

Home Assignment 1

Machine Learning A

FEDERICO FIORIO

September 12, 2022

1 Make Your Own (10 points)

1. Regarding the profile information, I think it would be interesting to find features which are correlated to the final grades.

These features could be:

- average grades in previous courses
- final grade in calculus
- final grade in statistics

So the sample space would be:

$$\mathcal{X} : \mathbb{R} \times \mathbb{N} \times \mathbb{N}$$

2. The label space would be the set of possible obtainable grades, considering the danish scale it is (maybe I'm wrong I don't know the danish scale yet, I'm sorry):

$$\mathcal{Y} : \{12, 10, 7, 4, 2, 0\}$$

3. The grades can be used as labels but also as numbers in order to perform regression, in case of regression, however when the algorithm will try to predict the results it will not obtain grades belonging to the \mathbb{N} set, so I will have to map each \mathbb{R} number predicted to the closest grade number after the prediction. For example 3,1 will become a 4.

Since I'm dealing with regression I can use the square-loss function:

$$l(y', y) = (y' - y)^2$$

to obtain grade that make sense, when I want to make a prediction the result has to be rounded to the closest grade

4. I would use the euclidean distance defined as:

$$d(x, x') = \sqrt{\sum_{i=1}^d (x_i - x'_i)^2}$$

but for KNN we could save computation just by doing

$$\sum_{i=1}^d (x_i - x'_i)^2$$

5. I would divide the given dataset into: S_{train} , S_{val} and S_{test} . On S_{train} I would train the various instances of the model, we can use KNN as example.

on S_{val} I would take the best K; so the best possible classifier based on the previously defined error function. (the one with minimized loss function, so here I am evaluating the models based on the error function)

On S_{test} I evaluate the performance of the algorithm on new samples based on the same error function (this time unbiased). This last step will give me an idea of how the algorithm behave when I try to deploy it.

6. yes I would expect some issues, maybe the data on which I trained the algorithm wasn't very representative, so the deployed algorithm would fail to generalize, in this case, I would need more data in order to try to improve the algorithm performances.

It might also happen that the algorithm learns too well the noise and patterns between the data points and at the end will fail to generalize, in this case, I would try to use simpler models in order to obtain better results.

2 Digits Classification with K Nearest Neighbors (45 points)

The validation error was computed as error rate, this means that for each validation set of size n for $n \in \{10, 20, 40, 80\}$ I computed the number of samples missclassified for each k and normalized by the size n .

So for example, in the first validation set (1 of 5) for $n = 10$, if I get for $k = 4$, 5 missclassified samples, it means that the error rate for $k = 4$ is $\frac{5}{10}$

The following are the 4 images for the 4 possible values of n , a line is represented by an array, an array has len 50 as the number of k requested, each number in the array represent the error rate for that k ($K = \text{number of the cell of the array} + 1$) and that specific validation set.

In each image there are the results of 5 (as the size of the set i) validation sets of size n .

Figure 1: validation sets of size 10

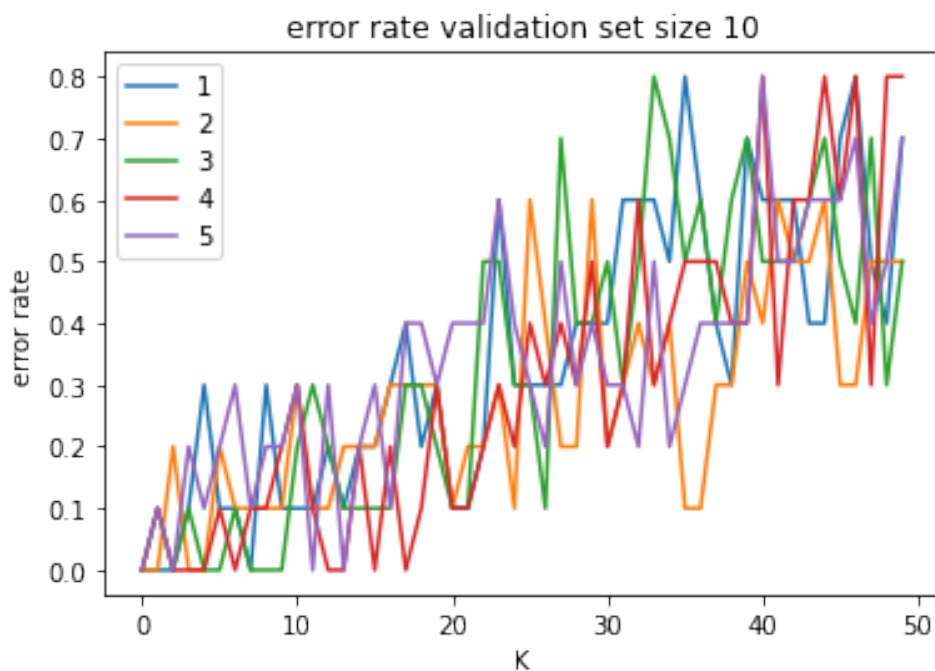


Figure 2: validation sets of size 20

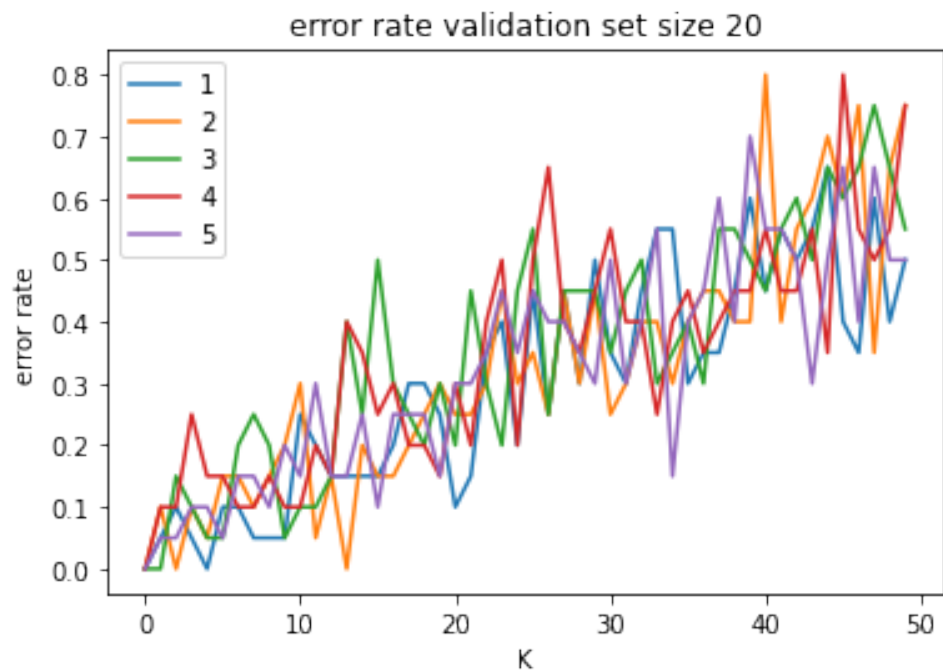


Figure 3: validation sets of size 40

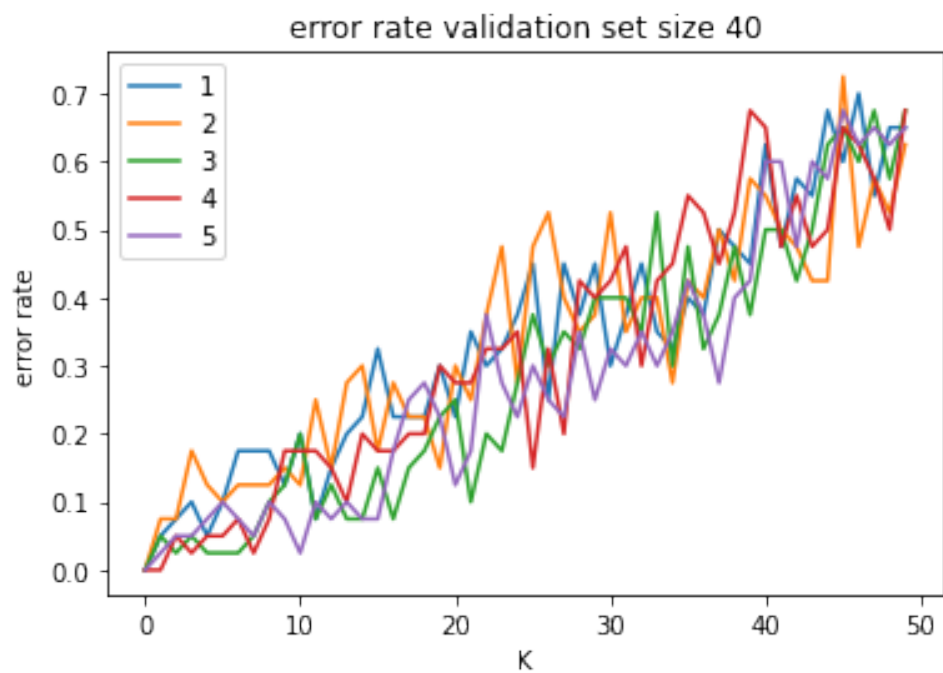


Figure 4: validation sets of size 80

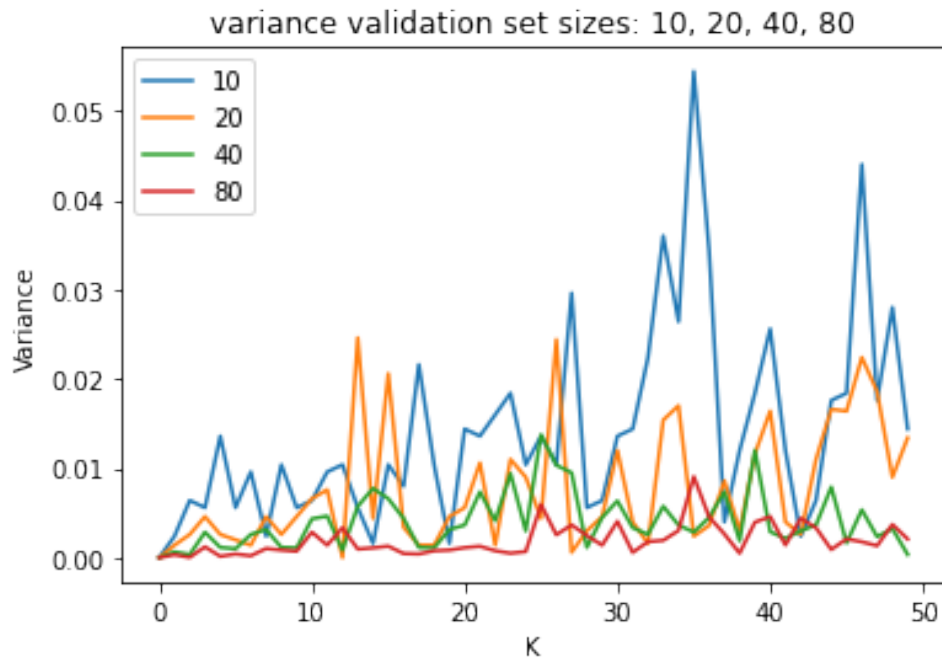
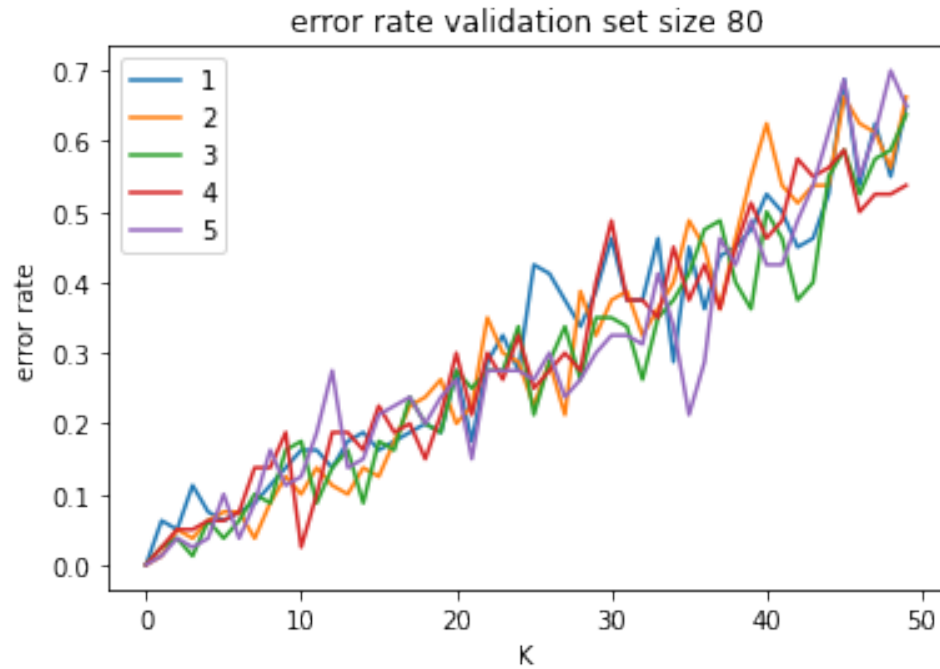


Figure 5: Variance over the 5 validation sets for each n as a function of K

- What can you say about fluctuations of the validation error as a function of n ?
I observed that as n increases the fluctuations of the validation error decrease.
- What can you say about the prediction accuracy of K-NN as a function of K ?
Once again as K increases the accuracy in prediction should decrease, but the highest

accuracy doesn't indicate the best K for K -NN, in fact in the extreme case of $K = 0$, where we take into consideration for prediction only the sample that we want to predict, we obtain an accuracy of 100% but, it is actually the worst possible model for deployment, because it won't evaluate new samples.

2.1 Task 2

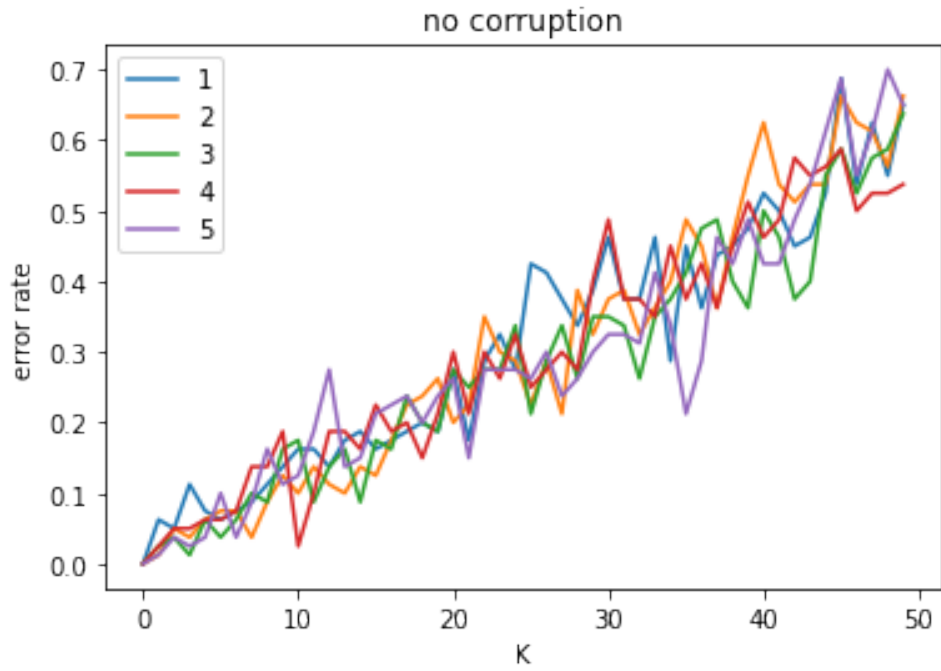


Figure 6: error rate for each validation set as a function of K, the error is computed as before, for each K: num-wrong-predictions / size-validation-set

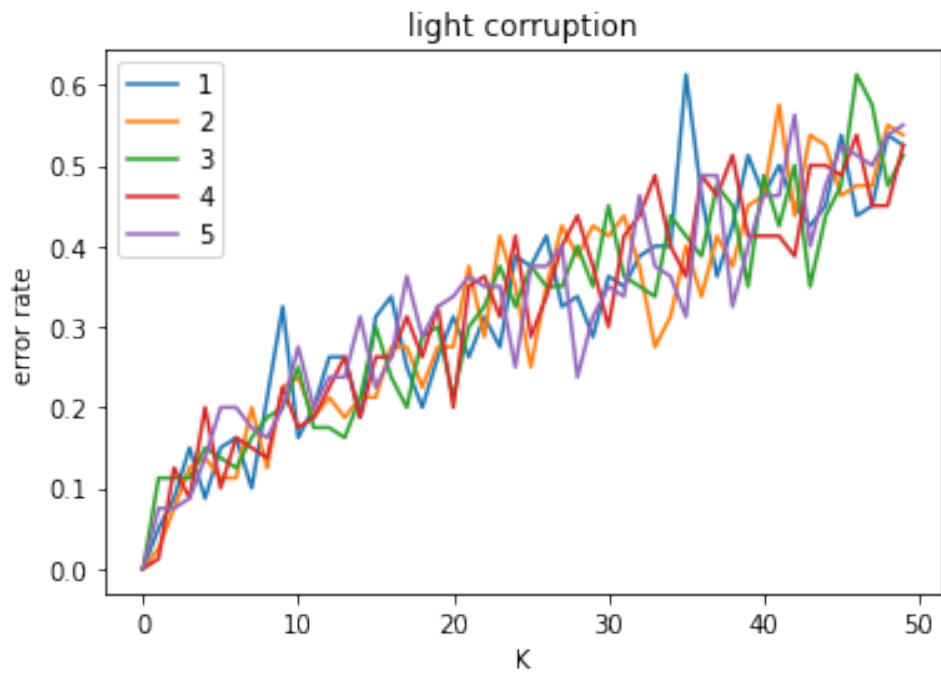


Figure 7: LIGHT CORRUPTION, error rate for each validation set as a function of K, the error is computed as before, for each K: num-wrong-predictions / size-validation-set

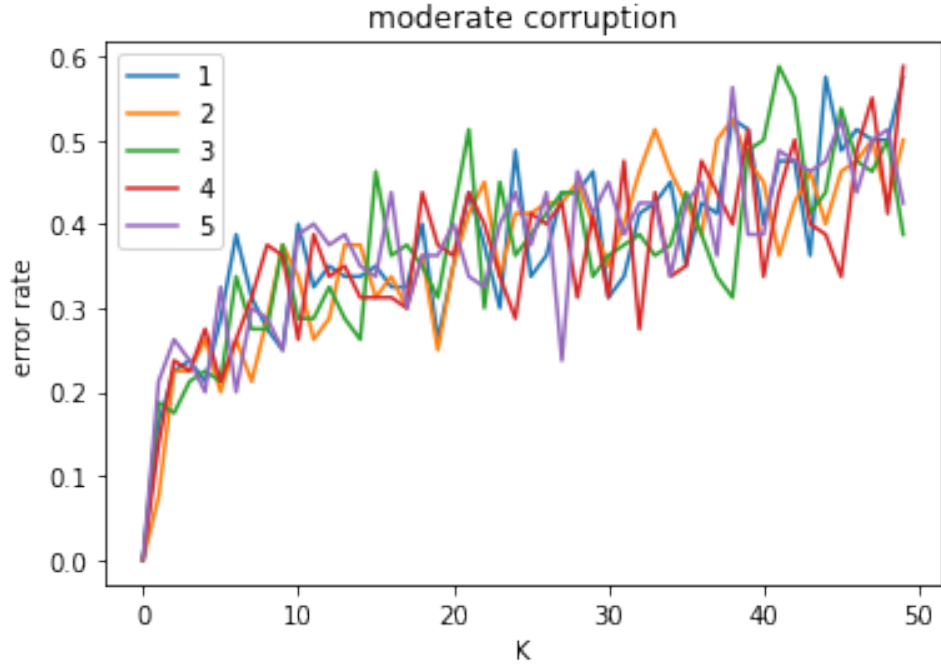


Figure 8: MODERATE CORRUPTION, error rate for each validation set as a function of K, the error is computed as before, for each K: num-wrong-predictions / size-validation-set

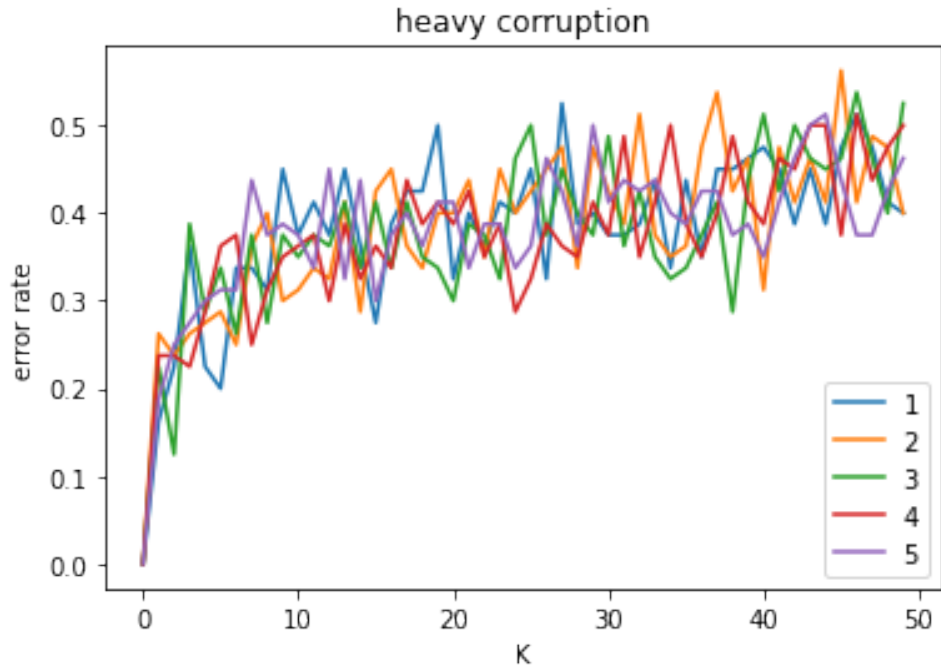


Figure 9: HEAVY CORRUPTION, error rate for each validation set as a function of K, the error is computed as before, for each K: num-wrong-predictions / size-validation-set

- Discuss how corruption magnitude influences the prediction accuracy of K-NN and the optimal value of K.

The higher the corruption, the lower tends to be the accuracy, with high corruption the algorithm tends to have worse overall performances compared to the inexistent corruption data.

With higher corruption the algorithm tend to have a balanced accuracy as K increases, I mean that after some K the algorithm doesn't really become worse as K increases, the accuracy values remain stable.

The optimal value of K tends still to be with a low value of K, however with low presence of corruption we can still find low amount of errors as K increase, for example, $k = 9$ in the first validation set of the first image of this second task.

In the high corruption data, the KNN struggle to find a good K for generalization and the errors for predictions tend to increase drastically as K increase and then start to stabilize as K increases even more

3 Linear Regression (45 points)

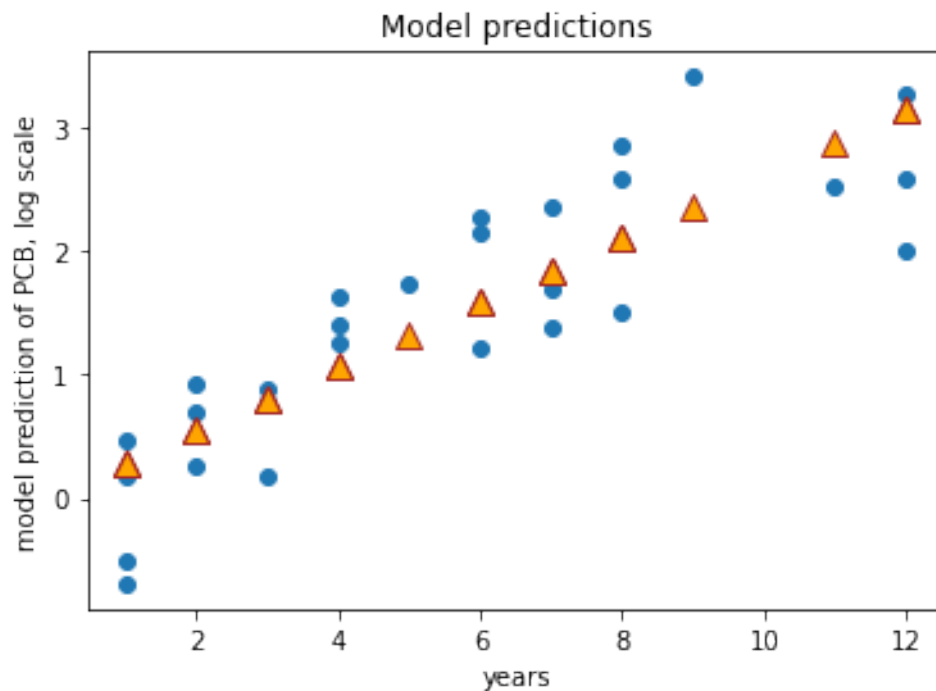


Figure 10: age vs model output, age vs log of PCB

```
#mean_squared_error
summ = 0
for i in range(len(X)):
    summ += math.pow((Y[i] - non_linear_model(X[i], a, b)) , 2)

mean_squared_error = summ/len(X)
mean_squared_error

34.83556116722035
```

Figure 11: MSE of non linear model

```
#parameters of regression model
a = np.array(w_star)[0][0]
b = np.array(w_star)[0][1]
a,b

(0.2591282395640714, 0.031472469714475815)
```

Figure 12: Parameters of regression model

4. Discuss this quantity. What does it mean if R^2 is 1 and especially if R^2 is 0? Can R^2 be negative?

if $R^2 = 1$ it means that the model can perfectly explain the data (the numerator goes to 0), in other words, how much a variable's behavior (independent variable) can explain the behavior of another variable (dependent variable).

if $R^2 = 0$ it means that most probably the regression model needs to be adjusted because through this model the independent variable cannot explain the dependent variable.

R^2 it's a statistical measure, and even if the R^2 values are high, we have no guarantees that the model once deployed will be perfect, same for low values of R^2 , we have no guarantees that the model is bad.

R^2 have values ranging from 0 to 1, thus it cannot be negative.

```
#mean_squared_error
summ = 0
for i in range(len(X)):
    summ += math.pow((Y_label[i] - non_linear_model(X[i], a, b)) , 2)

mean_squared_error = summ/len(X)
mean_squared_error

28.084390174944378
```

Figure 13: mean-squared error of model with transformed inputs

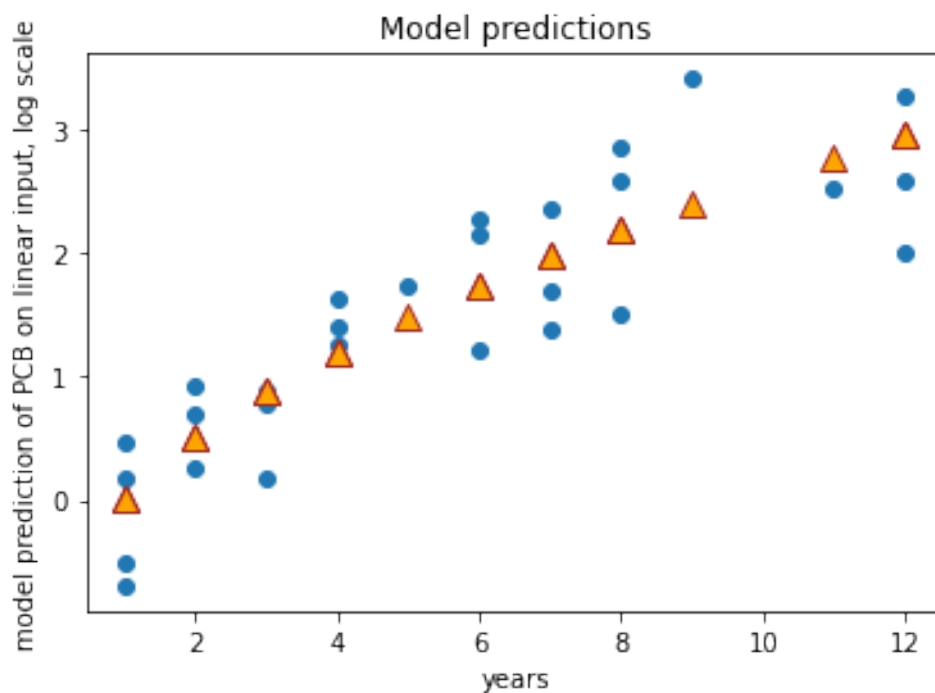


Figure 14: linear scale for years, target log-scale

R^2 of second model: 0.4816250669292409

The result shows an higher R^2 coefficient, this means that the second model should capture better the relation between the independent variable x and the dependent variable y