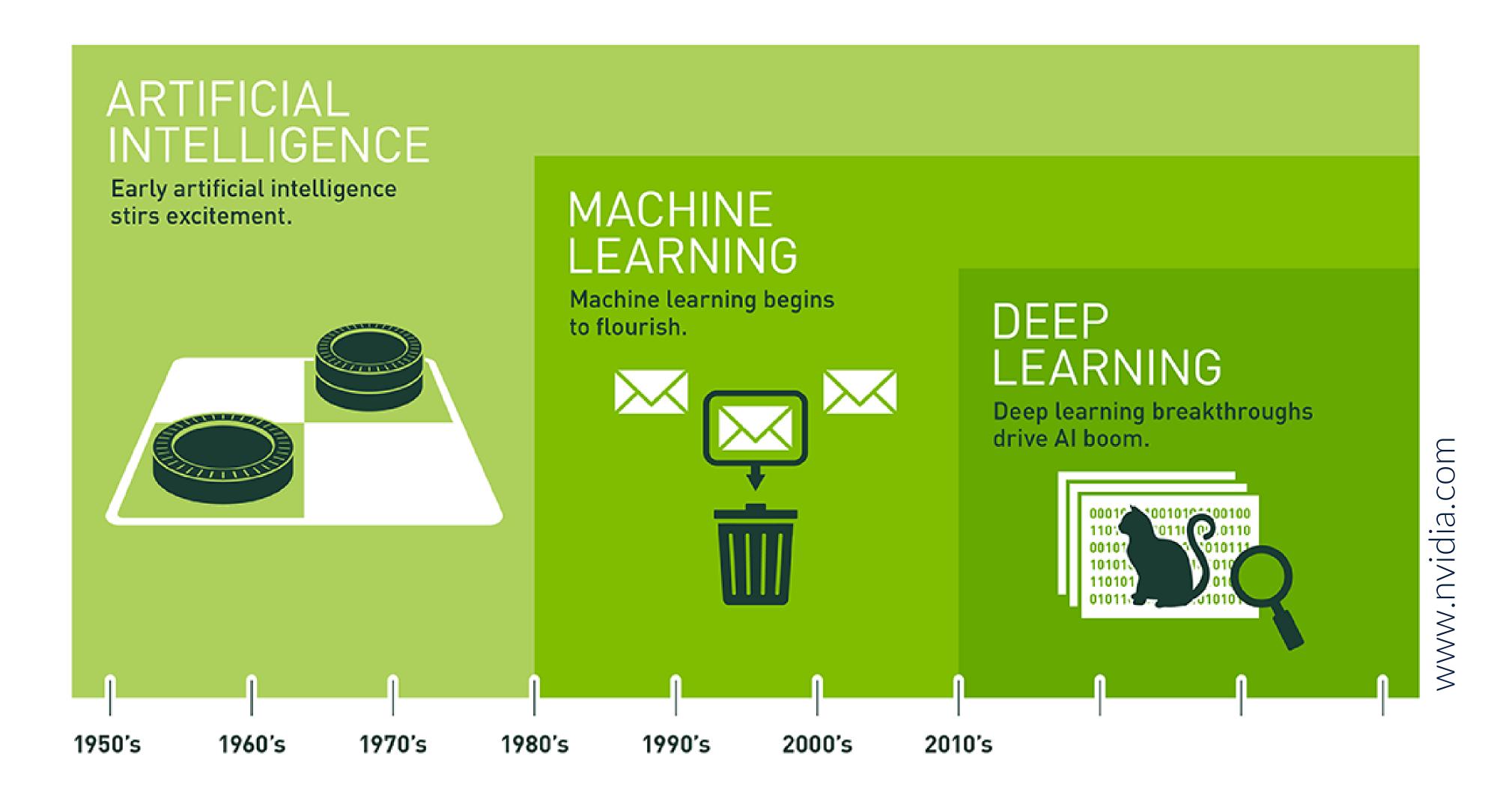
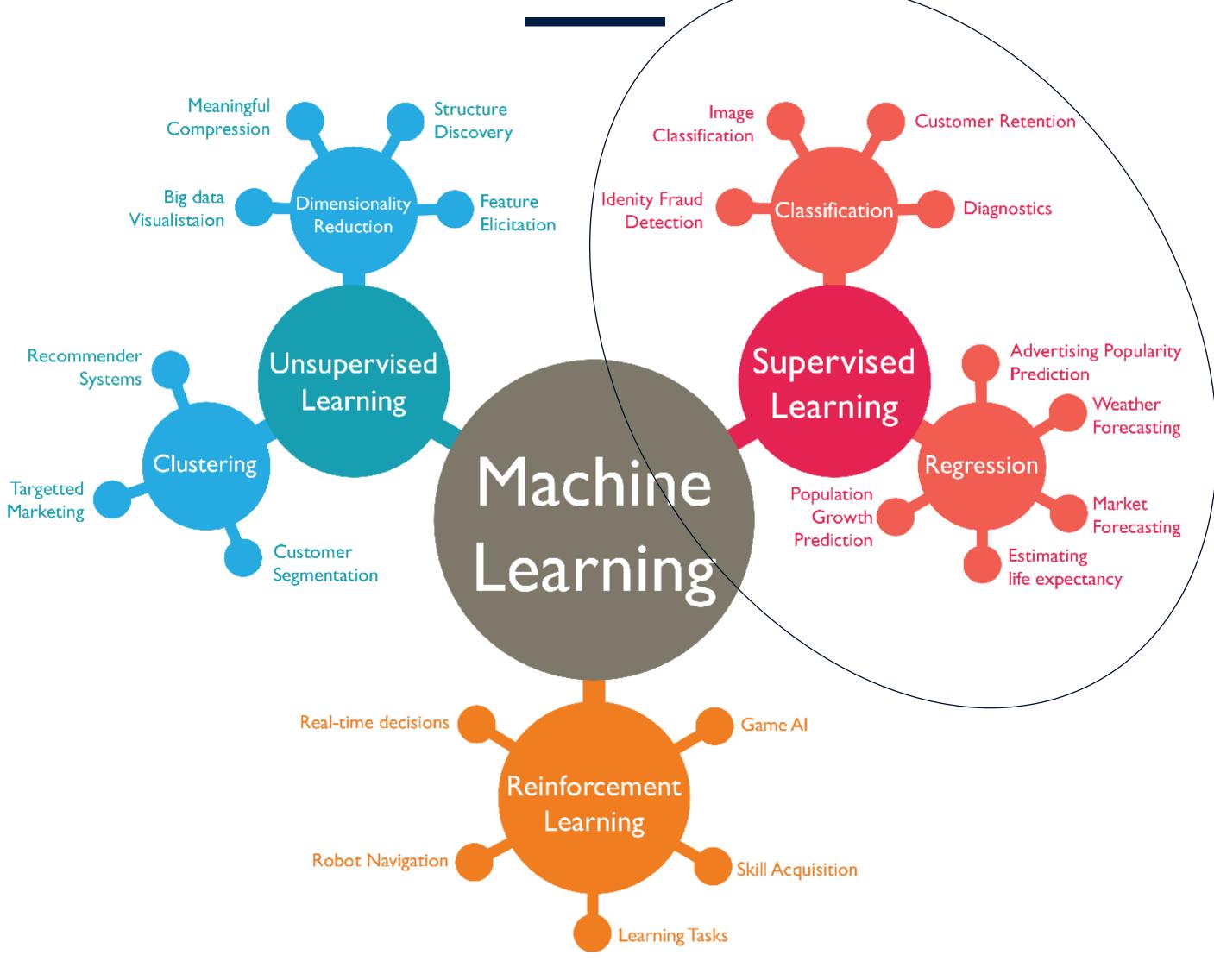


A glimpse of history



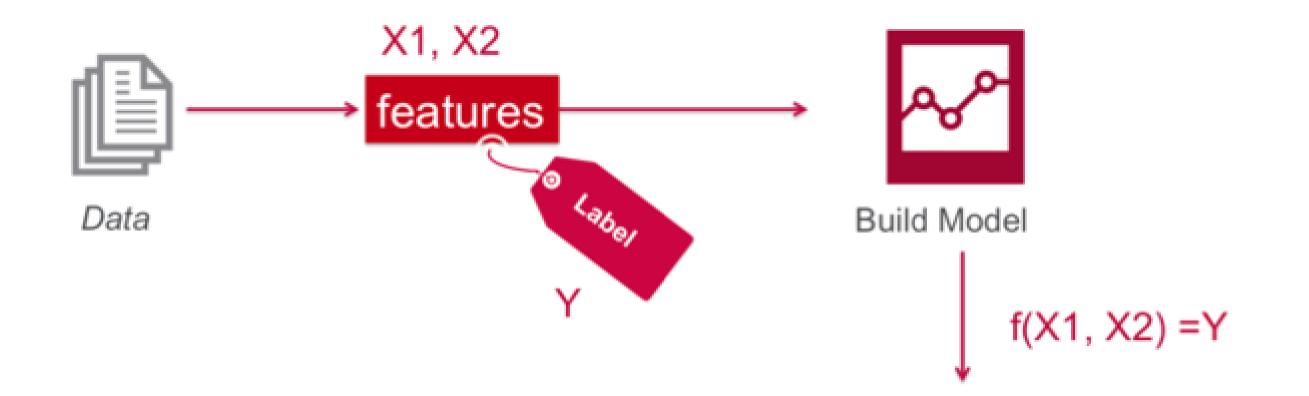




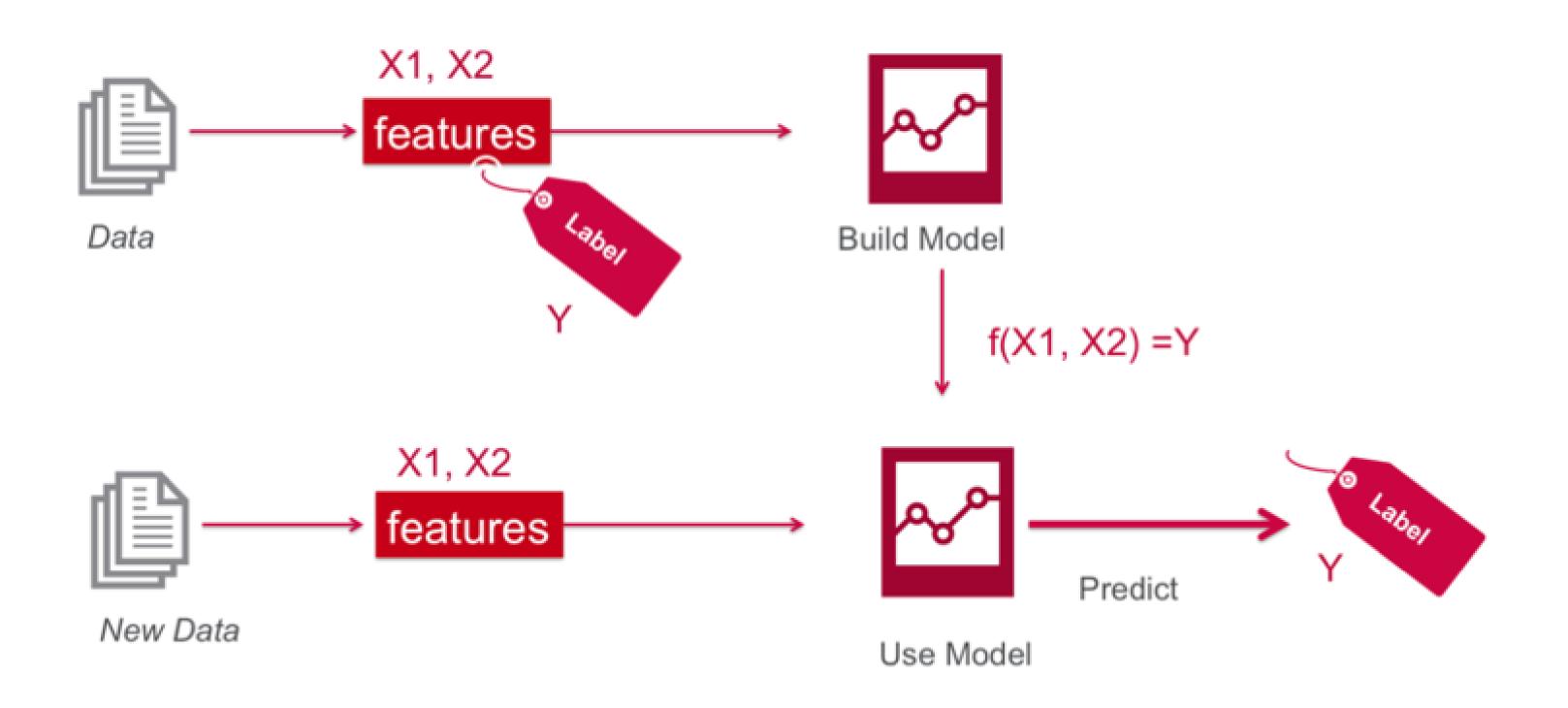




Supervised learning



Supervised learning



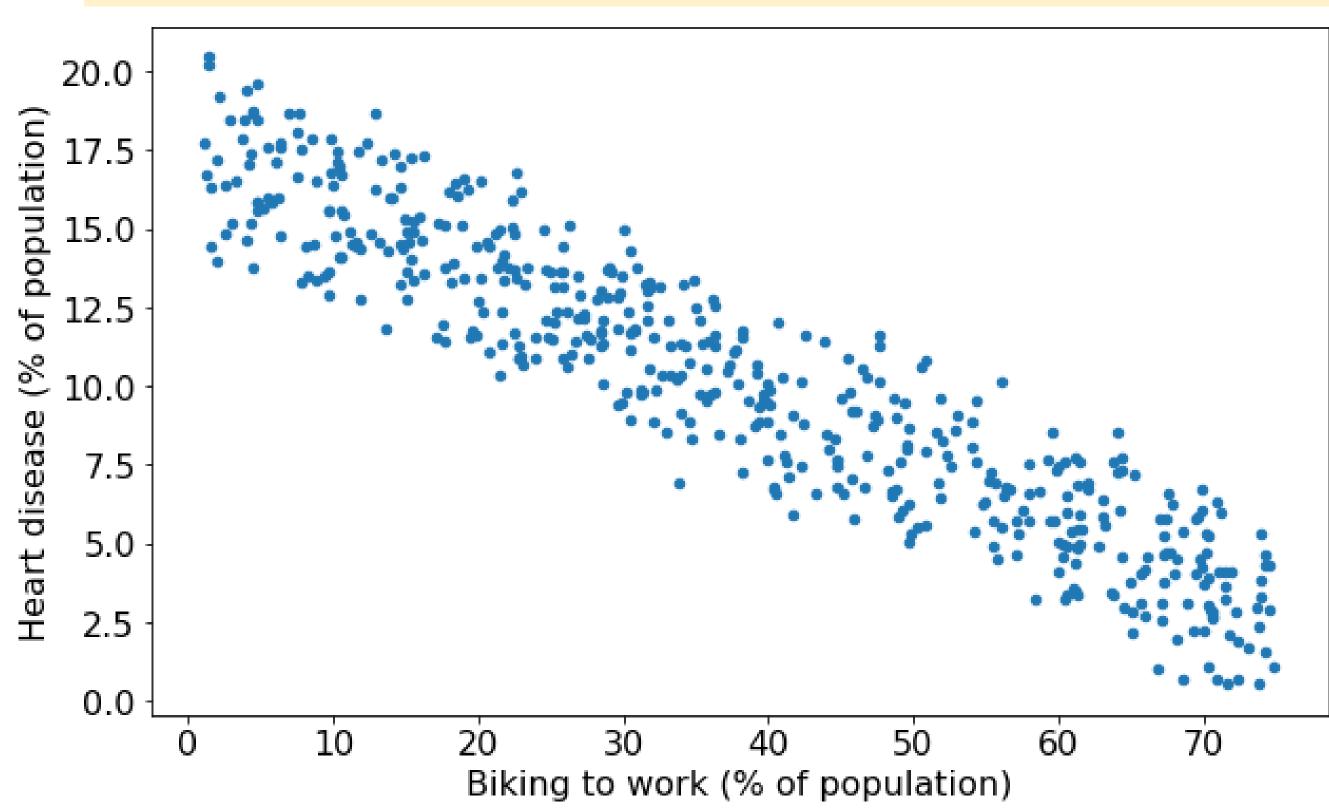




```
data = pd.read_csv("heartData.csv")
data.head()
```

	biking	smoking	heart.disease
0	30.801246	10.896608	11.769423
1	65.129215	2.219563	2.854081
2	1.959665	17.588331	17.177803
3	44.800196	2.802559	6.816647
4	69.428454	15.974505	4.062224

```
ax = data.plot.scatter('biking','heart.disease')
ax.set_xlabel("Biking to work (% of population)")
ax.set_ylabel("Heart disease (% of population)")
```





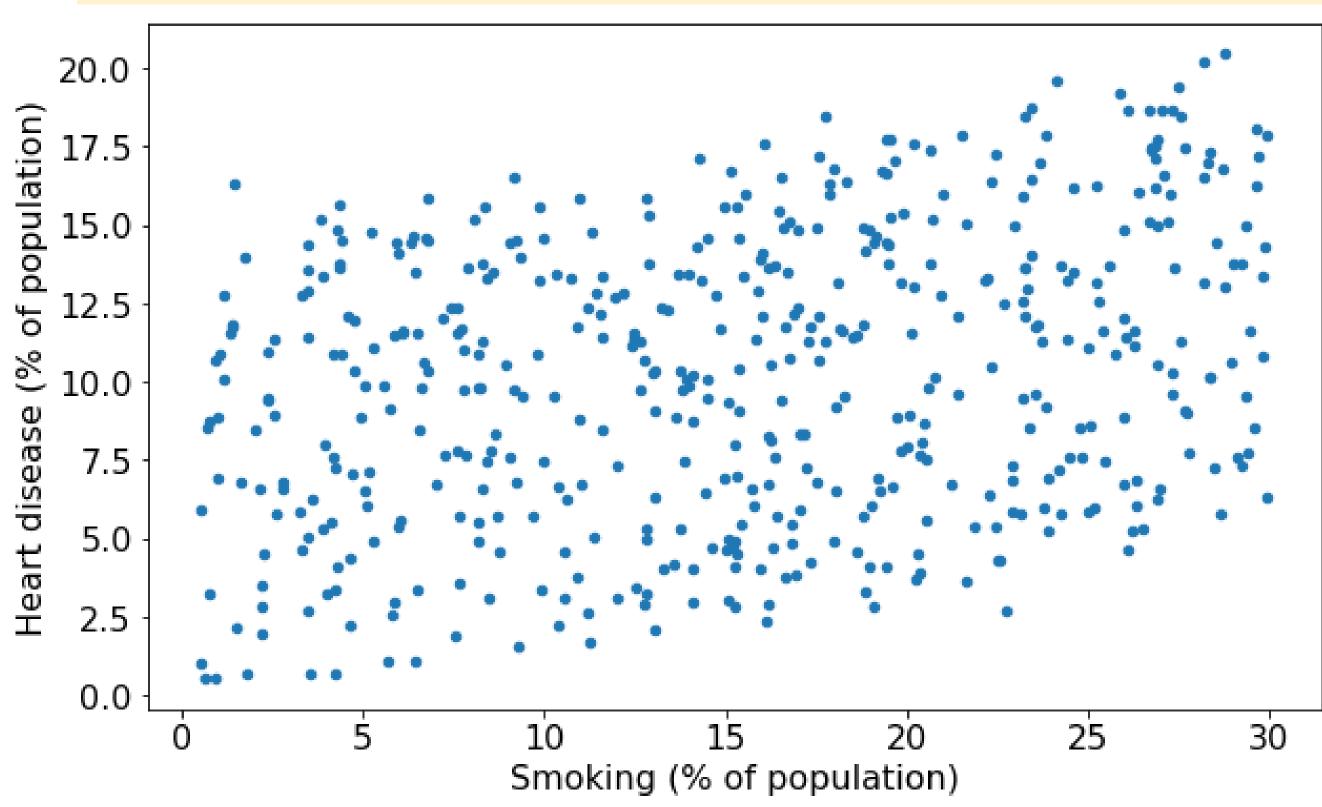


Heart disease vs Smoking

```
data = pd.read_csv("heartData.csv")
data.head()
```

	biking	smoking	heart.disease
0	30.801246	10.896608	11.769423
1	65.129215	2.219563	2.854081
2	1.959665	17.588331	17.177803
3	44.800196	2.802559	6.816647
4	69.428454	15.974505	4.062224

```
ax = data.plot.scatter(smoking','heart.disease')
ax.set_xlabel("Smoking (% of population)")
ax.set_ylabel("Heart disease (% of population)")
```







- Supervised learning model
- Linear & Univariate, hence the model is a straight line

$$y = w_0 + w_1 x + \epsilon$$
Iinear model error



- y is the Dependent Variable
- x is the Independent Variable
- w₁ is the **Slope** of the line
- w₀ is the **Coefficient** (or the y-intercept)

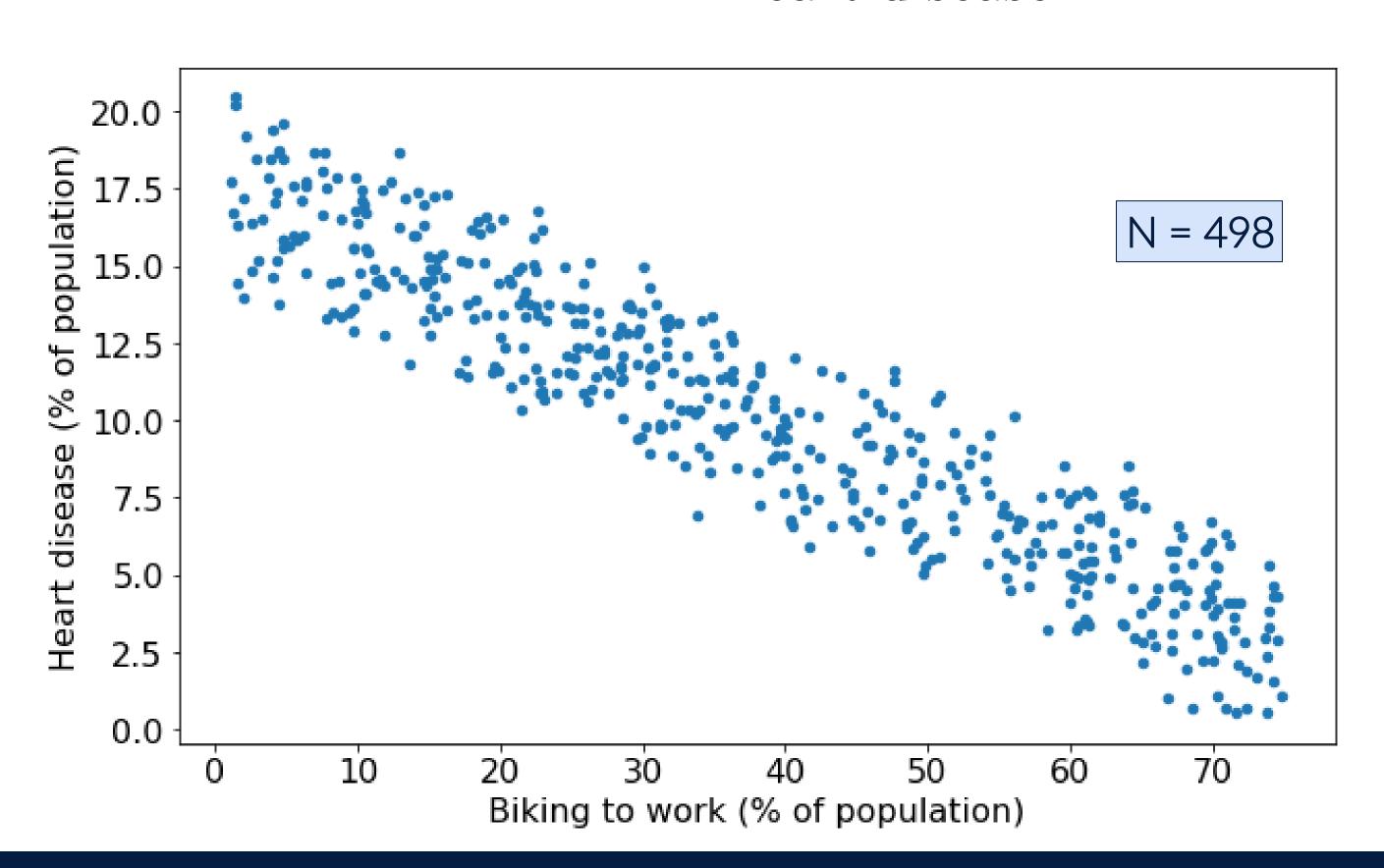
$$y = \underbrace{w_0 + w_1 x}_{\hat{y}} + \underbrace{w_1 x}_{\hat{y}} + \underbrace{\varepsilon}_{\hat{y}}$$

$$y = \underbrace{w_0 + w_1 x}_{\hat{y}} + \underbrace{\varepsilon}_{\sim \mathcal{N}(0, \sigma^2)}$$

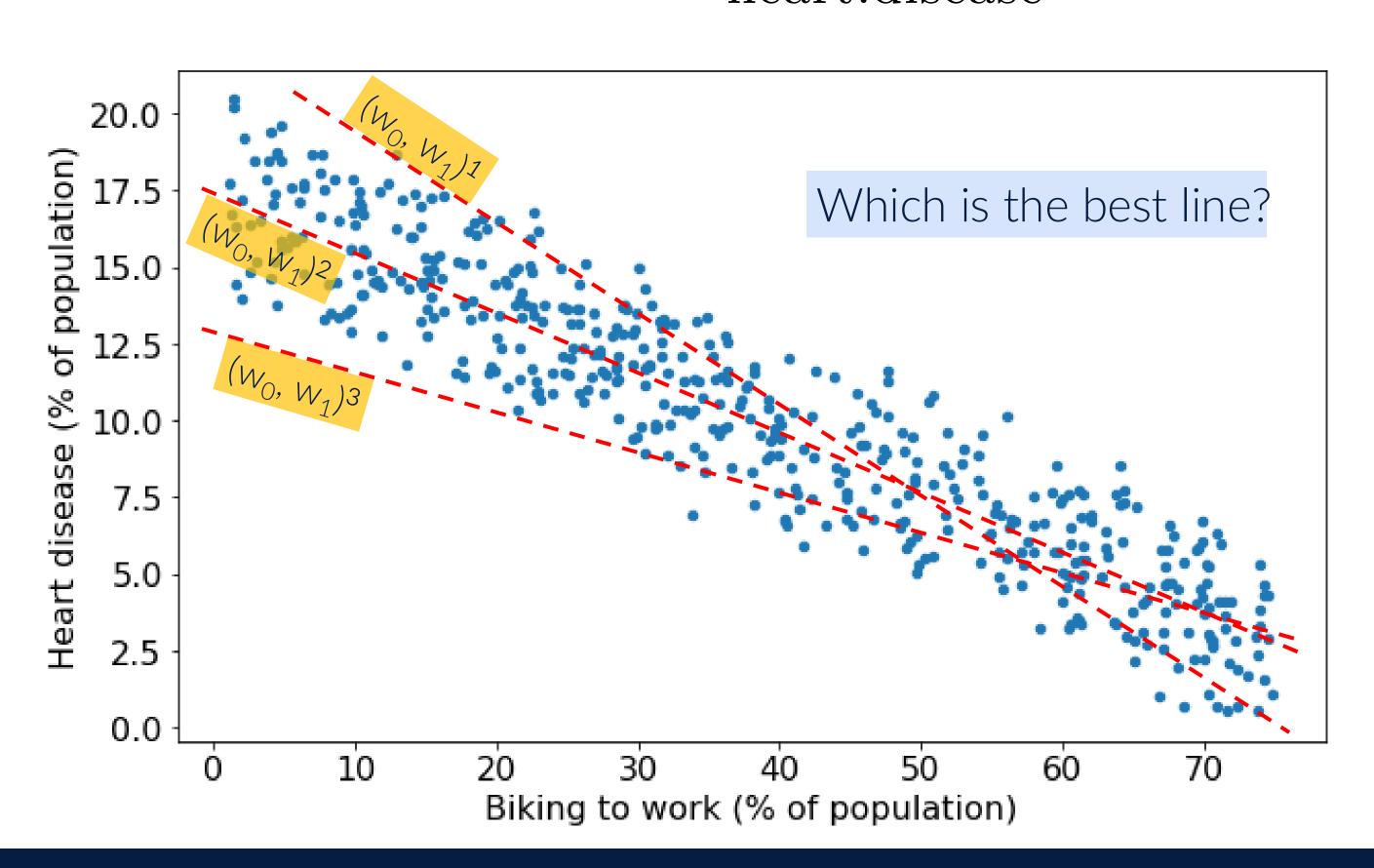
$$y = \underbrace{w_0 + w_1 x}_{\hat{y}} + \underbrace{\varepsilon}_{\sim \mathcal{N}(0, \sigma^2)}$$

$$\mathbf{A}$$
 $(x, y) =$

heart.disease =
$$w_0 + w_1$$
 biking $+\epsilon$ heart.disease



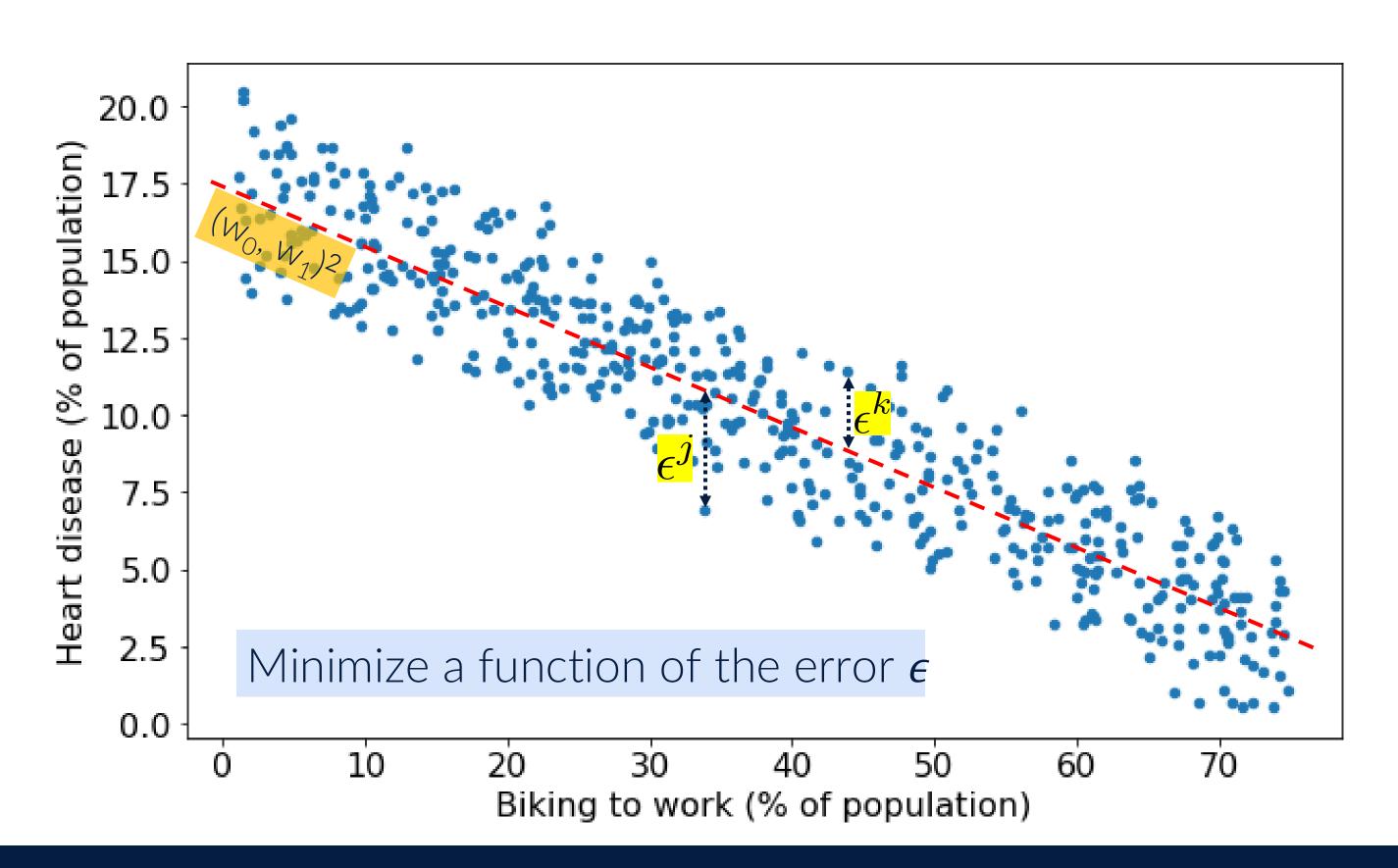
heart.disease =
$$w_0 + w_1$$
 biking $+\epsilon$ heart.disease



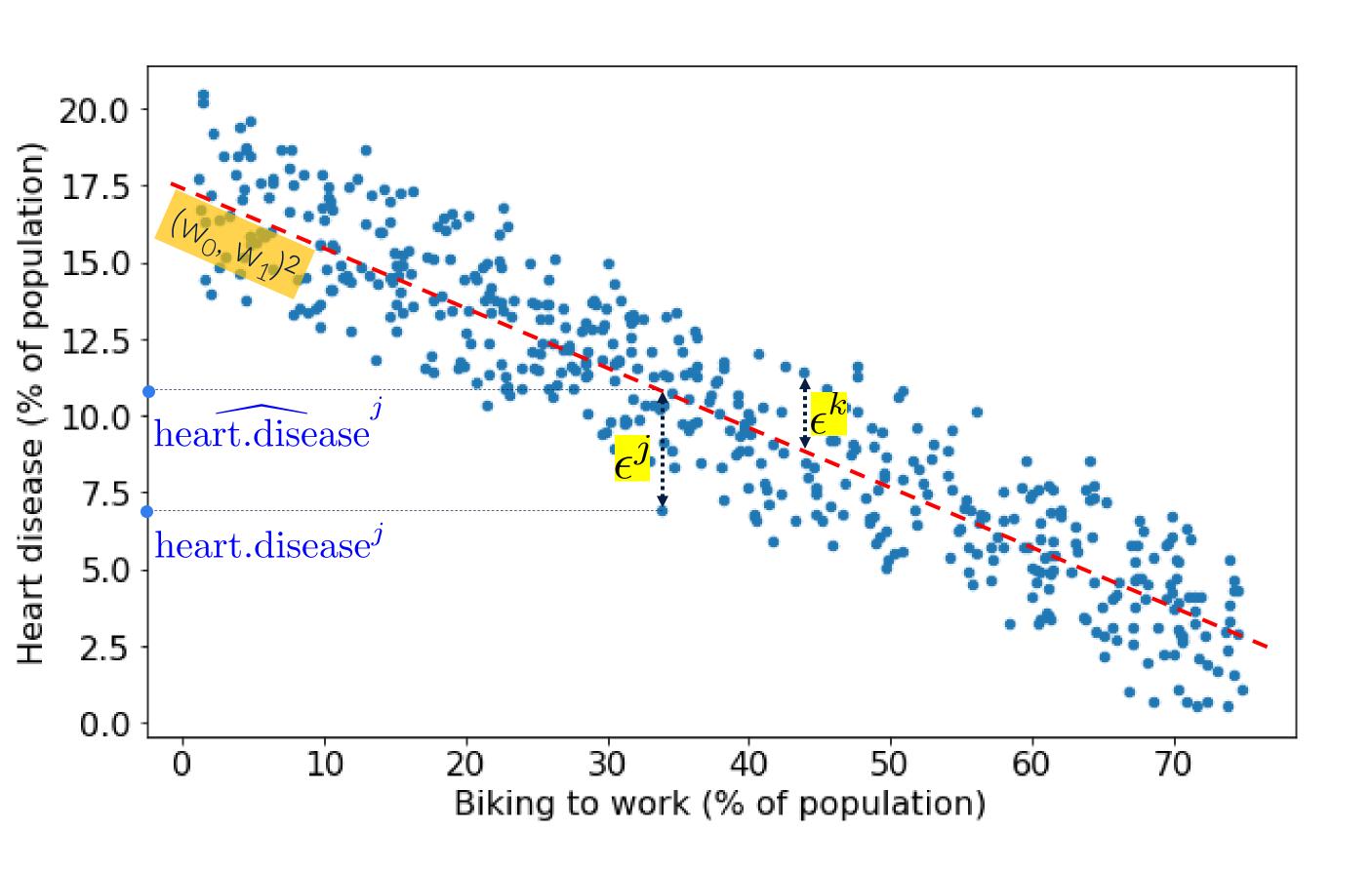




heart.disease =
$$w_0 + w_1$$
biking $+\epsilon$
heart.disease



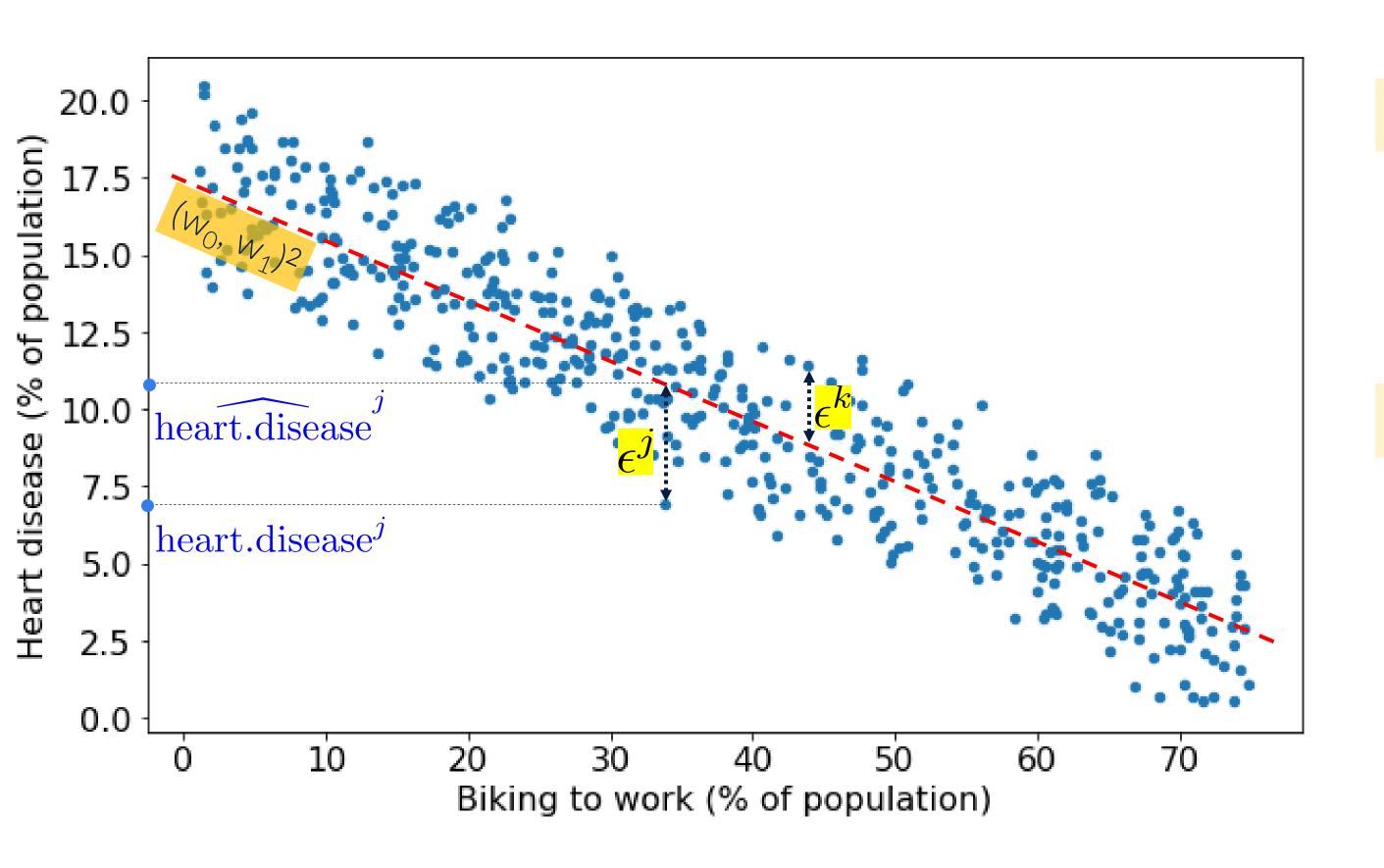
Error (or cost, loss) function



heart.disease =
$$w_0 + w_1$$
biking $+\epsilon$
heart.disease

$$\epsilon = \text{heart.disease} - \text{heart.disease}$$

Common cost functions for regression



Mean squared error (MSE)

$$\frac{1}{N} \sum_{i=1}^{N} \left(\widehat{\text{heart.disease}}^{i} - \widehat{\text{heart.disease}}^{i} \right)^{2}$$

Mean absolute error (MAE)

$$\frac{1}{N} \sum_{i=1}^{N} \left| \text{heart.disease}^{i} - \text{heart.disease}^{i} \right|$$

Where, heart.disease
$$i = w_0 + w_1 \text{biking}^i$$

The best model in the least squares

Ordinary least squares (OLS)

$$y = \underbrace{w_0 1 + w_1 x}_{\hat{y}} + \epsilon$$
 $\hat{y} = \tilde{\mathbf{X}} \mathbf{w} \text{ where, } \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \text{ and } \tilde{\mathbf{X}} = \begin{bmatrix} 1 & x^1 \\ 1 & x^2 \\ \vdots & \vdots \\ 1 & x^N \end{bmatrix}$

We seek to find \mathbf{w}^* that minimizes the MSE.

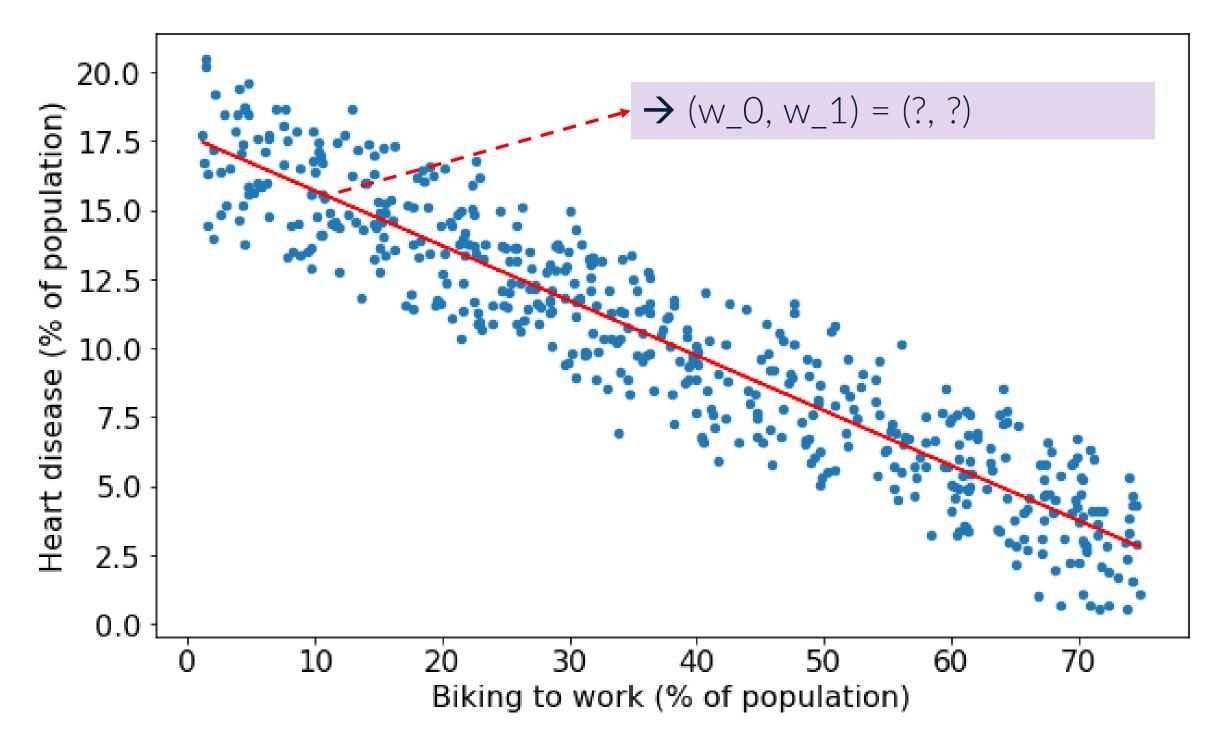
$$\mathbf{w}^* = \left(\mathbf{ ilde{ ilde{X}}^T ilde{ ilde{X}}}
ight)^{-1} \mathbf{ ilde{ ilde{X}}^T y}$$

An example in (pure) Python

Exercise: Find the best line in the least squares sense

using the formulation of the OLS solution

from the previous slide

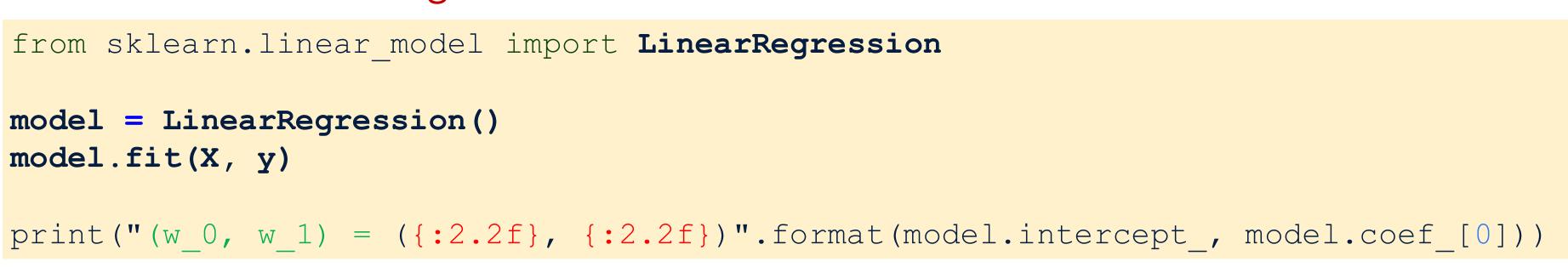


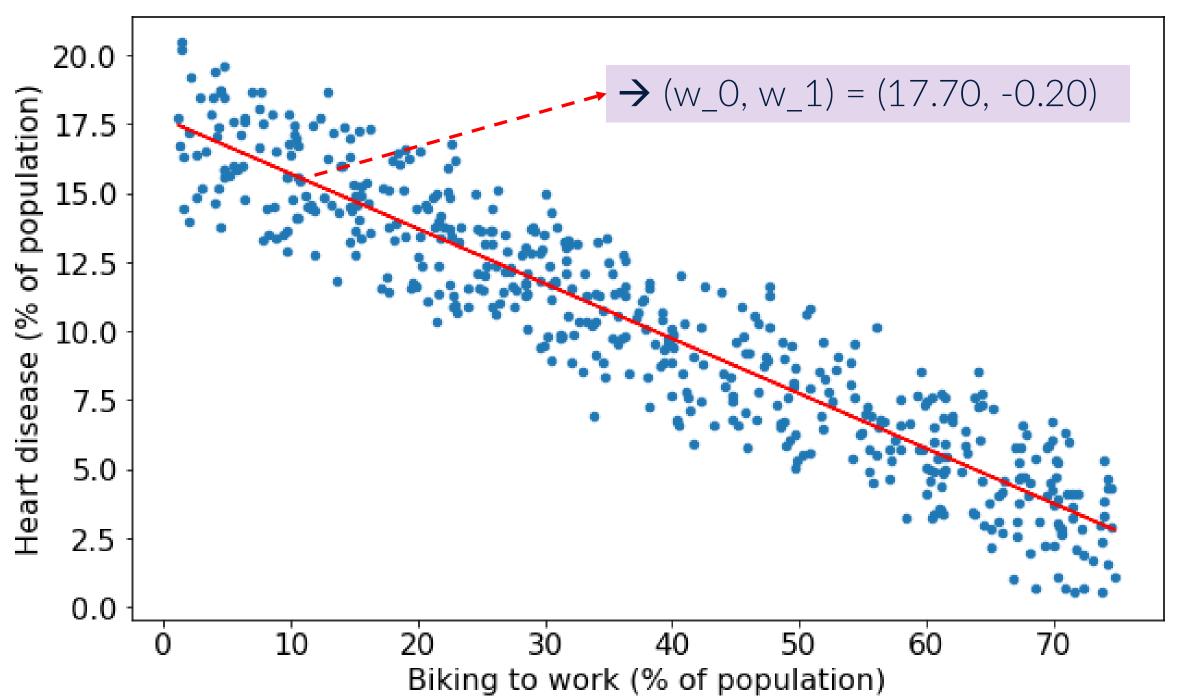
Now using Scikit-Learn

Data preparation

```
X = data.iloc[:,0:1].values
y = data.iloc[:,-1].values
```

w* via OLS method using Scikit-Learn







Extension of the univariate model

$$y = \underbrace{w_0 1 + w_1 x_1 + w_2 x_2 + \ldots + w_m x_m}_{\hat{y}} + \epsilon$$

$$\hat{y} = \tilde{\mathbf{X}}\mathbf{w} \text{ where, } \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_m \end{bmatrix} \text{ and } \tilde{\mathbf{X}} = \begin{bmatrix} 1 & x_1^1 & x_2^1 & \dots & x_m^1 \\ 1 & x_1^2 & x_2^2 & \dots & x_m^2 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_1^N & x_2^N & \dots & x_m^N \end{bmatrix}$$

OLS solution:
$$\mathbf{w}^* = \left(\mathbf{\tilde{X}^T \tilde{X}}\right)^{-1} \mathbf{\tilde{X}^T y}$$

Back to our example

Using both features / independent variables

heart.disease = $w_0 1 + w_1 \text{biking} + w_2 \text{smoking} + \epsilon$

heart.disease

Exercise: Find the best (hyper)plane in the least squares sense, and comment on the results.

biking



Three common metrics

• Mean Absolute Error
$$MAE = \frac{1}{N} \sum_{i=1}^{N} |y_i - \hat{y}_i|$$

• (Root) Mean Squared Error
$$MSE = \frac{1}{N} \sum_{i=1}^{N} (y_i - \hat{y}_i)^2$$

Since the error is squared, any prediction error is being heavily penalized

Linear Regression as Predictive Model

- An example
 - Let's take our multivariate model
 - Can we predict the proportion of **heart disearse** for (**biking** = 30%, **smoking** = 5.5%)?
 - Note the tuple (biking = 30%, smoking = 5.5%) does not exist in the original dataset = [[30.0, 5.5]]

```
prediction = model_multivariate.predict(X_pred)

print(">>> Prédiction maladie cardiaque: {:2.1f}%".format(prediction[0]))
```

>>> Prédiction maladie cardiaque: 10.0%



How well do we do though?

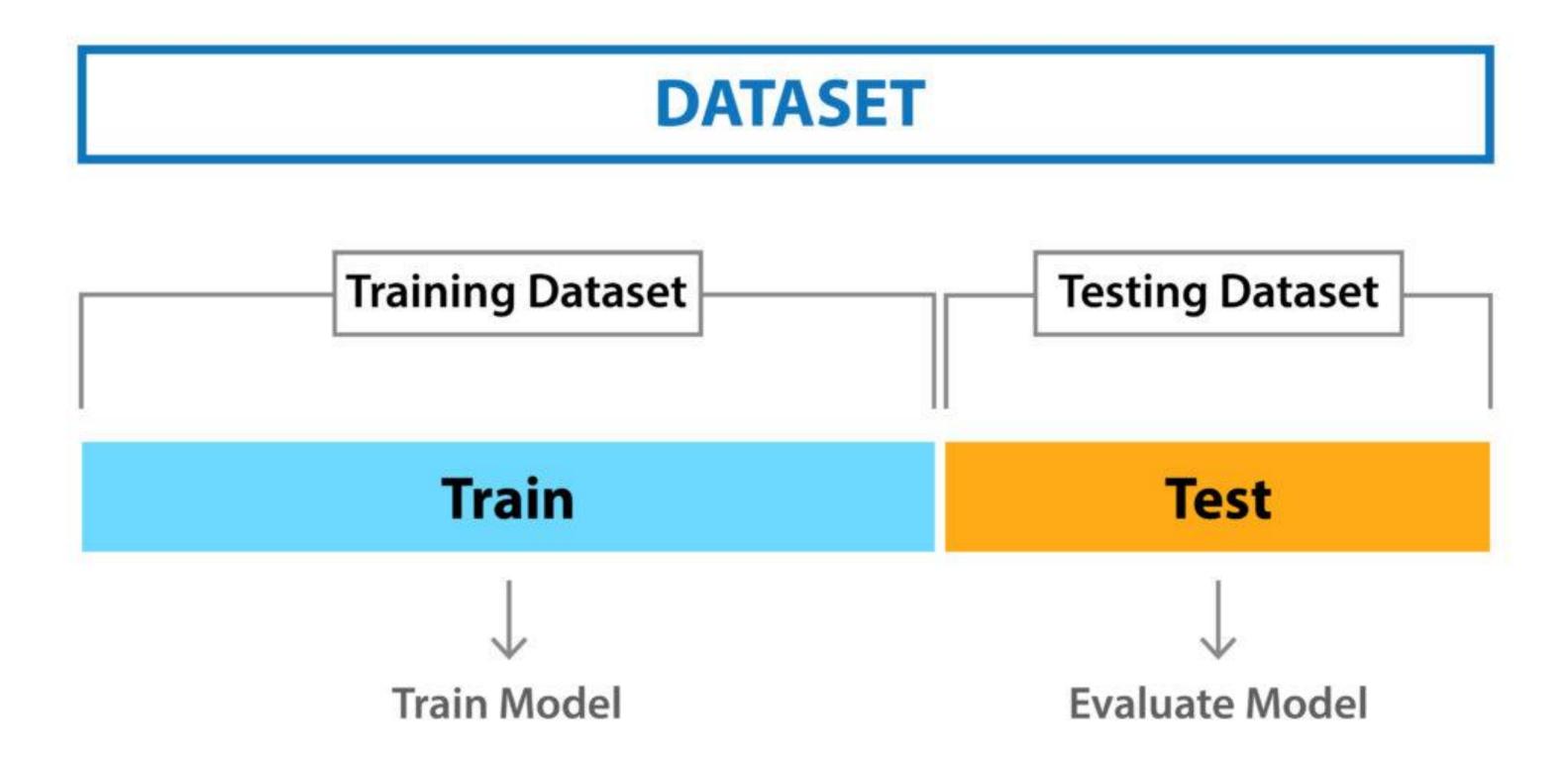
- The train-test split procedure is used to estimate the performance of machine learning algorithms when they are used to make predictions on data not used to train the model
 - Train Dataset: Used to fit the machine learning model
 - Test Dataset: Used to evaluate the fit machine learning model

Common split percentages include:

- Train: 80%, Test: 20%
- Train: 67%, Test: 33%
- Train: 50%, Test: 50%



The train-test split procedure



from sklearn.model_selection import train_test_split

x_train, x_test, y_train, y_test = train_test_split(X, y, test_size=0.2)





What did I learn?

- Supervised vs unsupervised learning
- Linear regression as supervised learning model
- Univariate vs multivariate linear regression
- Error / cost / loss function
- OLS solution: method for choosing the unknown parameters in a linear regression model by the principle of least squares (not least absolutes!)
- Performance evaluation of a regression model
 - Common metrics: R², MSE, MAE
- Train-test split procedure

