

# N-fold Cross-Validation

The **Leave-one-out Cross-Validation** method (or N-fold Cross-Validation) is a resampling and data partitioning method which aims at estimating the prediction error of a model.

In the context of regressions models, the prediction error  $\hat{\theta}_{pe}$  is defined as follows:

$$\hat{\theta}_{pe} = E[e_i^2] = E[(y_i - \hat{y}_i)^2] = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Then the **Root Mean Squared Error** (RMSE) is simply

$$RMSE = \sqrt{\hat{\theta}_{pe}} = \sqrt{\frac{1}{n} \sum_{i=1}^n e_i^2}$$

# N-fold Cross-Validation algorithm

To implement the leave-one-out cross-validation method for regression models, the general procedure is the following:

For  $i$  in  $1, \dots, n$

1. Compute  $\hat{y}_i^{(-i)}$  the predictions from the model fitted on all observations except for the  $i^{th}$  observation.
2. Compute the errors  $e_i = y_i - \hat{y}_i^{(-i)}$
3. Compute the RMSE as  $\sqrt{\frac{1}{n} \sum_{i=1}^n e_i^2}$

# Working example

**Example:** We generate  $n = 300$  artificial data points  $(x_i, y_i)$  from the following process:

$$x_i \sim N(5, 3)$$

$$y_i = x_i^2 + \epsilon_i \quad \text{with } \epsilon_i \sim N(0, 8)$$

Suppose that we wish to fit the two following models:

Model 1:  $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$

Model 2:  $y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \epsilon_i$

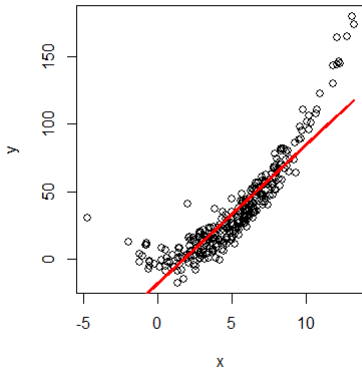
We will use the Leave-one-out Cross-Validation method in R and Python to select the best model, that is the model with the lowest RMSE.

# Code for scatterplots and fit in R

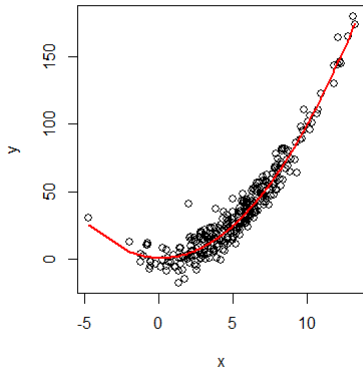
```
1 # 1. generate artificial data
2 set.seed(2023)
3 n <- 300
4 x <- rnorm(n = n, mean = 5, sd = 3)
5 y <- x^2 + rnorm(n = n, mean = 0, sd = 8)
6 data = data.frame(x = x, y = y)
7
8 # 2. plot the models that we want to fit
9 par(mfrow = c(1,2), pty = "s")
10
11 model1 <- lm(y ~ x)
12 plot(x = x, y = y, main = 'Linear model', cex = 1.1, pch = 1, lwd = 1.2)
13 yhat1 <- model1$coef[1] + model1$coef[2] * x
14 lines(x, yhat1, lw = 2.5, col = 'red')
15
16 model2 <- lm(y ~ x + I(x^2))
17 plot(x = x, y = y, main = 'Quadratic model', cex = 1.1, pch = 1, lwd = 1.2)
18 yhat2 <- model2$coef[1] + model2$coef[2] * x + model2$coef[3] * x^2
19 lines(x[order(x)], yhat2[order(x)], lw = 2.5, col = 'red')
```

# Plots in R

**Linear model**



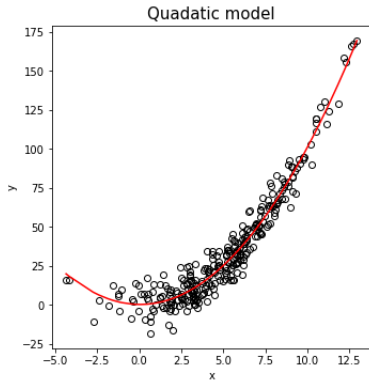
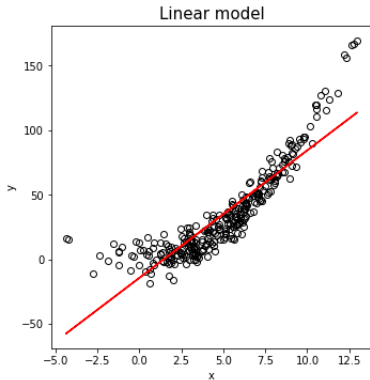
**Quadratic model**



# Code for scatterplots and fit in Python

```
1 # 1. generate artificial data
2 import numpy as np
3 np.random.seed(2023) ; n = 300
4 x = np.random.normal(loc = 5, scale = 3, size = n)
5 y = x**2 + np.random.normal(loc = 0, scale = 8, size = n)
6
7 # 2. plot the models that we want to fit
8 import matplotlib.pyplot as plt
9
10 plt.subplot(1, 2, 1)
11 plt.plot(x, y, 'o', fillstyle = 'none', color = 'black')
12 plt.title('Linear model', fontsize = 15)
13 plt.xlabel("x") ; plt.ylabel("y")
14 beta1, beta0 = np.polyfit(x, y, 1) ; yhat1 = beta0 + beta1*x
15 plt.plot(x, yhat1, color = 'red')
16
17 plt.subplot(1, 2, 2)
18 plt.plot(x, y, 'o', fillstyle = 'none', color = 'black')
19 plt.xlabel('x') ; plt.ylabel('y') ; plt.title('Quadratic model', fontsize = 15)
20 beta2, beta1, beta0 = np.polyfit(x, y, 2) ; yhat2 = beta0 + beta1*x + beta2*(x
    **2)
21 orders = np.argsort(x.ravel())
22 plt.plot(x[orders], yhat2[orders], color = 'red')
```

# Plots in Python



# N-fold Cross-Validation in R

```
1 # 3. Cross-Validation
2 # fit the models on leave-one-out samples
3
4 pred.cv.mod1 <- pred.cv.mod2 <- numeric(n)
5
6 for(i in 1:n) {
7
8   # quadratic model
9   mod1 = lm(y ~ x, subset = -i)
10  pred.cv.mod1[i] = predict(mod1, data[i,])
11
12  # quadratic model
13  mod2 = lm(y ~ x + I(x^2), subset = -i)
14  pred.cv.mod2[i] = predict(mod2, data[i,])
15
16 }
```



# N-fold Cross-Validation in Python

```
1 # 3. Cross-Validation
2 # fit the models on leave-one-out samples
3
4 import pandas as pd
5
6 data = pd.DataFrame({'x': x, 'y': y})
7 xn = data['x'].values.reshape(-1,1); yn = data['y'].values.reshape(-1,1)
8
9 from sklearn.preprocessing import PolynomialFeatures
10 from sklearn.model_selection import LeaveOneOut, cross_val_score
11
12 loocv = LeaveOneOut()
13
14 # linear model
15 mod1 = PolynomialFeatures(degree = 1, include_bias = False).fit_transform(xn)
16 mod11 = LinearRegression().fit(mod1, yn)
17 scoresmod1 = cross_val_score(mod11, mod1, yn,
18                               scoring = 'neg_mean_squared_error',
19                               cv = loocv)
20 # quadratic model
21 mod2 = PolynomialFeatures(degree = 2, include_bias = False).fit_transform(xn)
22 mod22 = LinearRegression().fit(mod2, yn)
23 scoresmod2 = cross_val_score(mod22, mod2, yn,
24                               scoring = 'neg_mean_squared_error',
25                               cv = loocv)
```

# Conclusion from R and Python outputs

In R:

```
1 MSE1 = (1/n) * sum((y - pred.cv.mod1)^2) # theta_hat_pe
2 MSE2 = (1/n) * sum((y - pred.cv.mod2)^2) # theta_hat_pe
3
4 # Root Mean Squared Error (RMSE)
5 sqrt(c(MSE1, MSE2))
6 # [1] 15.68599  7.99332
```

In Python

```
1 # Root Mean Squared Error (RMSE)
2 import statistics
3 import math
4
5 RMSE1 = math.sqrt(statistics.mean(abs(scoresmod1)))
6 RMSE2 = math.sqrt(statistics.mean(abs(scoresmod2)))
7 [RMSE1, RMSE2]
8 # [16.169293289892607, 7.873829105930071]
```

The second model (Quadratic) has the lowest RMSE and thus is preferred (to the Linear) in both cases. It was expected since the fit as described by the red line or curve on the graphs describe the generated data more accurately.

## Further reading and code

Rizzo, M.L. (2019). Statistical Computing with R, Second Edition (2nd ed.). Chapman and Hall/CRC.

<https://doi.org/10.1201/9780429192760>

The R Project for Statistical Computing:

<https://www.r-project.org/>

Python:

<https://www.python.org/>

Accessing R and Python code:

<https://github.com/JRigh/LOOCV-Cross-validation-for-regression>