

idealista
FICHeros



Fotos

html

kaggle

(CSV)

Fotocasa ~ IDEALISTA

CATASTRO
html

UNSTRUCTURED
DATA

Id

Area

BARRIO

Extensión ...

Casa 1

80

MORALEJA

✓

1.

1.

1.

Varios

Obj

1000000

200 000

STRUCTURED
DATA

```
# Data Preparation
```

```
x = df[['TotalSF']] # pandas DataFrame ←
```

```
y = df["SalePrice"] # pandas Series
```

DATAFRAME

objektuo = SERIES



$$y = f(x)$$

REG. LINEAL Characterísticas

$$y = a x_1 + b x_2 + c x_3$$

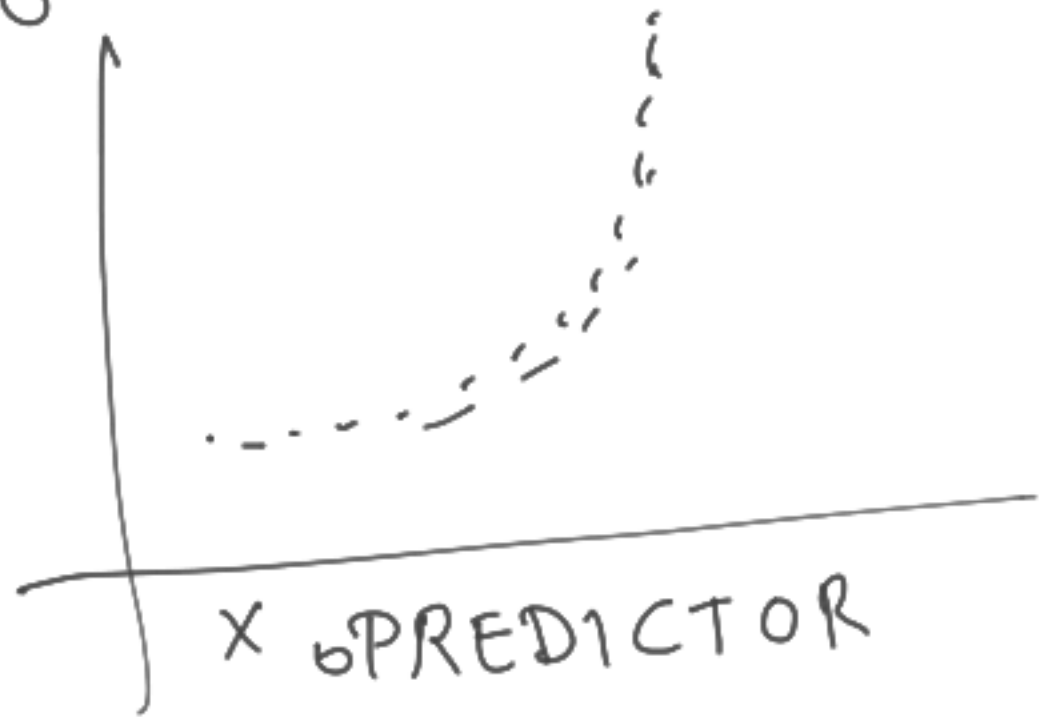
a, b, c

reg. fit(X, y) \rightarrow best a, b, c min OLS

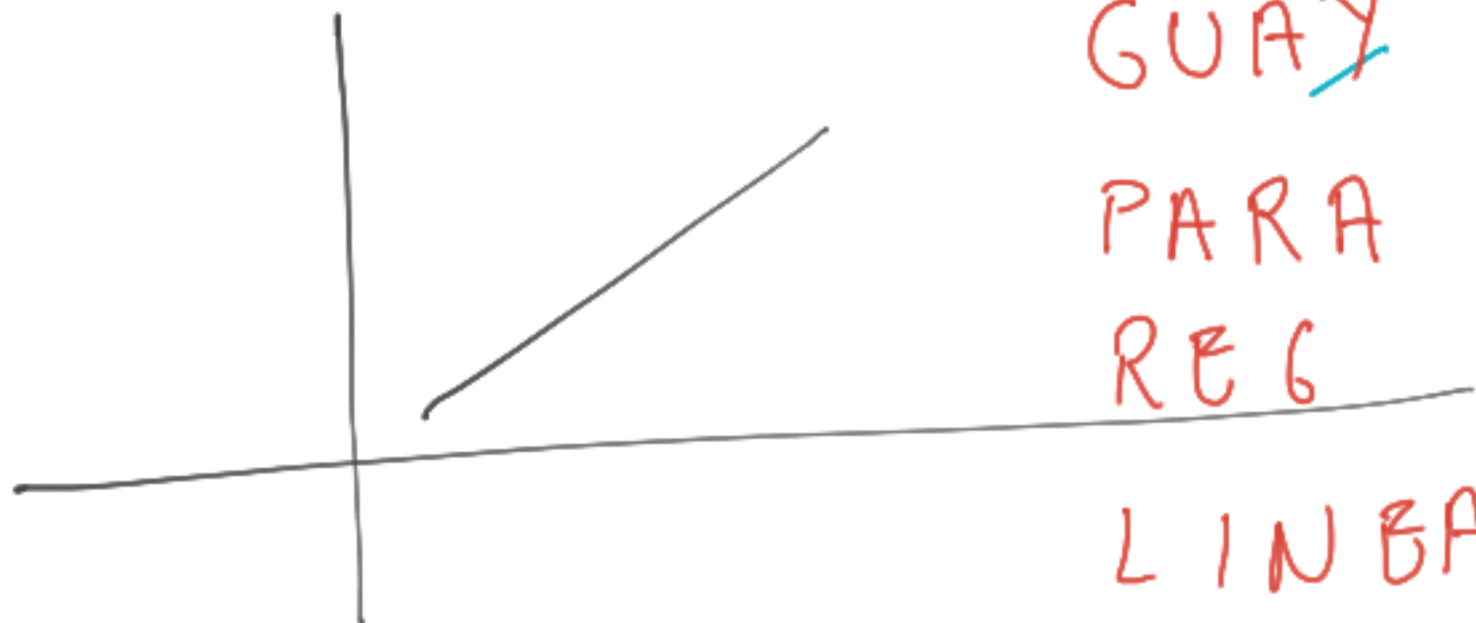
$$J = \min \sum (y - \hat{y})^2$$

$$J_2 = \min \sum |y - \hat{y}|$$

objetivo ANTES

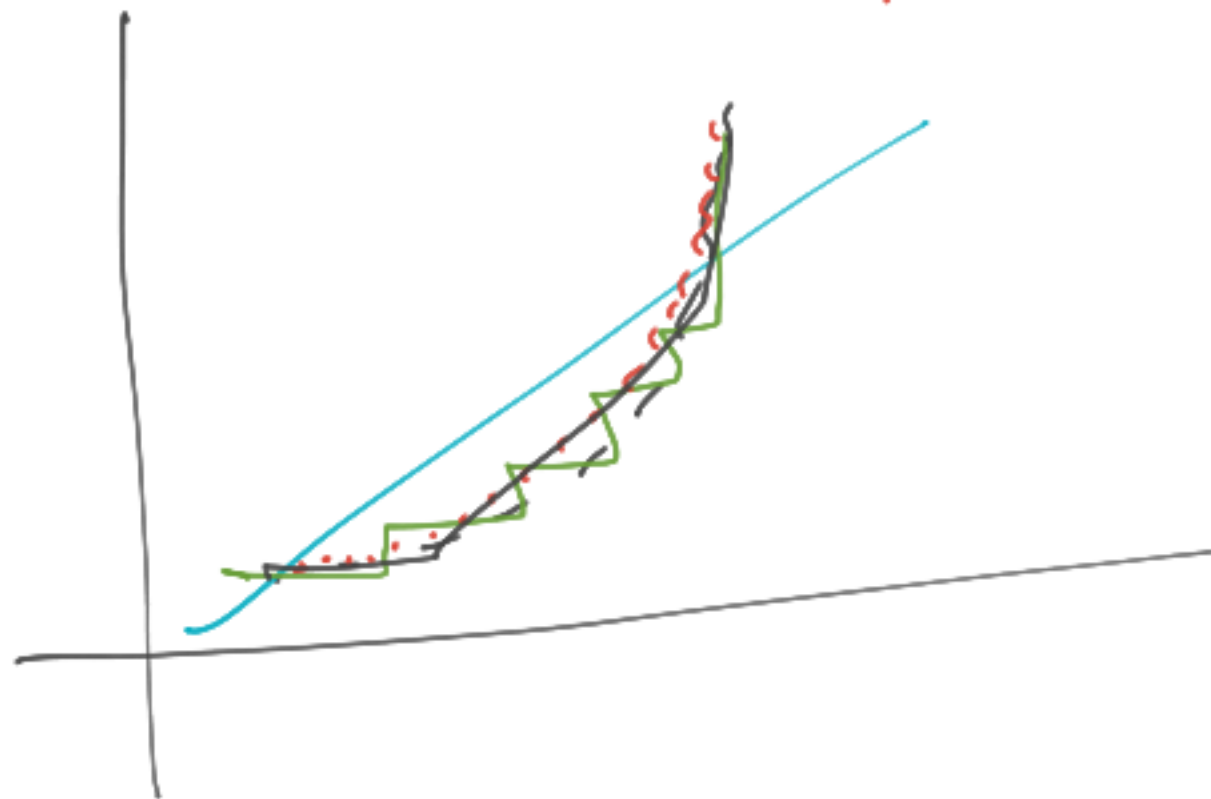


↓ log



GUAY
PARA
REG
LINEAL

DESPUES



cat.

A

B

A

B

C



cat = A

1

0

1

0

0

cat = B

0

1

0

1

0

cat = C

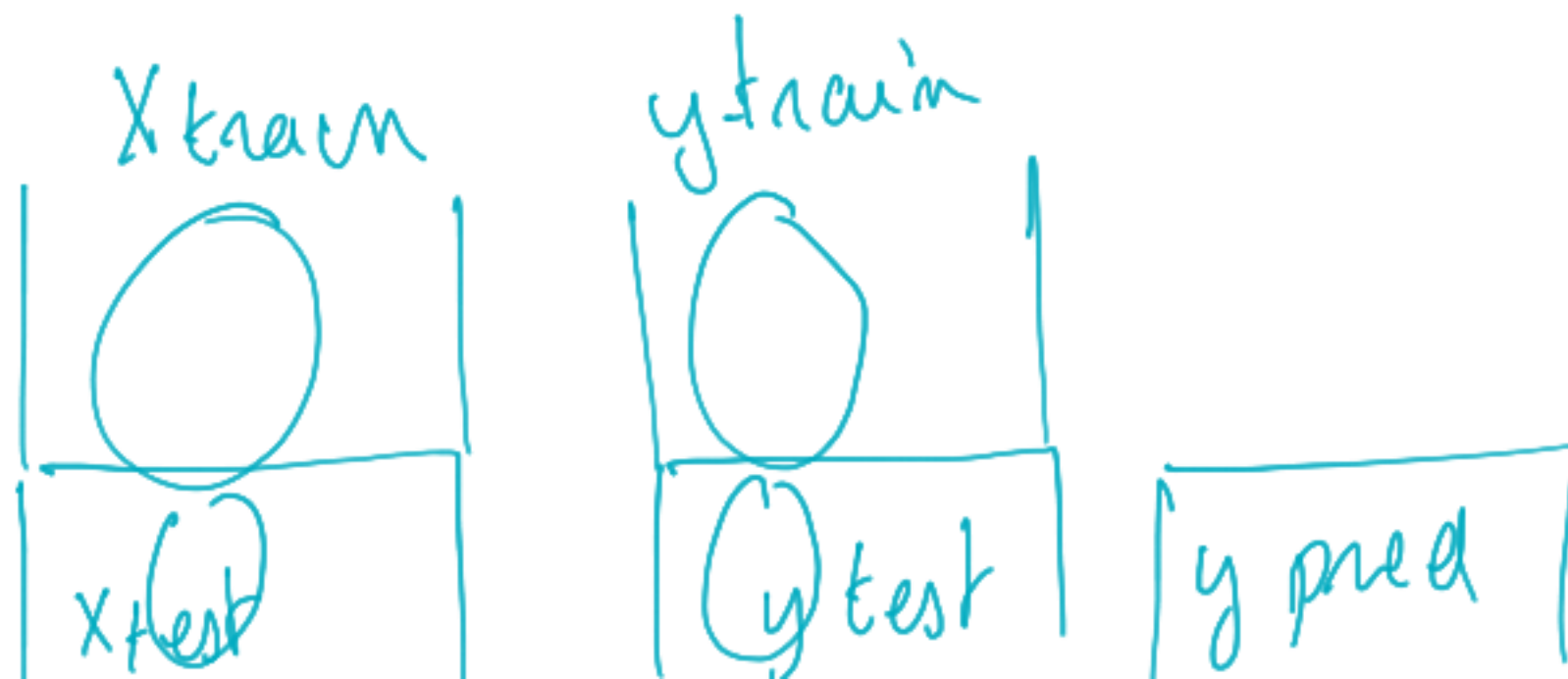
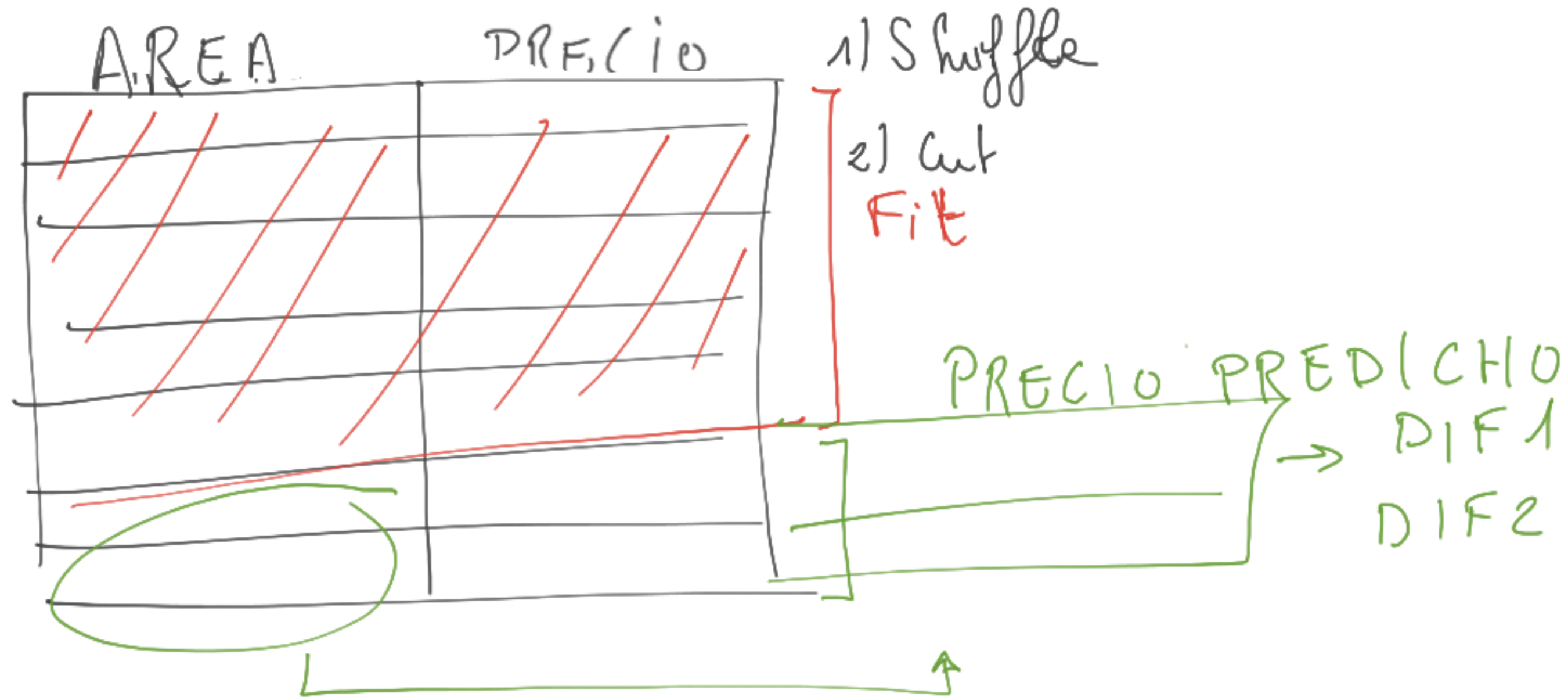
0

0

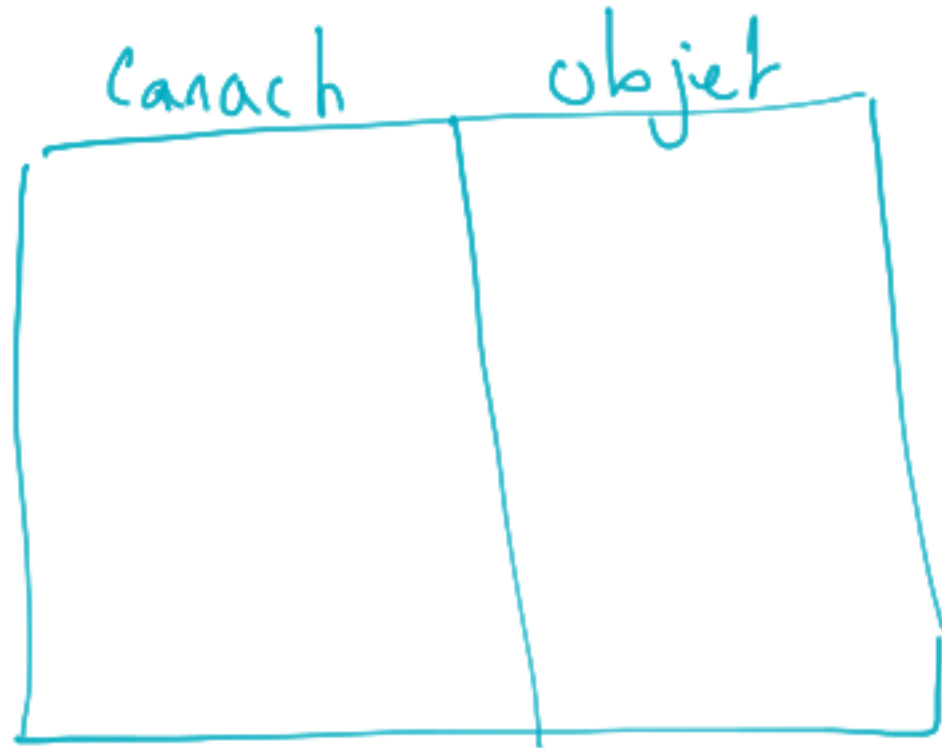
0

0

1



K nearest neighbors



No Fit X



1) Buscar K elementos + próximos + parecidos

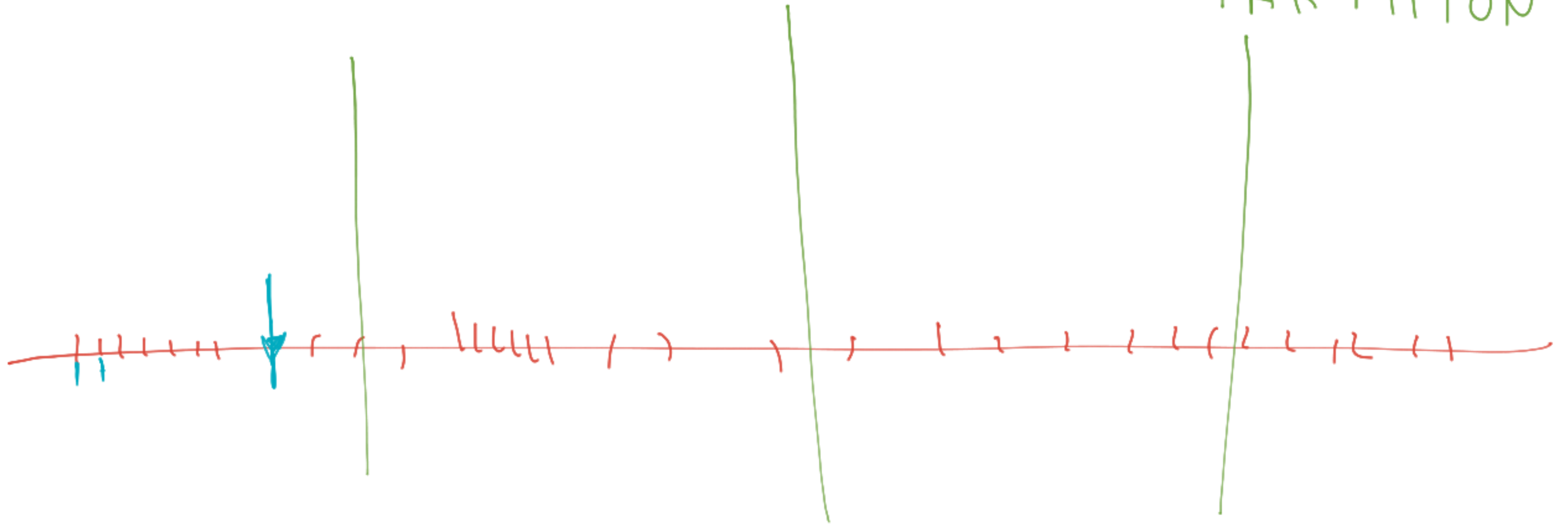
3 veijos

1 (TRAIN \rightarrow Se preino)



1 TRAIN

PARTITION



```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.20, random_state=42)
```

X _{train}	y _{train}
X _{test}	y _{test}

LIN REG

```
from sklearn.linear_model import LinearRegression
```

```
reg = LinearRegression()
```

```
reg.fit(X_train, y_train)
```

PRED TEST

```
y_pred = reg.predict(X_test)
```

```
from sklearn.metrics import mean_absolute_error
```

```
mean_absolute_error(y_test, y_pred)
```

```
from sklearn.metrics import mean_squared_error
```

RMSE: Lin Reg

```
np.sqrt(mean_squared_error(y_test, y_pred))
```

K Nei

```
from sklearn.neighbors import KNeighborsRegressor
```

```
regk = KNeighborsRegressor(n_neighbors=5)
```

```
regk.fit(X_train, y_train)
```

↓ PRED TEST

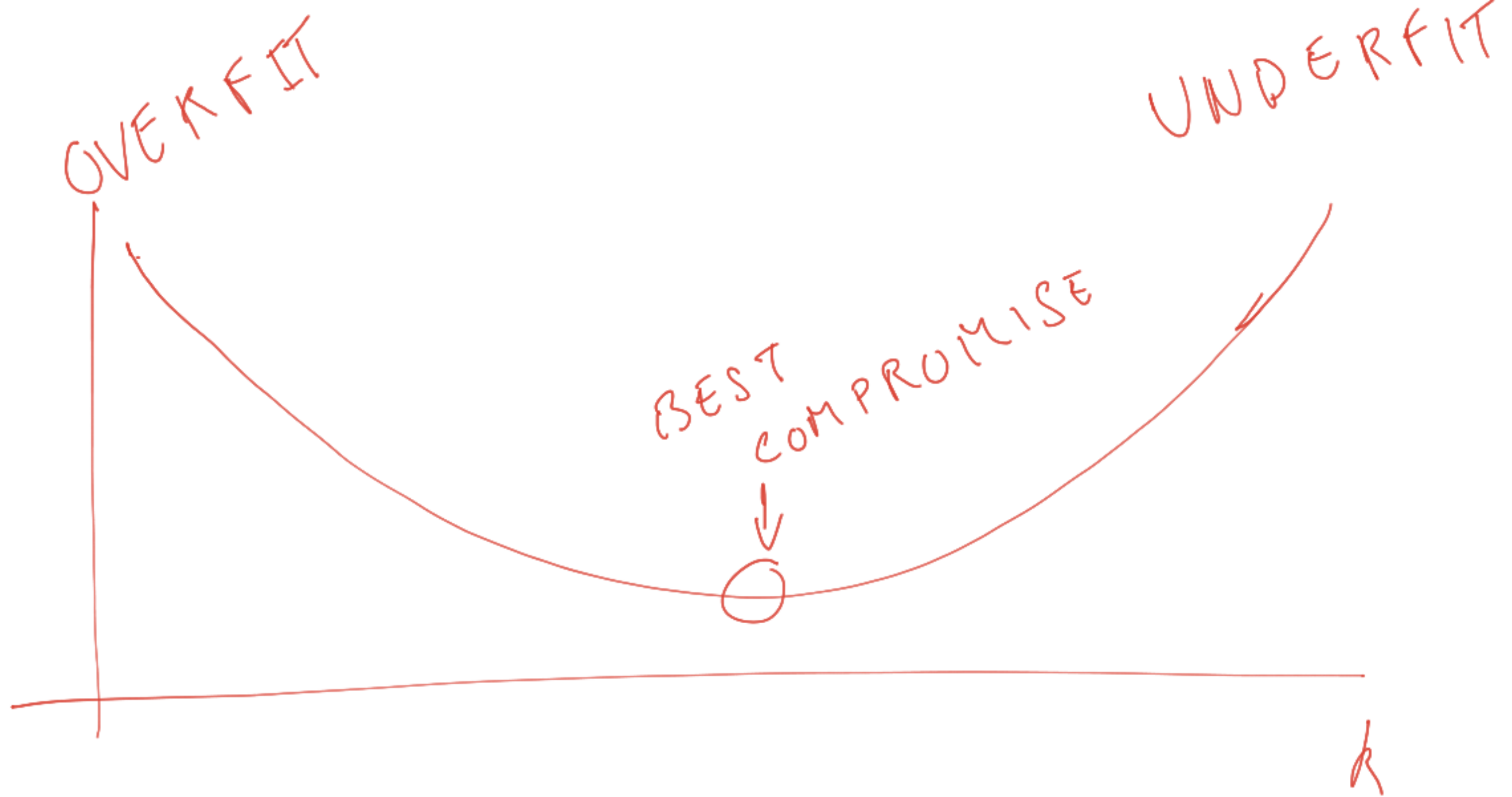
```
y_predk = regk.predict(X_test)
```

```
mean_absolute_error(y_test, y_predk)
```

RMSE: k Nei

```
np.sqrt(mean_squared_error(y_test, y_predk))
```

TEST		
	TEST	
TRAIN	TRAIN	TEST
Iter 1	Iter 2	Iter 3



```
from sklearn.model_selection import cross_val_score
```

```
# We calculate the metric for several subsets (determine by cv)
```

```
# With cv=5, we will have 5 results from 5 training/test
```

```
cross_val_score(reg,X,y,cv=5,scoring="neg_mean_squared_error")
```

← Evalua 1 modelo

```
from sklearn.model_selection import GridSearchCV
```

```
from sklearn.neighbors import KNeighborsRegressor
```

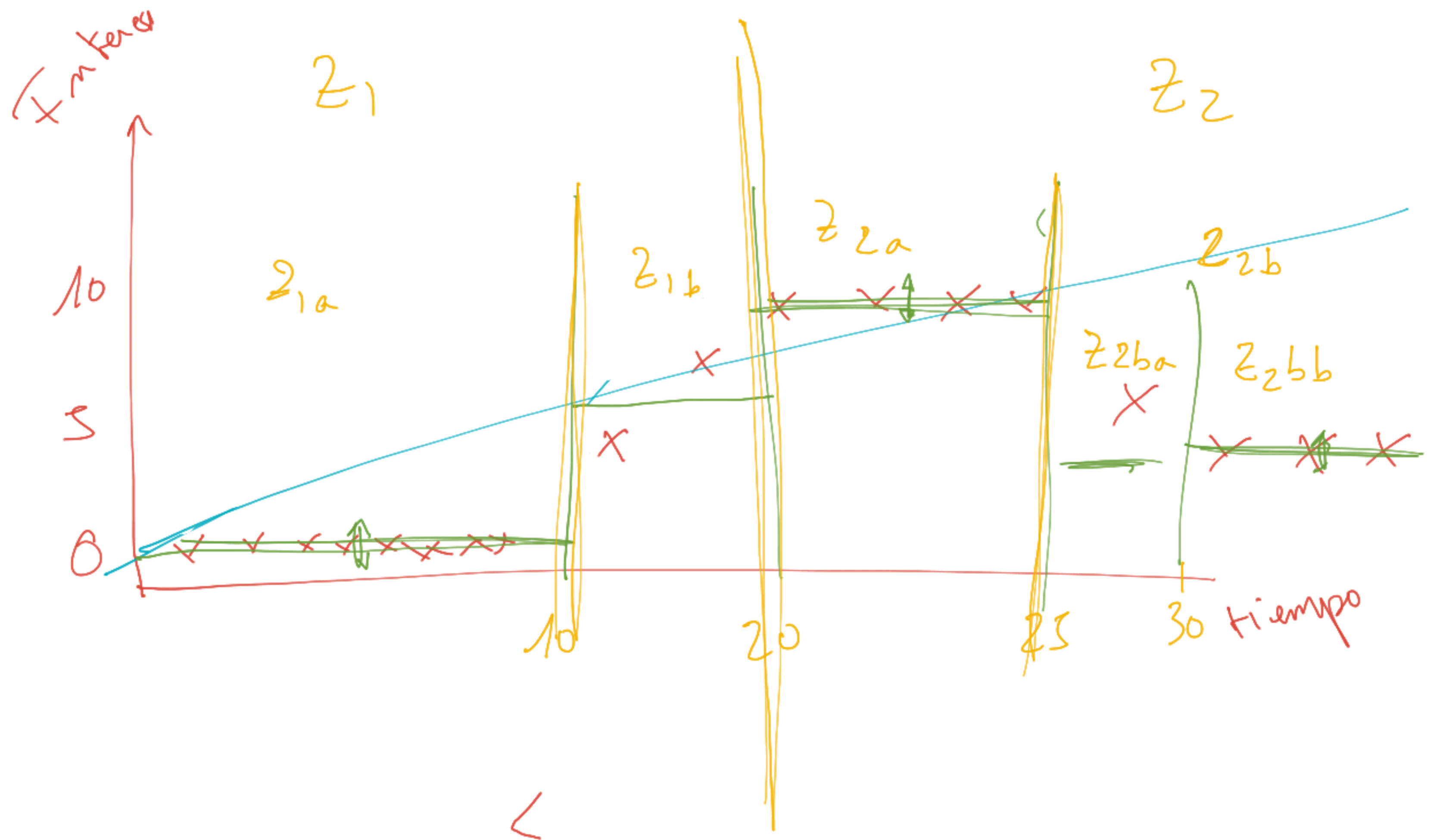
```
reg_test = GridSearchCV(KNeighborsRegressor(),
```

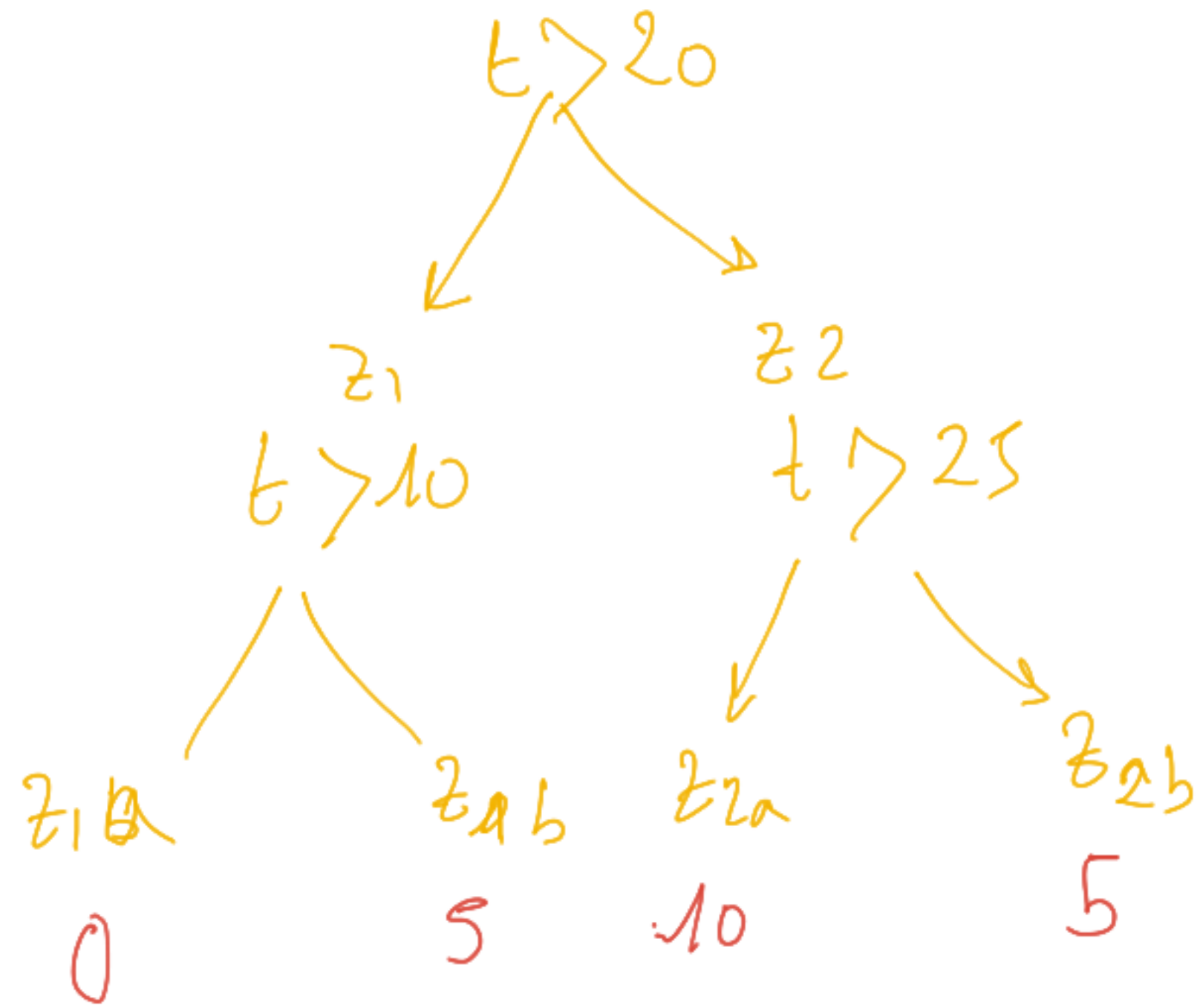
```
                        param_grid={"n_neighbors":np.arange(3,50)})
```

```
# Fit will test all of the combinations
```

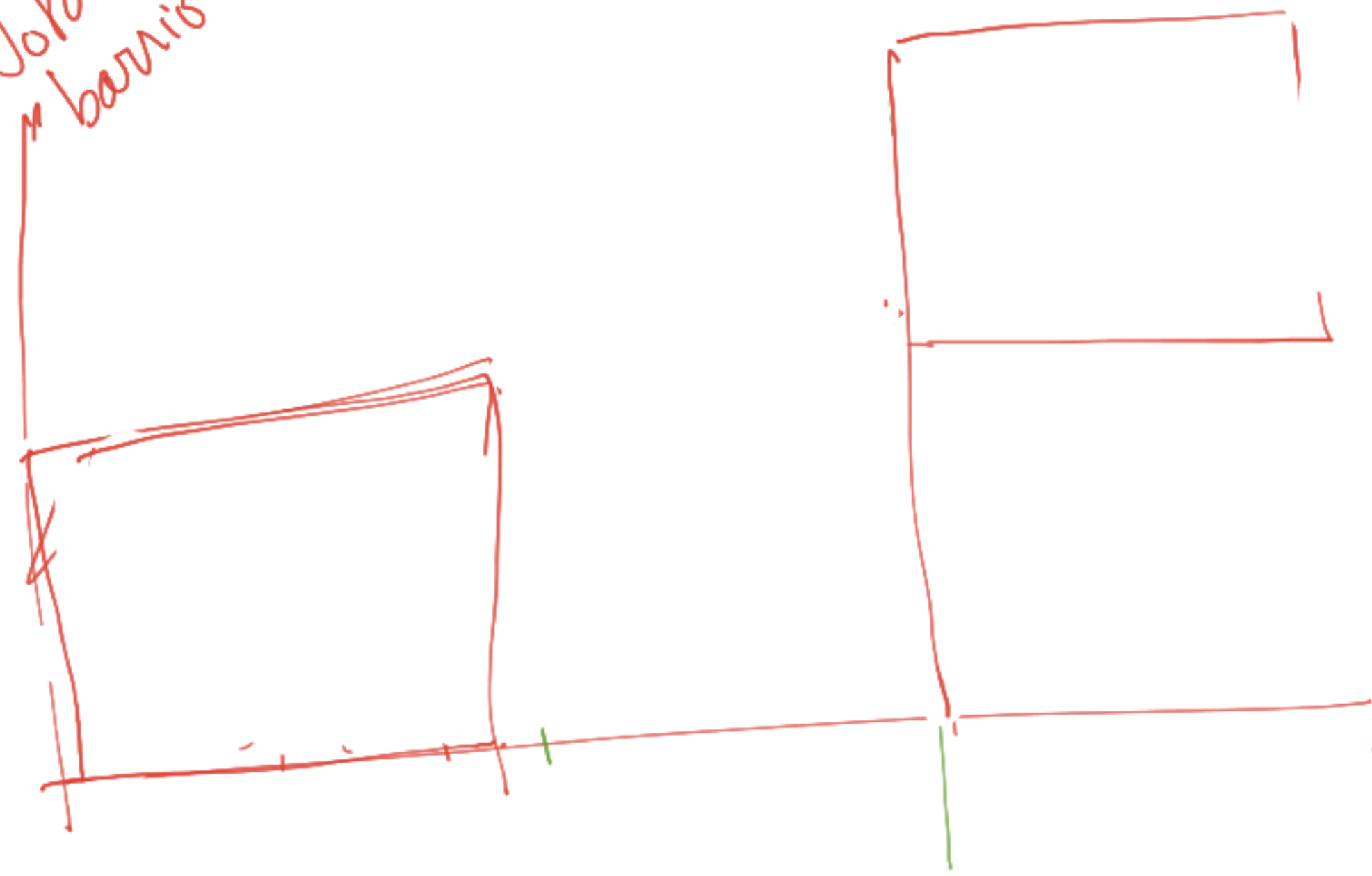
```
reg_test.fit(X,y)
```

← Evalua varios
modelos y se
QUEDA EL MEJOR.

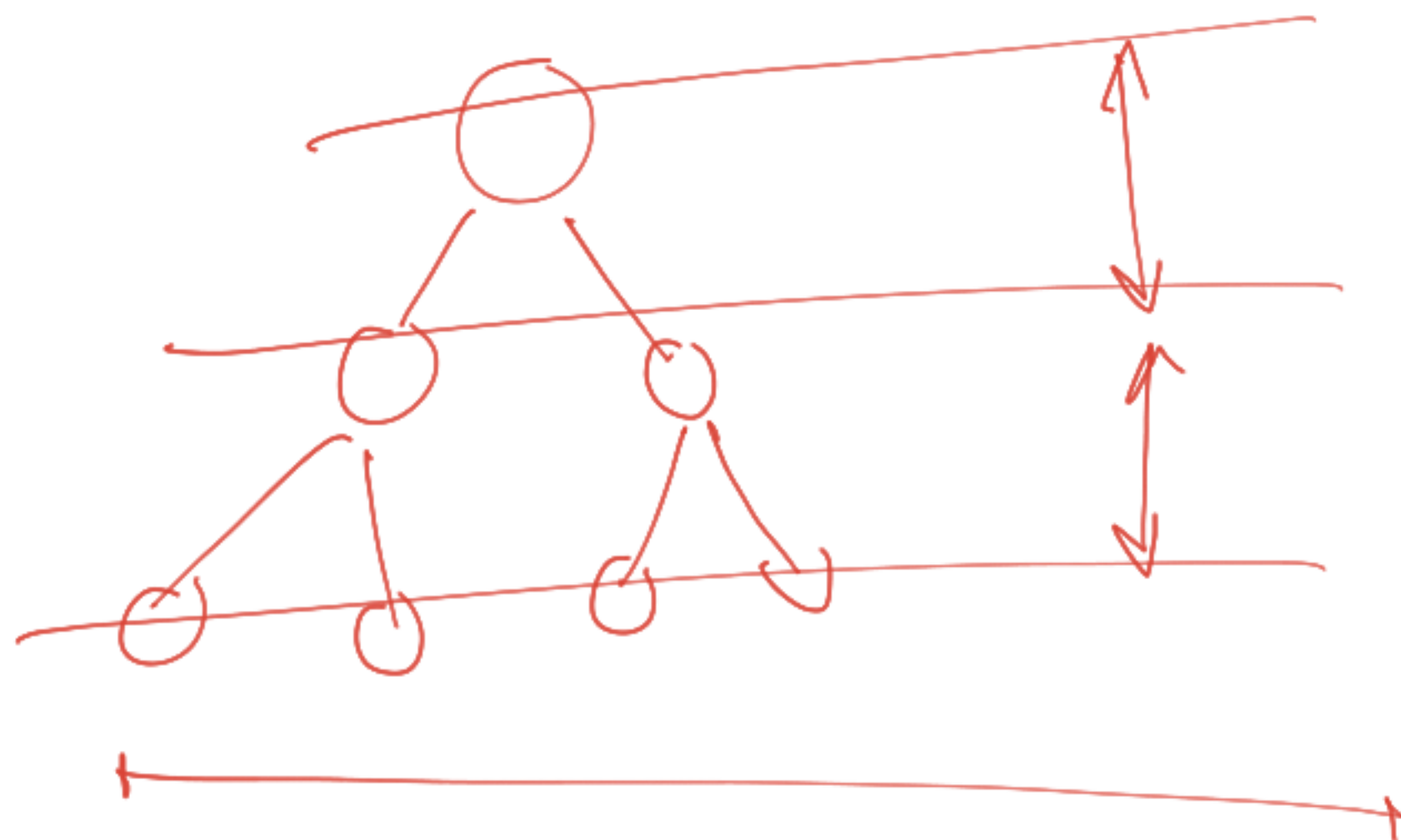




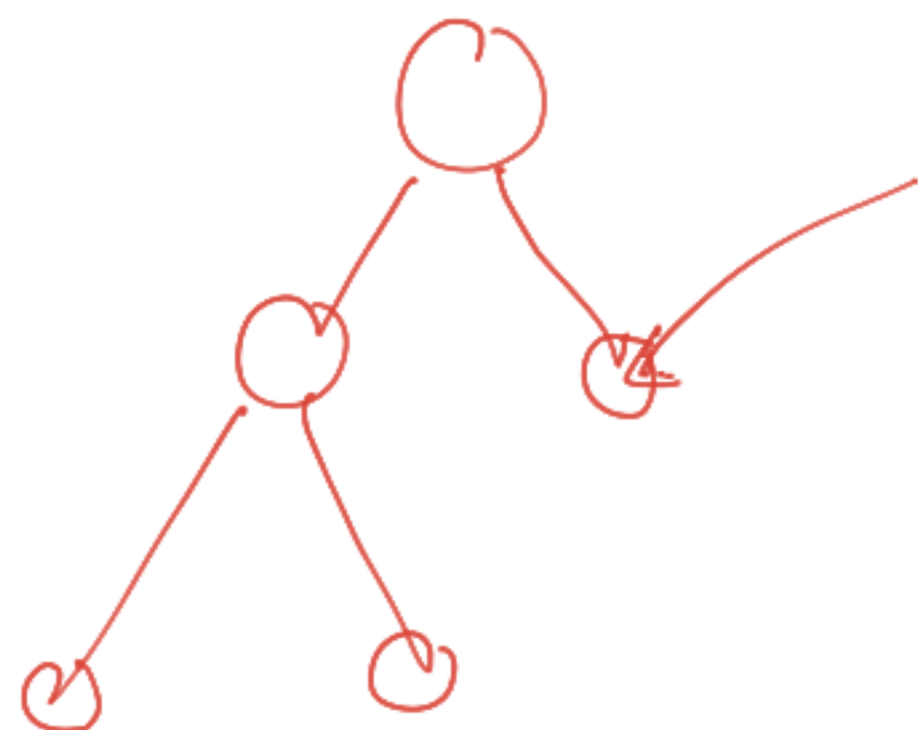
Noka
barris



Area-



Max Depth



num elem > min samples leaf.