

Computing PAI with resampling methods:

In Linear Regression:

```
In [ ]: import numpy as np
import pandas as pd

Old Data

In [ ]: np.random.seed(123)

# Quantitative Response
Y_old = np.random.normal(loc=50, scale=10, size=500)

# Quantitative variables
X1_old = np.random.normal(loc=30, scale=25, size=500)

# Binary variables
X2_old = np.random.uniform(low=0.0, high=1.0, size=500).round()

# Multiclass categorical variables
X3_old = np.random.uniform(low=0, high=4, size=500).round() # categories: 0,1,2,3,4

New Data (with a big change in X1 distribution)

In [ ]: np.random.seed(666)

# Quantitative Response
Y_new = np.random.normal(loc=50, scale=10, size=500)

# Quantitative variables
X1_new = np.random.normal(loc=15, scale=60, size=500)

# Binary variables
X2_new = np.random.uniform(low=0.0, high=1.0, size=500).round()

# Multiclass categorical variables
X3_new = np.random.uniform(low=0, high=4, size=500).round() # categories: 0,1,2,3,4

In [ ]: df_Old = pd.DataFrame( ["Y":Y_old, "X1": X1_old, "X2": X2_old, "X3": X3_old] )
df_New = pd.DataFrame( ["Y":Y_new, "X1": X1_new, "X2": X2_new, "X3": X3_new] )

In [ ]: from matplotlib import ggplot, aes, geom_line, geom_point, geom_histogram, geom_bar, geom_boxplot, scale_y_continuous, scale_x_continuous, labs, after_stat, geom_vline, scale_color_1
from mizani.formatters import percent_format

In [ ]: import array as arr

df_Old_New = pd.concat([df_Old, df_New])

repeat_Old = ["Old Data"]*len(df_Old)
repeat_New = ["New Data"]*len(df_New)

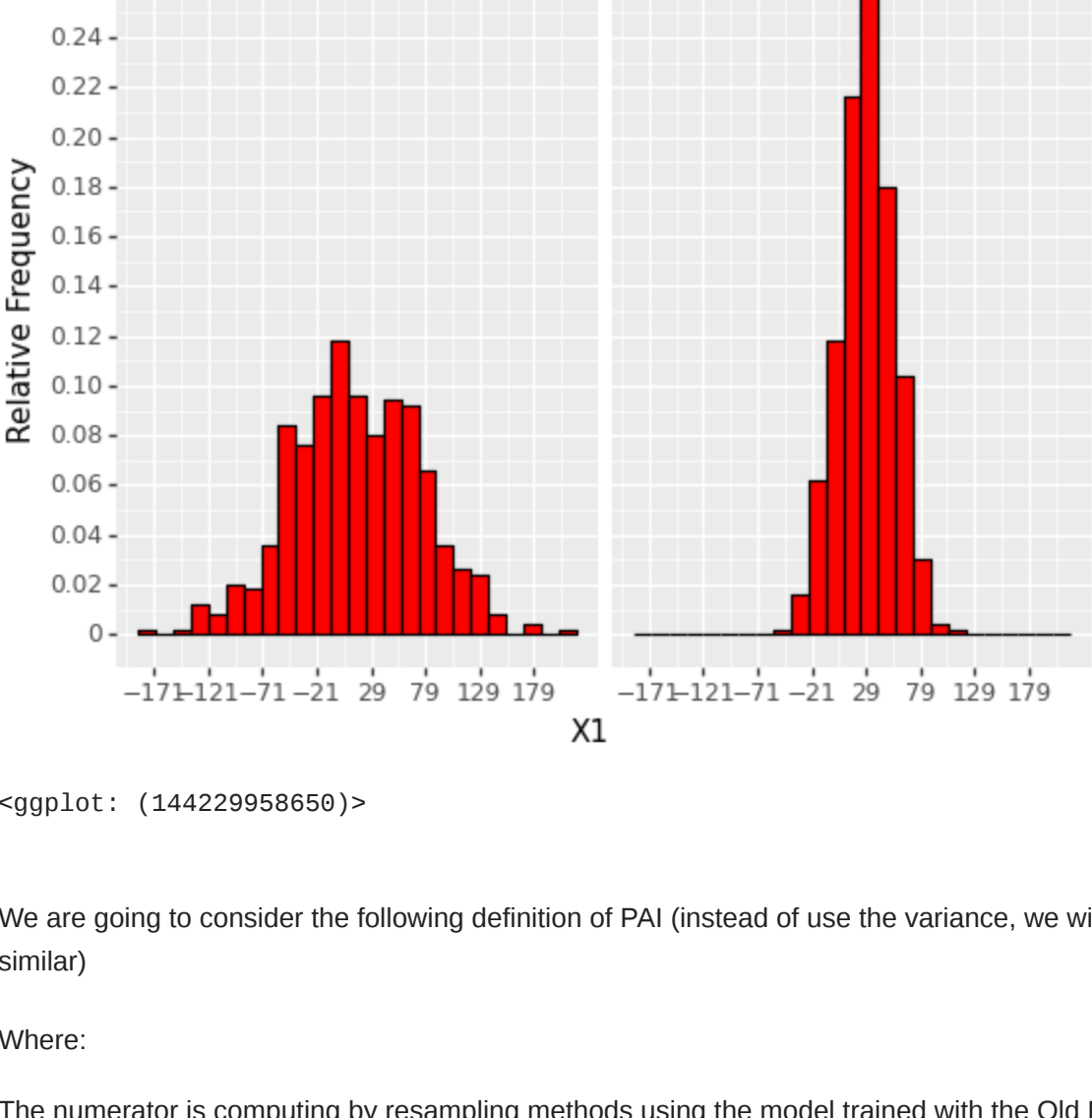
df_repeat_New = pd.DataFrame( {"group": repeat_New} )
df_repeat_Old = pd.DataFrame( {"group": repeat_Old} )

groups = pd.concat([df_repeat_Old, df_repeat_New])

df_Old_New_groups = pd.concat([df_Old_New, groups], axis=1 )
df_Old_New_groups

Out[ ]:
      Y      X1  X2  X3  group
0  39.143694  48.800842  1.0  1.0  Old Data
1  59.973454  31.741019  0.0  3.0  Old Data
2  52.829785  23.363859  0.0  4.0  Old Data
3  34.927053  53.239612  1.0  1.0  Old Data
4  44.213997  61.520532  1.0  3.0  Old Data
...
495  56.756956  12.495333  0.0  3.0  New Data
496  36.846327  75.862895  1.0  1.0  New Data
497  36.528625  -21.430239  0.0  0.0  New Data
498  62.143112  -56.088634  0.0  2.0  New Data
499  53.433602  10.075296  0.0  3.0  New Data
1000 rows x 5 columns

In [ ]: (
ggplot( df_Old_New_groups )
+ aes(x='X1', y = after_stat('width*density'))
+ geom_histogram(fill="red", color="black", bins = 25)
+ labs(x = "X1", y = "Relative Frequency")
+ scale_x_continuous( breaks = range(int(df_Old_New_groups['X1'].min()), int(df_Old_New_groups['X1'].max()), 50) )
+ scale_y_continuous( breaks = np.arange(0, 0.5, 0.02) )
+ facet_wrap('group')
)
```



Out[ ]: <ggplot: (144229958656)>

We are going to consider the following definition of PAI (instead of use the variance, we will use the standard deviation. if we would consider the PAI definition with the variance, the process to compute it would have been very similar)

Where:

The numerator is computing by resampling methods using the model trained with the Old Data (for the response and the predictors) but predicting the response variable using the New Data for the predictors.

The denominator is computing by resampling methods using the model trained Old Data (for the response and the predictors) and also predicting the response variable using the Old Data for the predictors.

$$PAI = \frac{\frac{1}{N} \sum_{i \in NewData} \widehat{Var}(\hat{y}_i)}{\frac{1}{n} \sum_{i \in OldData} \widehat{Var}(\hat{y}_i)}$$

So, PAI is the quotient between the mean variance of the response predictions, using the model trained with the Old data but using the New data for the predictors to get the response predictions and the mean variance of the response predictions, using the model trained with the Old data and also using the Old data for the predictors to get the response predictions

```
In [ ]: def varcharProcessing(X, varchar_process = "dummy_dropFirst"):
    dtypes = X.dtypes

    if varchar_process == "drop":
        X = X.drop(columns = dtypes[dtypes == np.object].index.tolist())

    elif varchar_process == "dummy":
        X = pd.get_dummies(X, drop_first=False)

    elif varchar_process == "dummy_dropFirst":
        X = pd.get_dummies(X, drop_first=True)

    else:
        X = pd.get_dummies(X, drop_first=True)

    X["intercept"] = 1
    cols = X.columns.tolist()
    cols = cols[1:] + cols[0:1]
    X = X[cols]

    return X

In [ ]: df_Old['X2'] = df_Old['X2'].astype('category')
df_Old['X3'] = df_Old['X3'].astype('category')

In [ ]: df_Old

Out[ ]:
      Y      X1  X2  X3
0  39.143694  48.800842  1.0  1.0
1  59.973454  31.741019  0.0  3.0
2  52.829785  23.363859  0.0  4.0
3  34.927053  53.239612  1.0  1.0
4  44.213997  61.520532  1.0  3.0
...
495  56.756956  12.495333  0.0  3.0
496  36.846327  75.862895  1.0  1.0
497  36.528625  -21.430239  0.0  0.0
498  62.143112  -56.796594  0.0  1.0
499  48.579493  2.142439  1.0  3.0
500 rows x 4 columns

In [ ]: from sklearn.linear_model import LinearRegression

In [ ]: B=100

y_predictions_Old_Data = np.zeros((B, len(df_Old)))

for i in range(0, B):

    df_Old_BOOT_SAMPLE = df_Old.sample(n=len(df_Old), random_state=i, replace=True) # i-th boot sample

    X_Old_BOOT_Sample = df_Old_BOOT_SAMPLE[['X1', 'X2', 'X3']]

    y_Old_BOOT_Sample = df_Old_BOOT_SAMPLE['Y']

    X_Old_BOOT_Sample = varcharProcessing(X_Old_BOOT_Sample, varchar_process = "dummy_dropFirst")

    # We train the model with i-th boot sample of the Old Data:
    Model_train_Old_data = LinearRegression().fit(X_Old_BOOT_Sample, y_Old_BOOT_Sample)

    # y predictions using Model_train_Old_data with the original Old Data for the predictors

    X_Old = df_Old[['X1', 'X2', 'X3']]

    X_Old = varcharProcessing(X_Old, varchar_process = "dummy_dropFirst")

    # La idea es que con cada iteracion el modelo cambia (los parametros) ya que ha sido entrenado con diferentes data set (las muestras bootstrap)
    # Pero las predicciones se hacen usando siempre el mismo data set (Old data), para asi asegurar que Vhat(y_i) es siempre la prediccio de la respueta
    # para el i-esimo individuo cambiara porque al re-entrenar el modelo con las distintas muestras boot cambian los parametros de este, pero no cambia
    # el vector x_i de valores de los predictores del individuo i, con los que tambien se genera la prediccion

    # Por tanto si hubiese mucha variabilidad entre los valores obtenidos de Vhat(y_i) = Vhat(beta) x x_i esto se deberia no a cambios en x_i (puesto que
    # en este problema no cambia), si no a cambios en Vhat(beta) debidos a las diferentes muestras boot usadas para entrenar el modelo.
    # Si en lugar del modelo de regresion lineal multiple usasemos otro de la misma indole (basicamente un modelo que pueda ser entrenado con unos datos
    # y pueda realizar predicciones sobre una respuesta usando otros datos), la idea seria la misma.

    y_predictions_Old_Data[i, :] = Model_train_Old_data.predict(X_Old)
```

The  $(i, k)$  element of the  $B \times n$  matrix  $y\_predictions\_Old\_Data$  is  $\hat{y}_k$  (the  $Y$  estimation for the  $k$ -th individual of the sample) when the model is trained with the  $i$ -th boot sample of Old\_Data\_Set.

In other words:

$$y\_predictions\_Old\_Data = \begin{pmatrix} \hat{y}_1^1 & \hat{y}_1^2 & \dots & \hat{y}_1^n \\ \hat{y}_2^1 & \hat{y}_2^2 & \dots & \hat{y}_2^n \\ \vdots & \vdots & \ddots & \vdots \\ \hat{y}_B^1 & \hat{y}_B^2 & \dots & \hat{y}_B^n \end{pmatrix}$$

Where:

$\hat{y}_k^i$  is the response prediction for the  $k$ -th element of the original Old\_Data\_Set, obtained using the model trained with the  $i$ -th boot sample of Old\_Data\_Set

$n = \text{len}(\text{Old\_Data\_Set})$

$B = n^{\circ}$  of boot samples

```
In [ ]: y_predictions_Old_Data

Out[ ]: array([[49.42508598, 50.11472721, 48.97935105, ..., 50.27502977,
        51.17601889, 47.4028062],
       [48.87265439, 50.14695669, ..., 49.86794534, ..., 49.2737147,
        49.40472178, ..., 49.47620932],
       [48.66018825, 50.95861995, 48.82043962, ..., 50.59214347,
        48.66699903, ..., 49.84584848],
       ...,
       [49.19061918, 51.17188553, 50.27997571, ..., 49.42699339,
        49.36788126, 51.11558389],
       [48.4354486, 50.99945593, 50.51106095, ..., 50.99576546,
        47.97509545, ..., 51.38434712],
       [49.80485734, 50.35848299, 50.47602223, ..., 49.57213947,
        50.16686803, ..., 49.0074667]])

We compute the standard deviation of each column of the matrix , and we get an estimation of Var(y_i) for i = 1, ..., n:

So, the i-th value of the following vector is

Var(y_i)

In [ ]: # compute the variance by cols in an array

y_predictions_Old_Data.var(axis=0)

Out[ ]: array([2.3736694, 1.83531851, 1.49999452, 2.02032802, 2.03471412,
        2.59662433, 2.32653866, 2.10161238, 2.51028022, 2.48036944,
        1.62799249, 1.81614304, 2.08080234, 2.46520527, 1.90376617,
        1.39392336, 2.15432889, 2.81813501, 1.87711112, 1.85349095,
        1.52218726, 2.01087883, 2.40951067, 2.01183162, 1.66672321,
        1.94795425, 1.9963901, 2.14494738, 2.13831145, 2.2681182,
        1.93510707, 2.03019929, 1.80845261, 1.88852273, 1.72965494,
        2.29470579, 2.2313705, 2.21600616, 2.52404316, 2.00561608,
        2.22137462, 1.90636041, 1.95101852, 2.20370171, 1.90985279,
        2.45884646, 1.65150189, 1.86239492, 2.05955444, 2.08133178,
        2.09107018, 2.09175392, 2.06124553, 2.37485157, 1.88746683,
        2.40074097, 2.03045088, 2.8090999, 2.41083172, 2.5622219,
        1.13638402, 1.41088064, 1.94105678, 2.08061339, 2.22737029,
        1.7952498, 2.36554675, 1.88416718, 2.11025988, 2.34593621,
        2.43904217, 2.30040853, 1.673633, 2.058072, 1.80890906,
        1.84127261, 1.79411597, 2.04758911, 1.44696274, 2.19625316,
        2.15356721, 2.43663391, 1.84033922, 2.06567795, 1.90737313,
        2.16427552, 2.0374608, 2.13522344, 2.32175074, 1.57980151,
        1.72140516, 2.01295333, 2.30454453, 2.15839802, 2.03020604,
        2.28739972, 1.59712926, 2.37949883, 2.34071699, 1.86147695,
        1.32172721, 2.14539257, 1.94091567, 1.78402293, 2.46771354,
        1.99270509, 1.53240014, 2.3358246, 2.07521755, 2.55004932,
        1.65392279, 2.12119525, 2.28159141, 1.65037548, 2.31523818,
        2.13680489, 1.87511976, 1.7570361, 2.0454583, 1.53014661,
        2.13877488, 2.54344998, 1.58554634, 2.33854072, 1.67445211,
        2.02464654, 2.3954862, 2.05940408, 2.00809808, 2.44392023,
        2.61004032, 1.92477361, 2.53486839, 2.43077422, 2.31594316,
        2.23053751, 1.50278065, 1.65479316, 2.53723677, 1.76171706,
        1.8731612, 1.88706053, 1.97285293, 1.54857729, 1.94464846,
        1.82073409, 1.35377844, 1.80424385, 1.87433806, 2.00780601,
        1.9845866, 2.33536142, 1.99318704, 2.0761117, 1.76830416,
        2.22150964, 1.70504076, 2.62562896, 1.50452379, 1.75079702,
        2.10999094, 2.03692916, 2.23573177, 2.44186908, 1.97346309,
        1.13635514, 2.25379509, 2.10235741, 1.85254657, 2.14479637,
        2.0238252, 1.99094124, 2.27144346, 1.93958067, 2.04749002,
        1.86235156, 2.04260279, 2.41948454, 1.92063609, 2.12542531,
        2.0462215, 2.13274718, 1.9406355, 2.28056378, 2.02460662,
        2.09574509, 2.32417108, 2.23425804, 1.95356526, 1.70782206,
        2.7923403, 2.14299365, 2.00925918, 1.73751139, 2.13069998,
        2.04781348, 2.06361869, 1.72952094, 1.74286309, 1.69793902,
        1.95868322, 2.53518962, 1.34778637, 1.95703641, 2.12084068,
        2.0163571, 2.10727252, 1.91302864, 1.83466435, 2.16943251,
        2.14040875, 2.3763879, 2.48511886, 1.89130727, 2.40598257,
        1.91444116, 2.47567096, 1.7401907, 1.8994908, 2.06063375,
        2.06548666, 2.39147523, 2.13028001, 2.2895939, 2.11337613,
        2.0628776, 1.67736474, 1.6982976, 2.12359273, 1.45225914,
        2.05293631, 2.09933561, 2.00132161, 2.62721961, 1.98032174,
        1.7617149, 2.26029724, 2.69430207, 1.90852748, 1.82384422,
        1.87950783, 2.55704616, 2.0702351, 1.76886125, 1.74954093,
        2.69517098, 2.03615693, 2.60928485, 2.28406125, 2.17860311,
        1.84848552, 2.43062639, 1.997738, 1.98180909, 1.97346309,
        1.79334553, 1.93679284, 1.80742437, 1.74573569, 1.99522728,
        1.9001478, 2.27293834, 1.98647476, 1.8150405, 1.90720433,
        2.94197196, 2.18354261, 1.78138715, 1.67709711, 2.02725808,
        2.44078022, 2.03094014, 2.4509071, 2.14778403, 2.2515794,
        1.4898767, 2.01469702, 2.38008713, 2.07756995, 2.26563375,
        2.12764036, 1.8845882, 1.97192778, 2.30337225, 2.1012967,
        2.1682583, 2.1454583, 1.73591251, 2.0886238, 1.9728476,
        1.7450414, 1.56012745, 2.14078333, 1.70893376, 2.05642478,
        2.6825844, 1.79935841, 2.17114618, 1.46349202, 1.65274671,
        2.16626081, 2.09234108, 1.93827301, 1.96406104, 2.17485225,
        1.9573505, 2.00590074, 2.16019626, 2.38724007, 2.22956641,
        2.26797791, 2.50453138, 2.18175008, 1.78979546, 1.50465655,
        1.87395415, 1.9592566, 2.48935403, 2.26473888, 1.63286279,
        2.21744857, 2.13937529, 2.06175399, 2.01967007, 2.06383905,
        2.25060181, 2.14090625, 2.19456829, 2.463169743, 2.16412957,
        2.20508091, 2.32417108, 2.23425804, 1.95356526, 1.70782206,
        2.7923403, 2.14299365, 2.00925918, 1.73751139, 2.13069998,
        2.04781348, 2.06361869, 1.72952094, 1.74286309, 1.69793902,
        1.95868322, 2.53518962, 1.34778637, 1.95703641, 2.12084068,
        2.0163571, 2.10727252, 1.91302864, 1.83466435, 2.16943251,
        2.14040875, 2.3763879, 2.48511886, 1.89130727, 2.40598257,
        1.91444116, 2.47567096, 1.7401907, 1.8994908, 2.06063375,
        2.06548666, 2.39147523, 2.13028001, 2.2895939, 2.11337613,
        2.0628776, 1.67736474, 1.6982976, 2.12359273, 1.45225914,
        2.05293631, 2.09933561, 2.00132161, 2.62721961, 1.98032174,
        1.7617149, 2.26029724, 2.69430207, 1.90852748, 1.82384422,
        1.87950783, 2.55704616, 2.0702351, 1.76886125, 1.74954093,
        2.69517098, 2.03615693, 2.60928485, 2.28406125, 2.17860311,
        1.84848552, 2.43062639, 1.997738, 1.98180909, 1.97346309,
        1.79334553, 1.93679284, 1.80742437, 1.74573569, 1.99522728,
        1.9001478, 2.27293834, 1.98647476, 1.8150405, 1.90720433,
        2.94197196, 2.18354261, 1.78138715, 1.67709711, 2.02725808,
        2.44078022, 2.03094014, 2.4509071, 2.14778403, 2.2515794,
        1.4898767, 2.01469702, 2.38008713, 2.07756995, 2.26563375,
        2.12764036, 1.8845882, 1.97192778, 2.30337225, 2.1012967,
        2.1682583, 2.1454583, 1.73591251, 2.0886238, 1.9728476,
        1.7450414, 1.56012745, 2.14078333, 1.70893376, 2.05642478,
        2.6825844, 1.79935841, 2.17114618, 1.46349202, 1.65274671,
        2.16626081, 2.09234108, 1.93827301, 1.96406104, 2.17485225,
        1.9573505, 2.00590074, 2.16019626, 2.38724007, 2.22956641,
        2.26797791, 2.50453138, 2.18175008, 1.78979546, 1.50465655,
        1.87395415, 1.9592566, 2.48935403, 2.26473888, 1.63286279,
        2.21744857, 2.13937529, 2.06175399, 2.01967007, 2.06383905,
        2.25060181, 2.14090625, 2.19456829, 2.463169743, 2.16412957,
        2.20508091, 2.32417108, 2.23425804, 1.95356526, 1.70782206,
        2.7923403, 2.14299365, 2.00925918, 1.73751139, 2.13069998,
        2.04781348, 2.06361869, 1.72952094, 1.74286309, 1.69793902,
        1.95868322, 2.53518962, 1.34778637, 1.95703641, 2.12084068,
        2.0163571, 2.10727252, 1.91302864, 1.83466435, 2.16943251,
        2.14040875, 2.3763879, 2.48511886, 1.89130727, 2.40598257,
        1.91444116, 2.47567096, 1.7401907, 1.8994908, 2.06063375,
        2.06548666, 2.39147523, 2.13028001, 2.2895939, 2.11337613,
        2.0628776, 1.67736474, 1.6982976, 2.12359273, 1.45225914,
        2.05293631, 2.09933561, 2.00132161, 2.62721961, 1.98032174,
        1.7617149, 2.26029724, 2.69430207, 1.90852748, 1.82384422,
        1.87950783, 2.55704616, 2.0702351, 1.76886125, 1.74954093,
        2.69517098, 2.03615693, 2.60928485, 2.28406125, 2.17860311,
        1.84848552, 2.43062639, 1.997738, 1.98180909, 1.97346309,
        1.79334553, 1.93679284, 1.80742437, 1.74573569, 1.99522728,
        1.9001478, 2.27293834, 1.98647476, 1.8150405, 1.90720433,
        2.94197196, 2.18354261, 1.78138715, 1.67709711, 2.02725808,
        2.44078022, 2.03094014, 2.4509071, 2.14778403, 2.2515794,
        1.4898767, 2.01469702, 2.38008713, 2.07756995, 2.26563375,
        2.12764036, 1.8845882, 1.97192778, 2.30337225, 2.1012967,
        2.1682583, 2.1454583, 1.73591251, 2.0886238, 1.9728476,
        1.7450414, 1.56012745, 2.14078333, 1.70893376, 2.05642478,
        2.6825844, 1.79935841, 2.17114618, 1.46349202, 1.65274671,
        2.16626081, 2.09234108, 1.93827301, 1.96406104, 2.17485225,
        1.9573505, 2.00590074, 2.16019626, 2.38724007, 2.22956641,
        2.26797791, 2.50453138, 2.18175008, 1.78979546, 1.50465655,
        1.87395415, 1.9592566, 2.48935403, 2.26473888, 1.63286279,
        2.21744857, 2.13937529, 2.06175399, 2.01967007, 2.06383905,
        2.25060181, 2.14090625, 2.19456829, 2.463169743, 2.16412957,
        2.20508091, 2.32417108, 2.23425804, 1.95356526, 1.70782206,
        2.7923403, 2.14299365, 2.00925918, 1.73751139, 2.13069998,
        2.04781348, 2.06361869, 1.72952094, 1.74286309, 1.69793902,
        1.95868322, 2.53518962, 1.34778637, 1.95703641, 2.12084068,
        2.0163571, 2.10727252, 1.91302864, 1.83466435, 2.16943251,
        2.14040875, 2.3763879, 2.48511886, 1.89130727, 2.40598257,
        1.91444116, 2.47567096, 1.7401907, 1.8994908, 2.06063375,
        2.06548666, 2.39147523, 2.13028001, 2.2895939, 2.11337613,
        2.0628776, 1.67736474, 1.6982976, 2.12359273, 1.45225914,
        2.05293631, 2.09933561, 2.00132161, 2.62721961, 1.98032174,
        1.7617149, 2.26029724, 2.69430207, 1.90852748, 1.82384422,
        1.87950783, 2.55704616, 2.0702351, 1.76886125, 1.74954093,
        2.69517098, 2.03615693, 2.60928485, 2.28406125, 2.17860311,
       
```