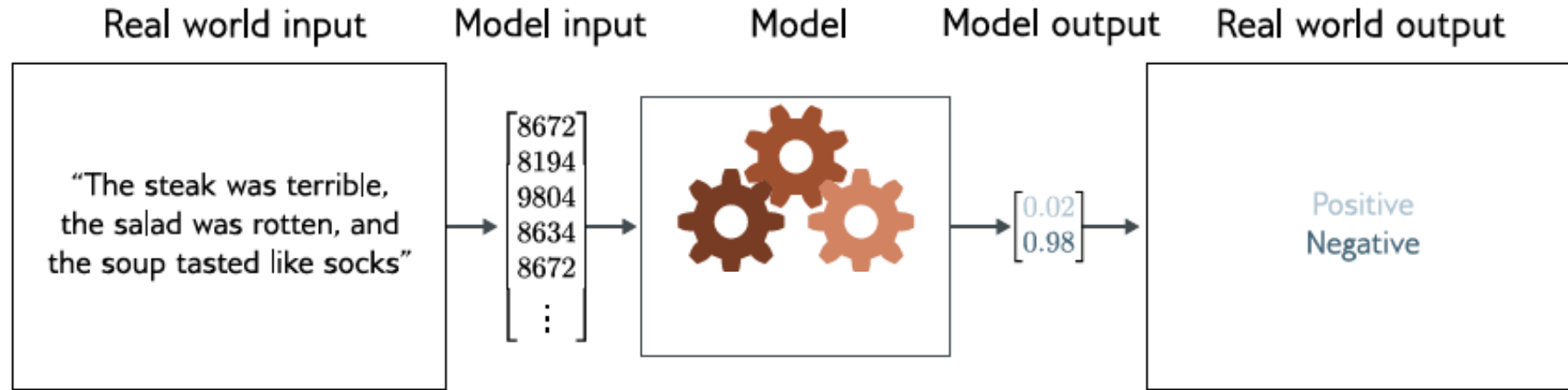
Aprendizaje supervisado



Predicción de una **salida continua**: valor de un inmueble

Regresión

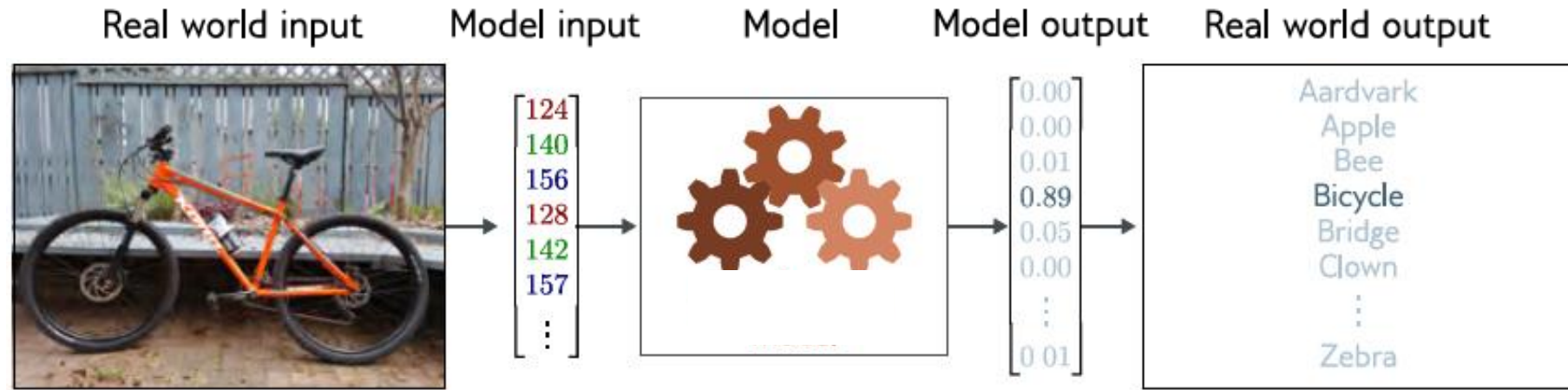
Aprendizaje supervisado



Predicción de una **salida discreta** (categoría, **binaria**)

Clasificación

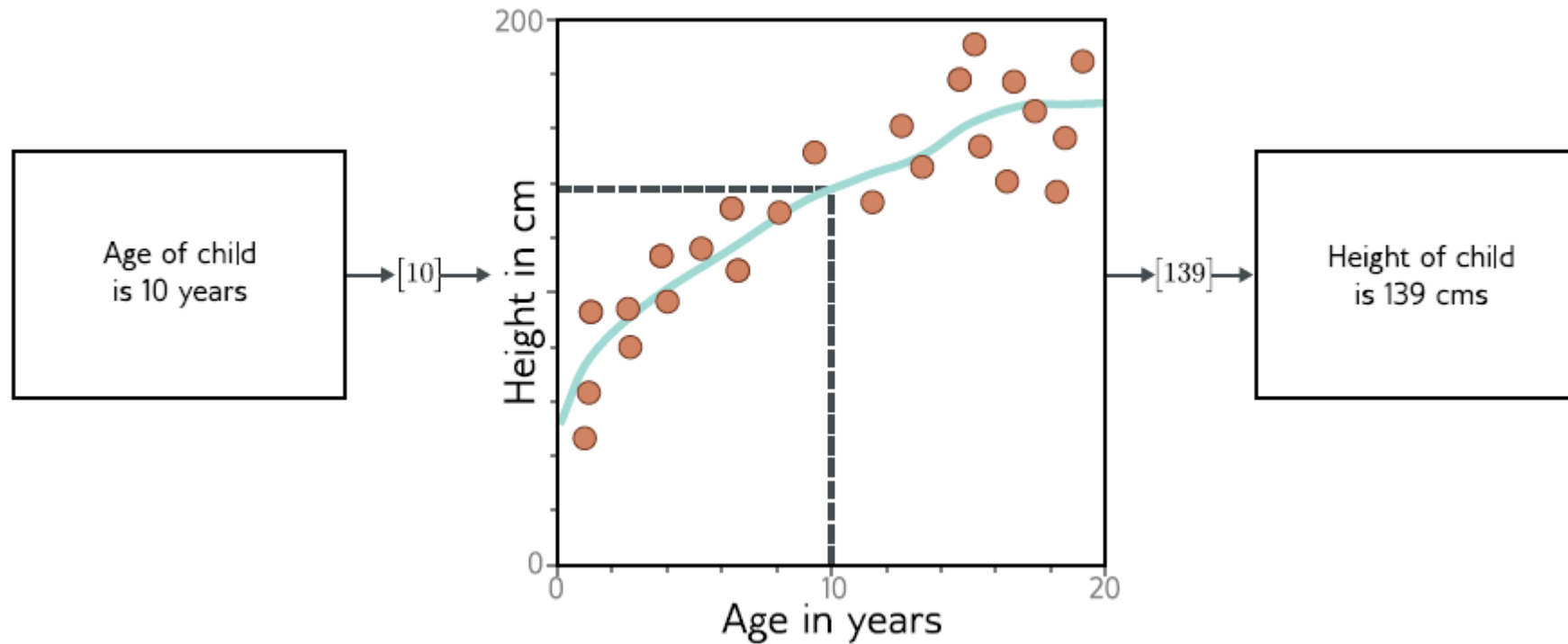
Aprendizaje supervisado



Predicción de una **salida discreta** (categoría, **multiclase**)

Clasificación

Aprendizaje supervisado – Ej. Regresión



Modelo de aprendizaje de máquina: familia de ecuaciones que permiten obtener una **salida** a partir de una **entrada**

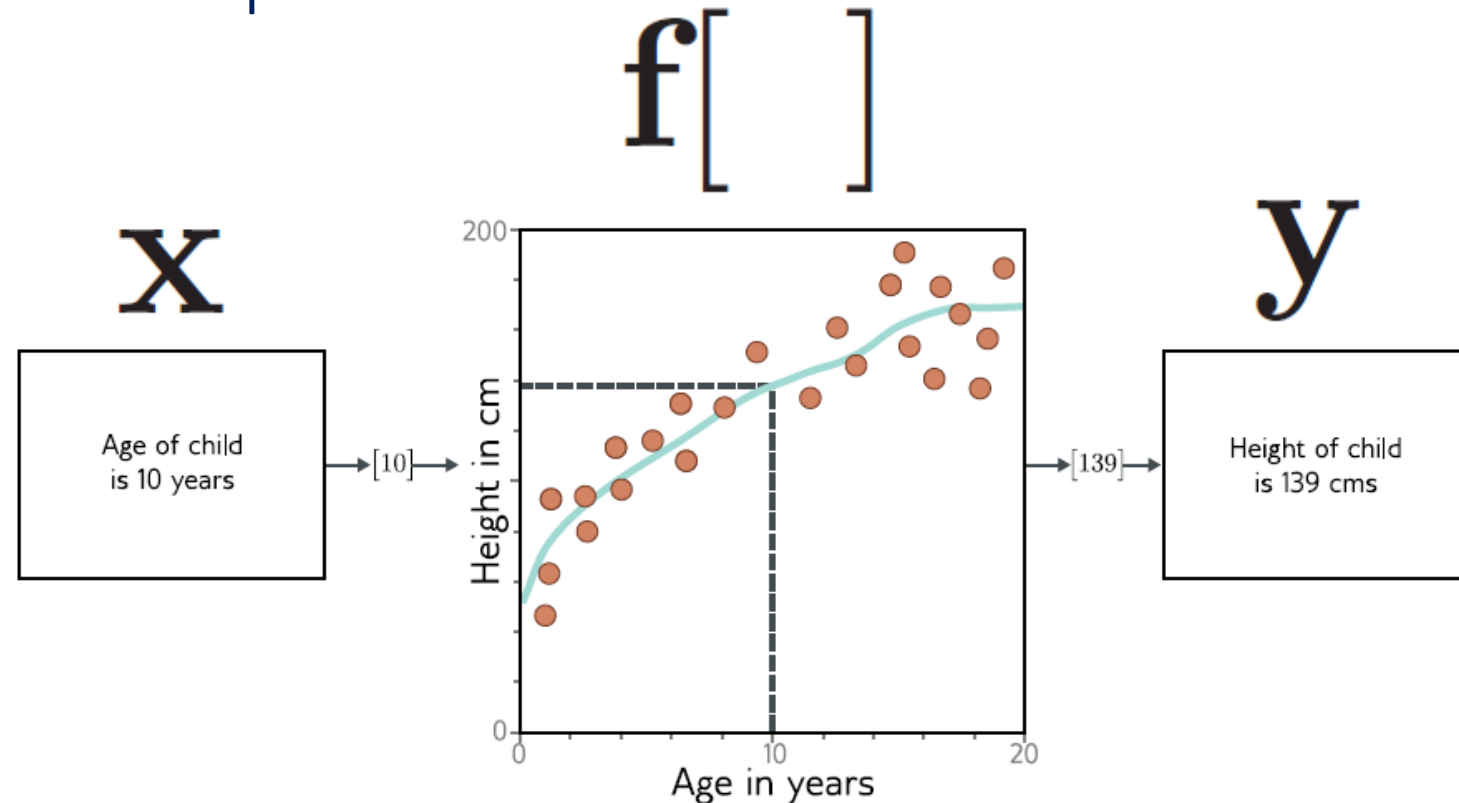
Aprendizaje supervisado

$$y = f[x]$$

f : Modelo de aprendizaje de máquina

x : entrada

y : salida



Aprendizaje supervisado

$$\mathbf{y} = \mathbf{f}[\mathbf{x}, \phi]$$

\mathbf{f} : Modelo de aprendizaje de máquina

\mathbf{x} : entrada

\mathbf{y} : salida

ϕ : parámetros que definen el modelo

Entrenamiento

Determinar el **mejor valor** de los **parámetros** considerando un **conjunto de datos** (de entrenamiento)

Aprendizaje supervisado - Entrenamiento

$$\mathbf{y} = \mathbf{f}[\mathbf{x}, \phi]$$

Datos de entrenamiento: $\{\mathbf{x}_i, \mathbf{y}_i\}$

Observaciones indexadas por i

\mathbf{x} : entrada observada

\mathbf{y} : salida observada (etiqueta)

Supervisado

Se requieren **ejemplos** pasados de **parejas entrada-salida** para entrenar el modelo

Aprendizaje supervisado - Entrenamiento

¿Cómo cuantificar qué tan bueno es un modelo $\mathbf{f}[\mathbf{x}, \phi]$?

Función de pérdida (loss): $\mathbf{L}[\phi]$

Valores **bajos** si con los parámetros el modelo captura **bien** la relación entre las entradas y salidas de los datos de entrenamiento

Aprendizaje supervisado - Entrenamiento

¿Cómo cuantificar qué tan bueno es un modelo $\mathbf{f}[\mathbf{x}, \phi]$?

Función de pérdida (loss): $\mathbf{L}[\phi]$

Mejor valor de los parámetros: $\hat{\phi} = \operatorname{argmin}_{\phi} [\mathbf{L}[\phi]]$

* $\mathbf{L}[\{\mathbf{x}_i, \mathbf{y}_i\}, \phi]$: depende de los datos de entrenamiento

Aprendizaje supervisado: regresión

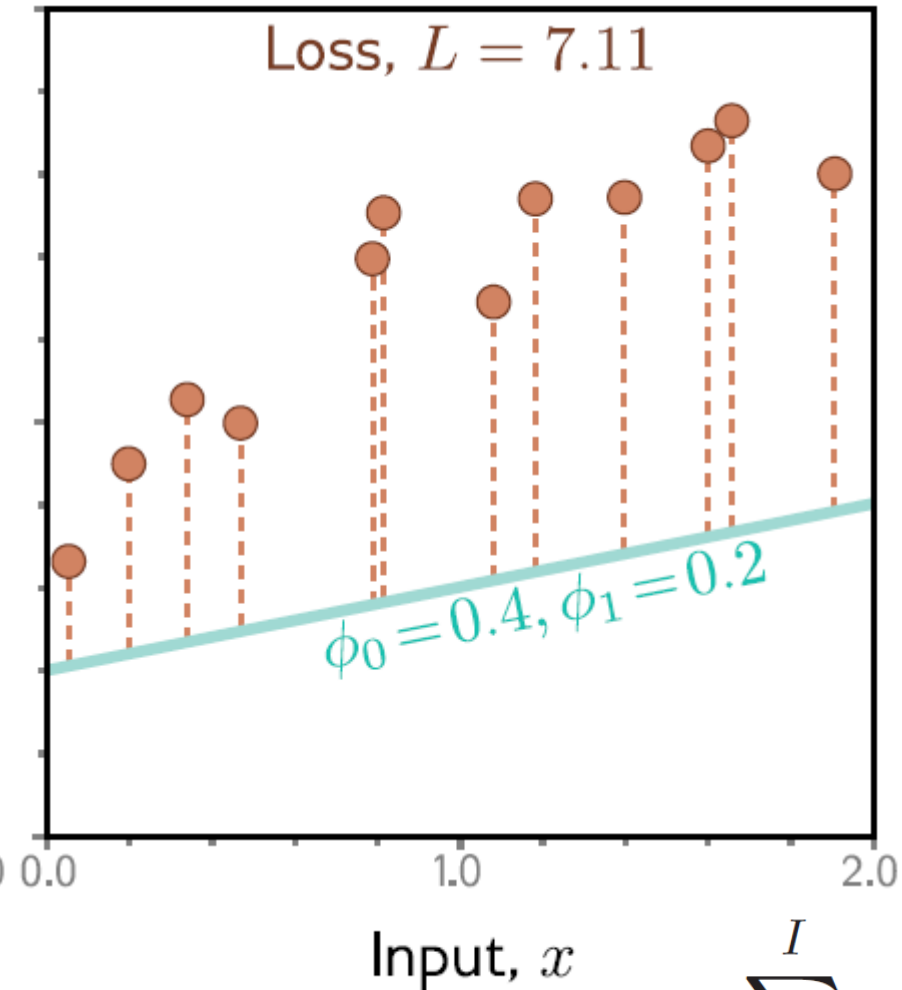
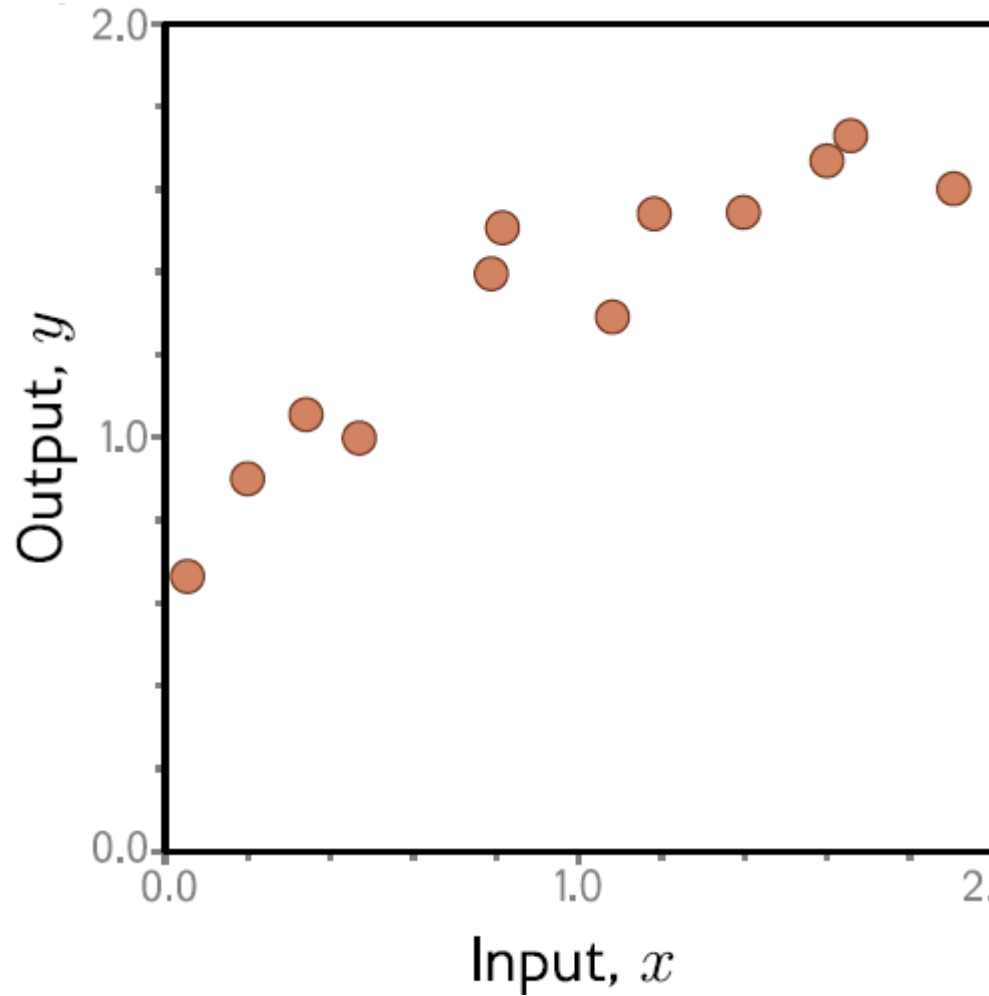
Entrenamiento: ejemplo Regresión Simple

Modelo: $y = f[x, \phi] = \phi_0 + \phi_1 x$

Función de pérdida:

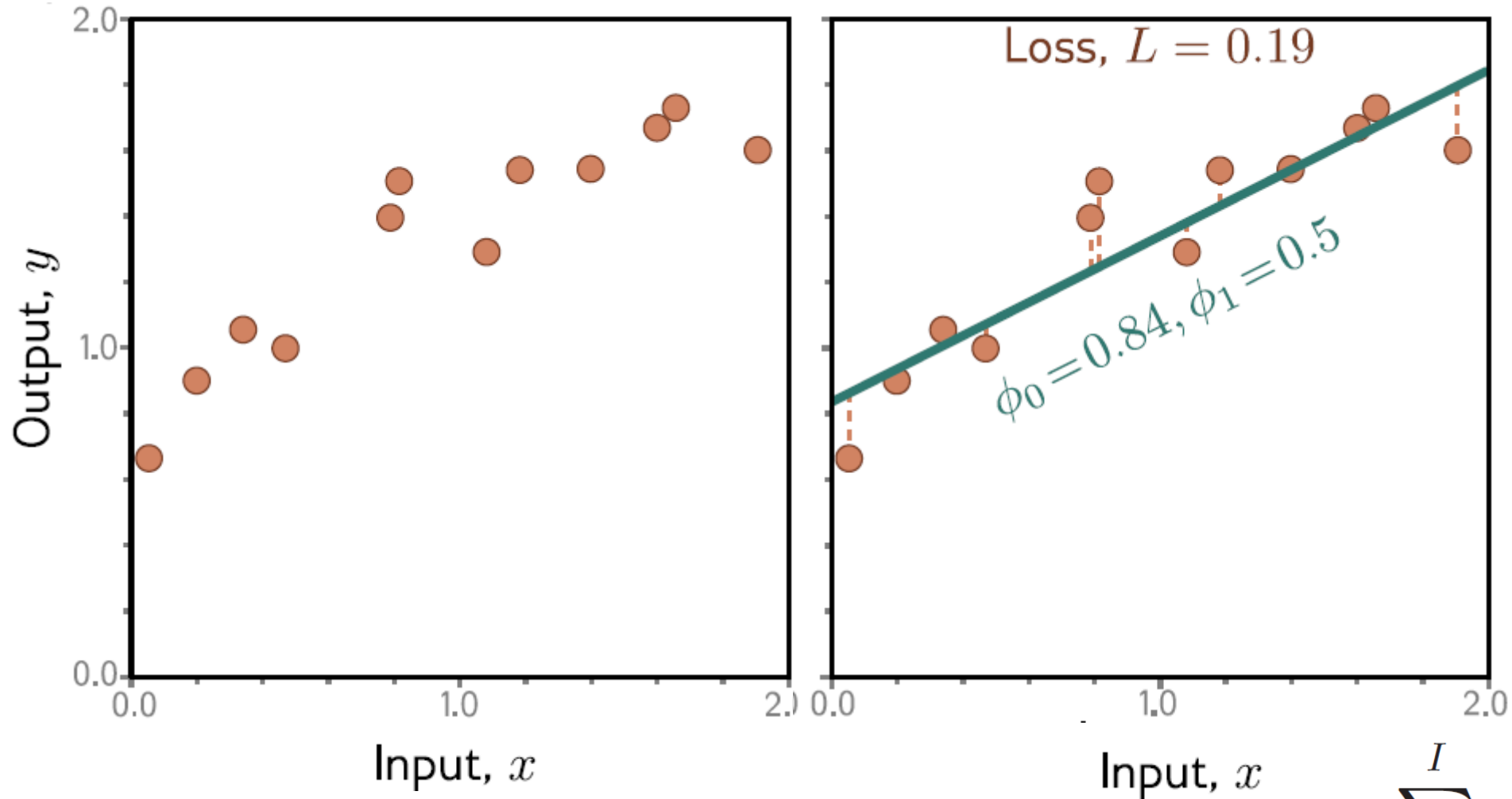
$$\begin{aligned} L[\phi] &= \sum_{i=1}^I (f[x_i, \phi] - y_i)^2 \\ &= \sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2 \end{aligned}$$

Entrenamiento: ejemplo Regresión Simple



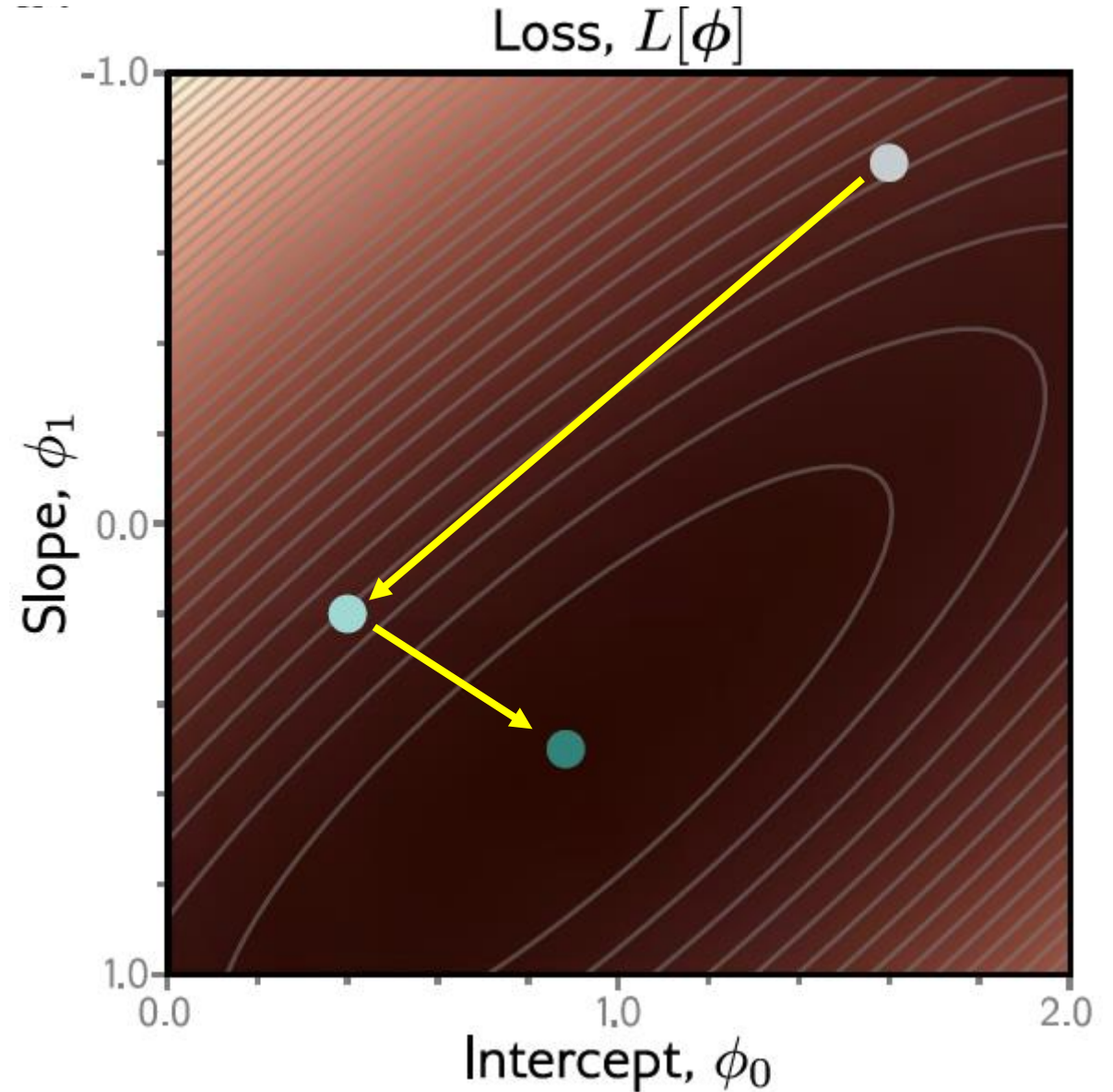
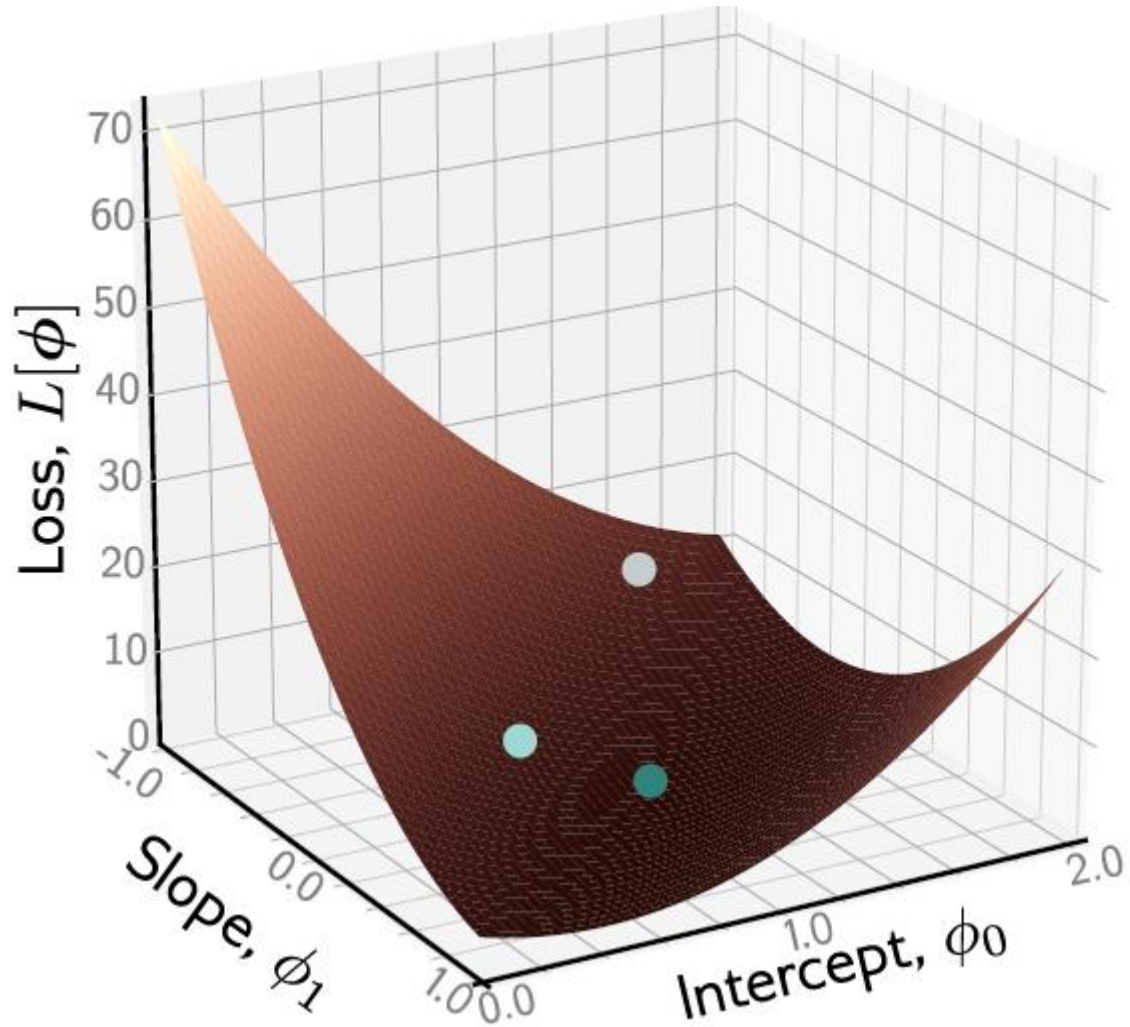
$$\sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2$$

Entrenamiento: ejemplo Regresión Simple



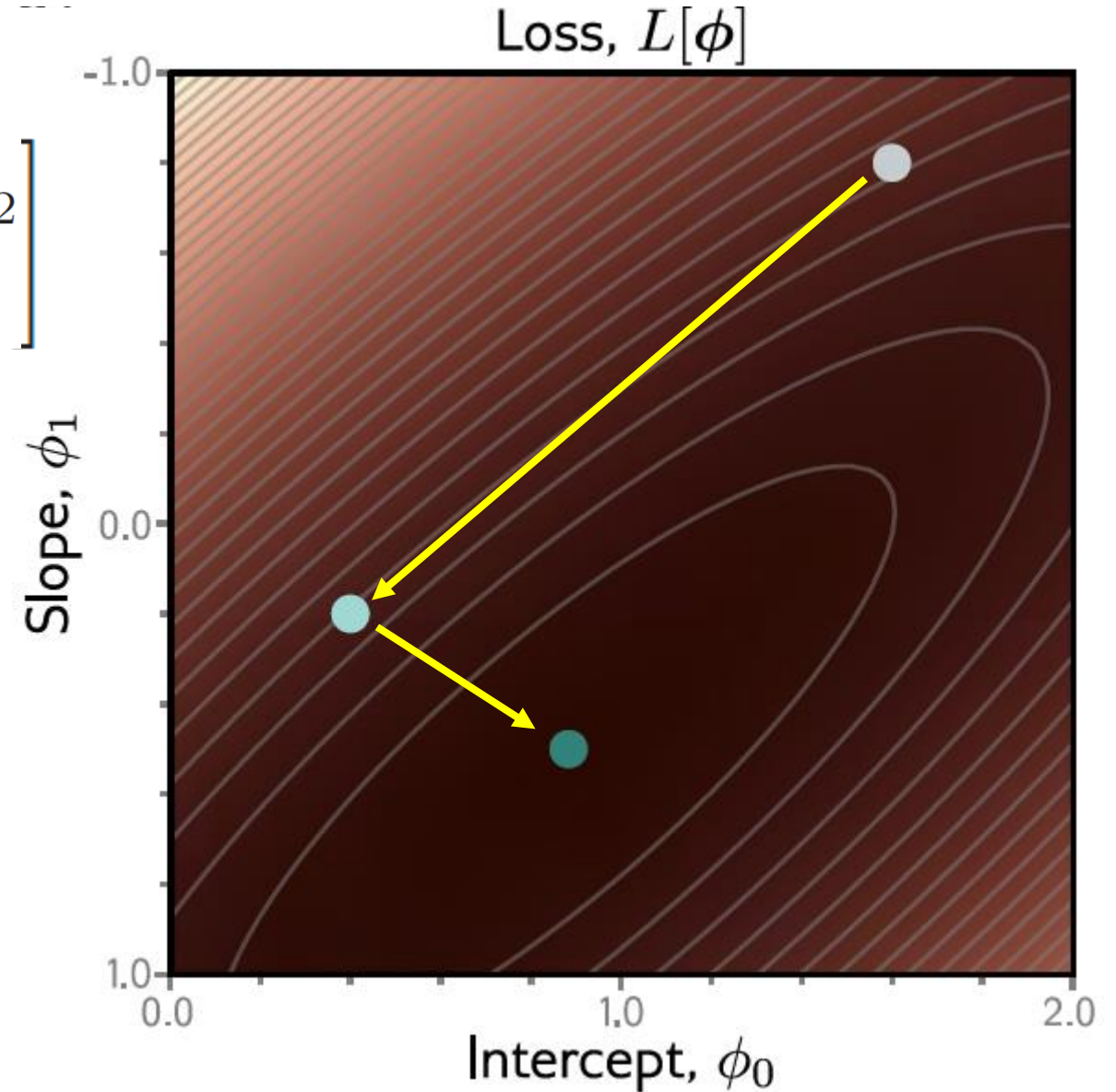
$$\sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2$$

Entrenamiento: ejemplo Regresión Simple

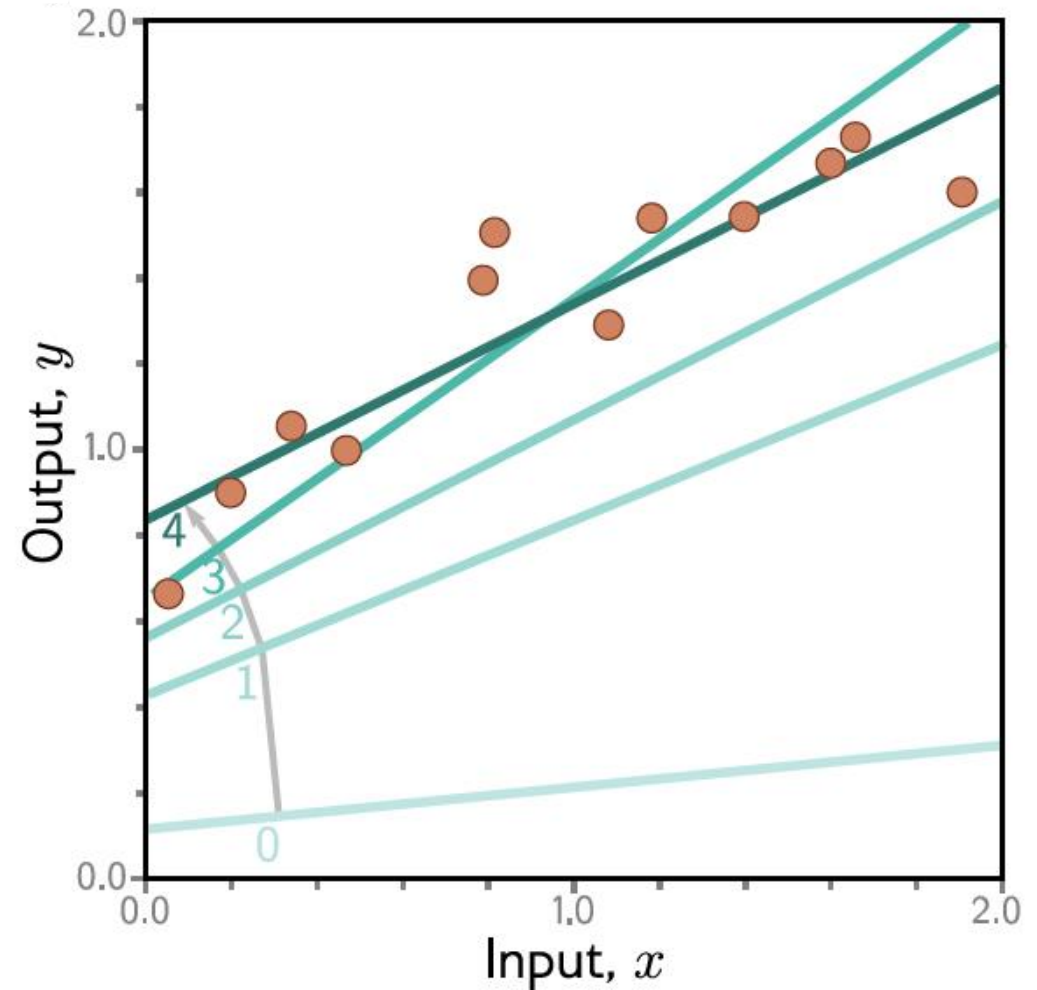
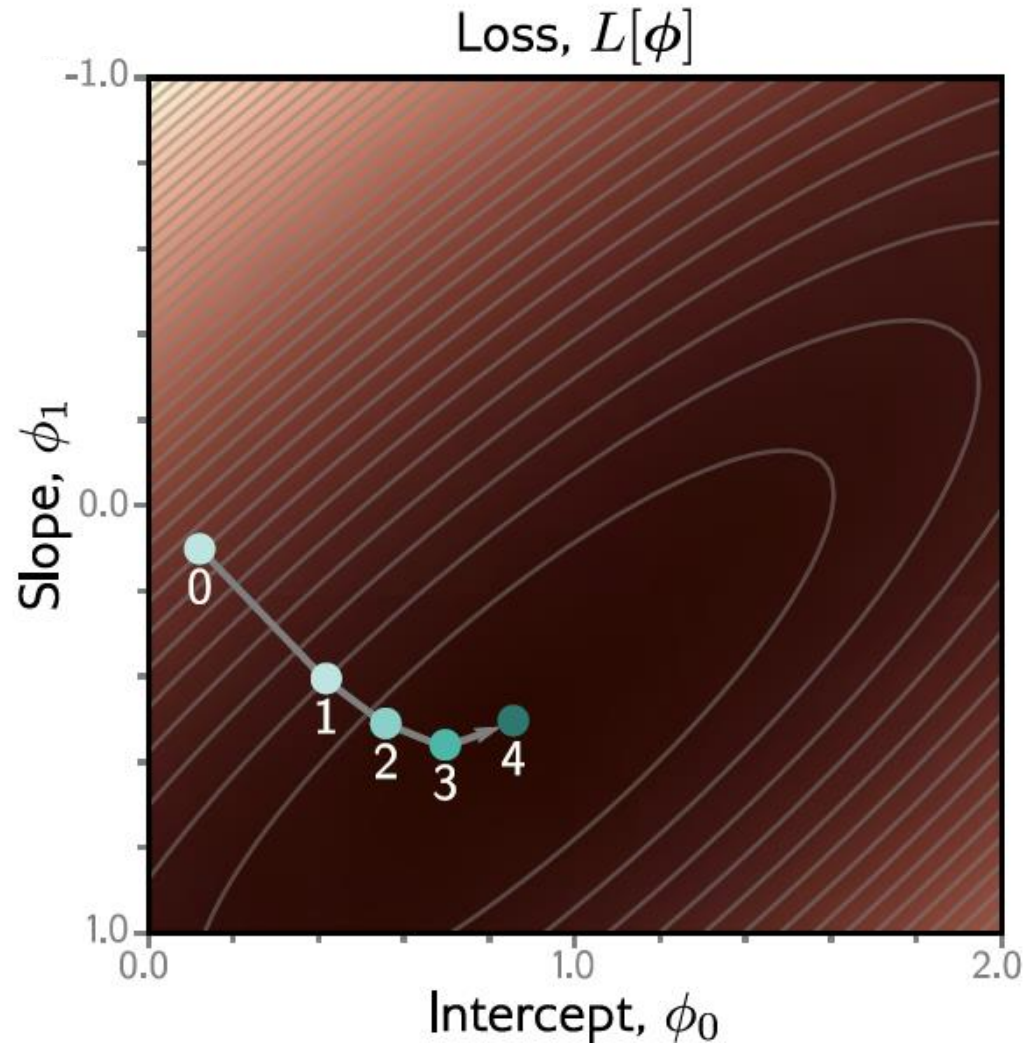


Entrenamiento: ejemplo Regresión Simple

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} \left[\sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2 \right]$$



Entrenamiento: ejemplo Regresión Simple



Entrenamiento, sobreajuste y subajuste

Aprendizaje supervisado - Entrenamiento

Entrenamiento con datos (de entrenamiento)

$$\hat{\phi} = \operatorname{argmin}_{\phi} \left[L[\phi] \right]$$

Capacidad de **generalizar** a otros datos: **datos de prueba**

Datos

Aprendizaje supervisado - Entrenamiento

Entrenamiento con datos (de entrenamiento)

$$\hat{\phi} = \operatorname{argmin}_{\phi} \left[L[\phi] \right]$$

Capacidad de **generalizar** a otros datos: **datos de prueba**

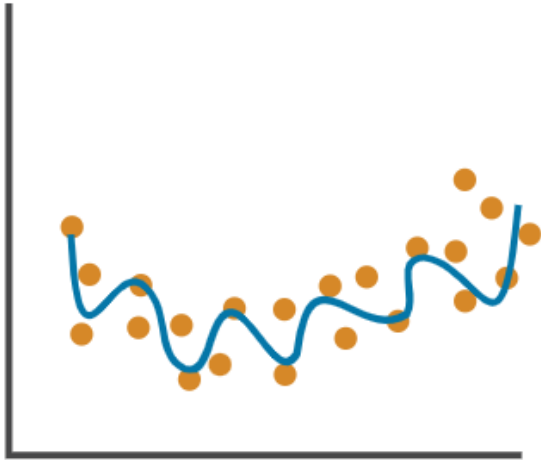
Datos de
entrenamiento

Datos de
prueba

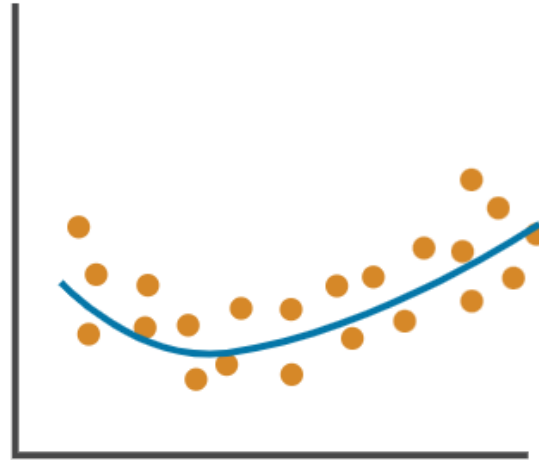
Sobreajuste
vs
subajuste

Aprendizaje supervisado - Entrenamiento

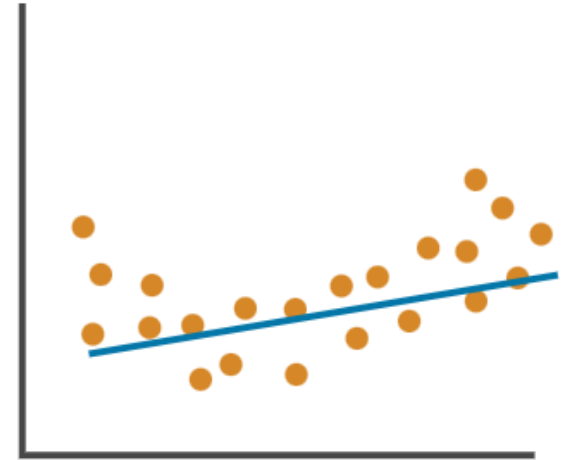
Sobreajuste



Ajuste



Subajuste



Validación cruzada

Usar **todos** los datos para entrenamiento.

Validación cruzada en k grupos (k-fold cross validation).

Realizar k veces la separación, con subconjuntos diferentes.



Datos de entrenamiento

Datos de prueba

Validación cruzada

Usar **todos** los datos para entrenamiento.

Validación cruzada en k grupos (k-fold cross validation).

Realizar k veces la separación, con subconjuntos diferentes.

Datos de prueba

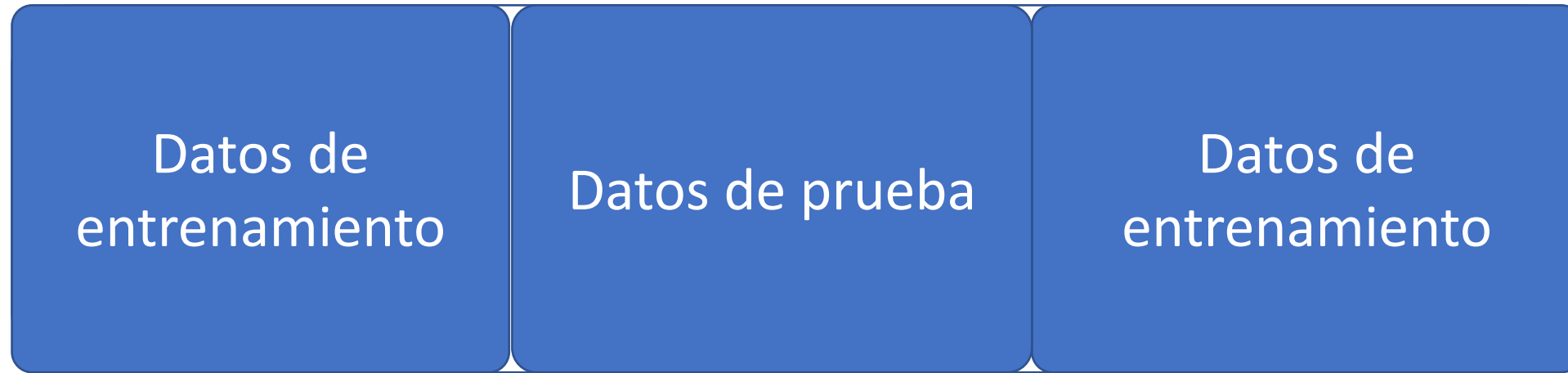
Datos de entrenamiento

Validación cruzada

Usar **todos** los datos para entrenamiento.

Validación cruzada en k grupos (k-fold cross validation).

Realizar k veces la separación, con subconjuntos diferentes.



Aprendizaje como Optimización

Medida de **desempeño (loss)**: función **objetivo**

Espacio de búsqueda:

- **Variables de decisión**: parámetros del modelo
- **Restricciones** sobre los parámetros

$$\hat{\phi} = \operatorname{argmin}_{\phi} \left[L[\phi] \right]$$

Modelos de regresión en Python (Scikit Learn)

Regresión en Python – Scikit learn

<https://scikit-learn.org/>

Aprendizaje supervisado:

https://scikit-learn.org/stable/supervised_learning.html

fit() + predict()

Regresión en Python – Scikit learn

Train test split:

https://scikit-learn.org/stable/modules/classes.html#module-sklearn.model_selection

[https://scikit-](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html)

[learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html)

Regresión en Python – Scikit learn

Cross validation:

https://scikit-learn.org/stable/modules/classes.html#module-sklearn.model_selection

[https://scikit-](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.cross_val_score.html#sklearn.model_selection.cross_val_score)

[learn.org/stable/modules/generated/sklearn.model_selection.cross_val_score.html#skl](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.cross_val_score.html#sklearn.model_selection.cross_val_score)

[earn.model_selection.cross_val_score](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.cross_val_score.html#sklearn.model_selection.cross_val_score)

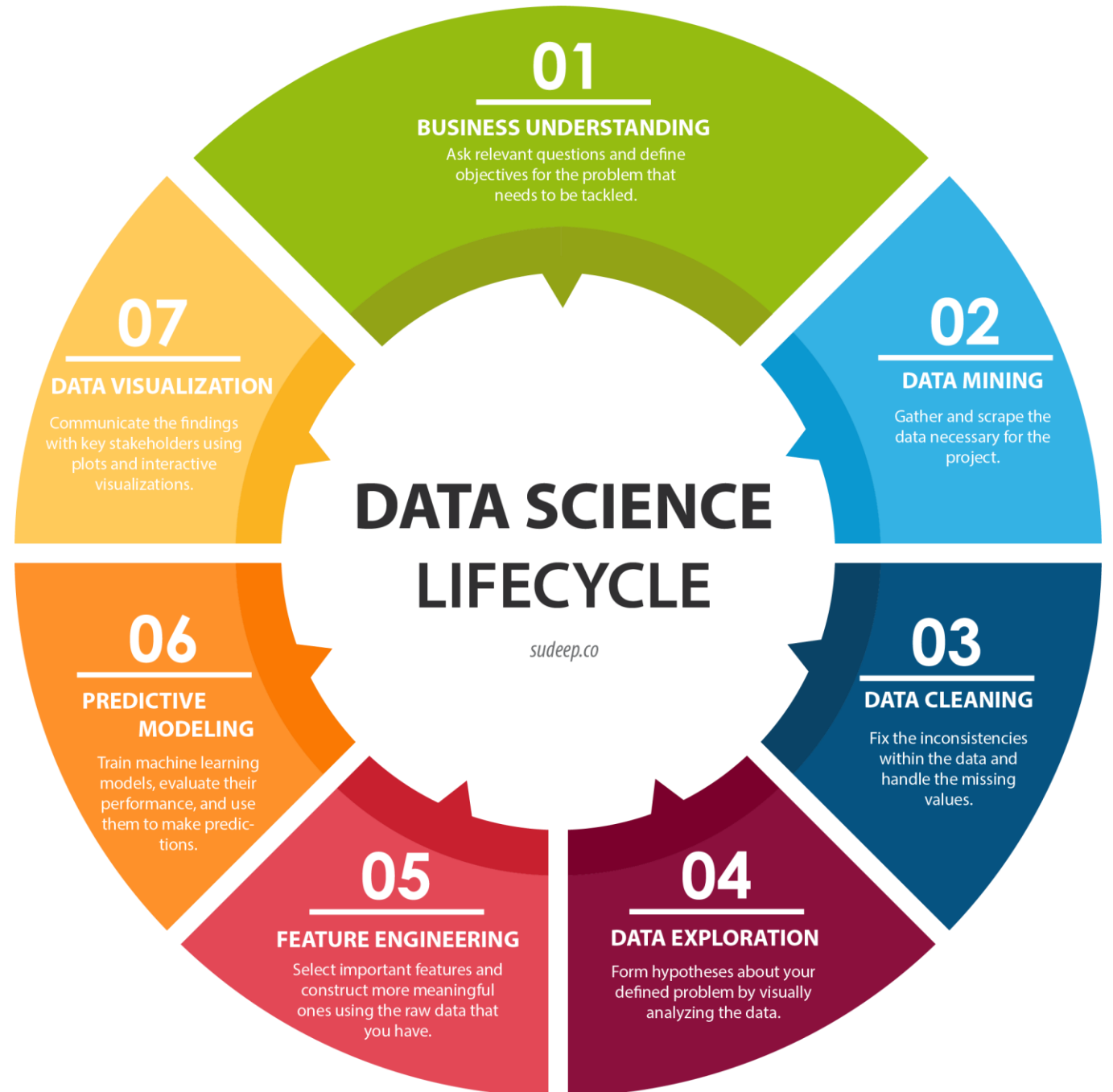
Regresión en Python – Statsmodels

API para regresión con OLS:

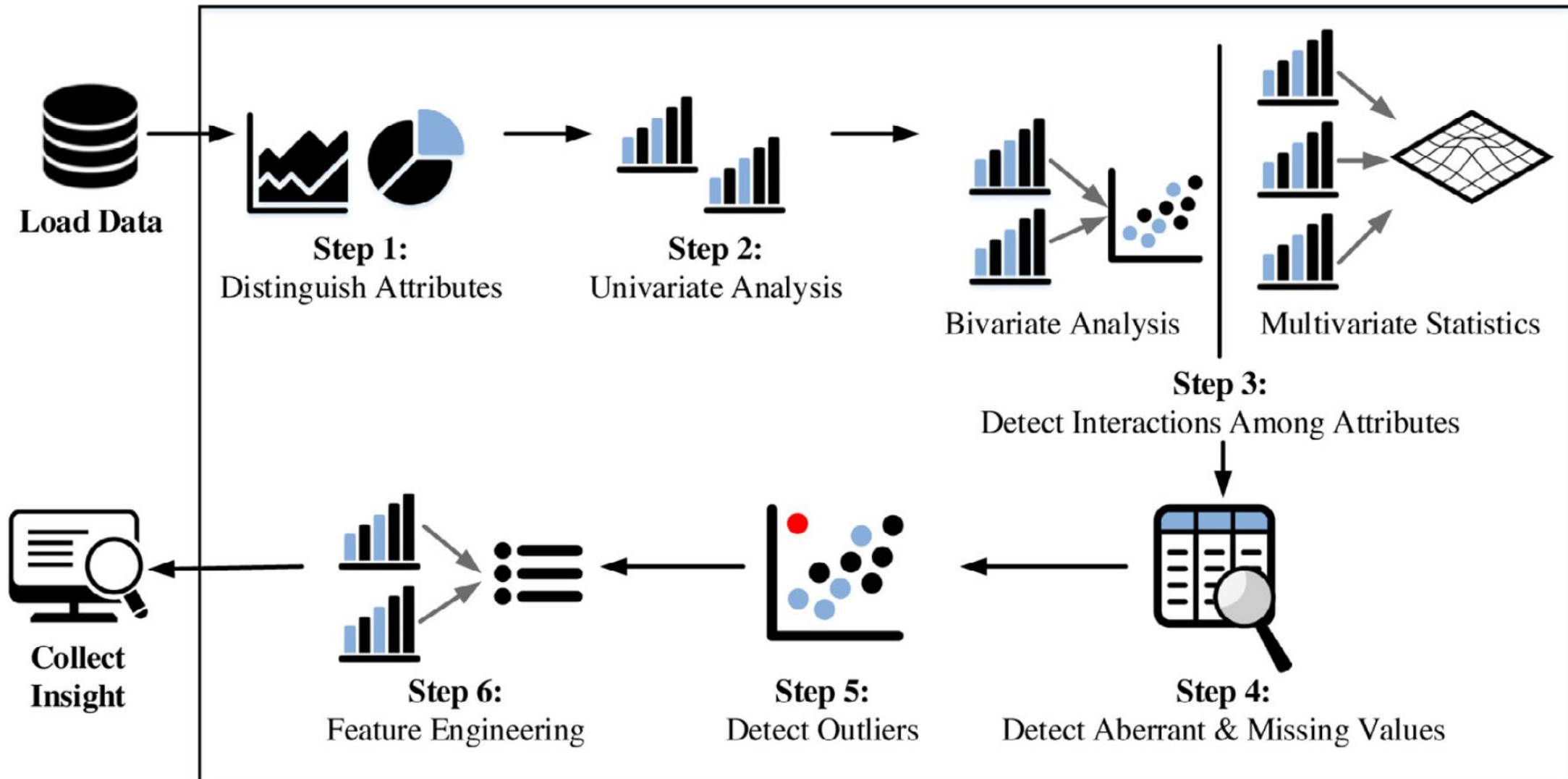
https://www.statsmodels.org/dev/generated/statsmodels.regression.linear_model.OLS.html

Metodologías para Análisis Exploratorio de Datos

Análisis Exploratorio (EDA) en Ciencia de Datos



Pasos de EDA



Lectura

A. Gosh (2018)

A comprehensive review of tools for exploratory analysis of tabular industrial datasets.

Sección 3.

