



Gépi tanulás python alapokon

dmlab

Data science módszerek python környezetben

Regresszió

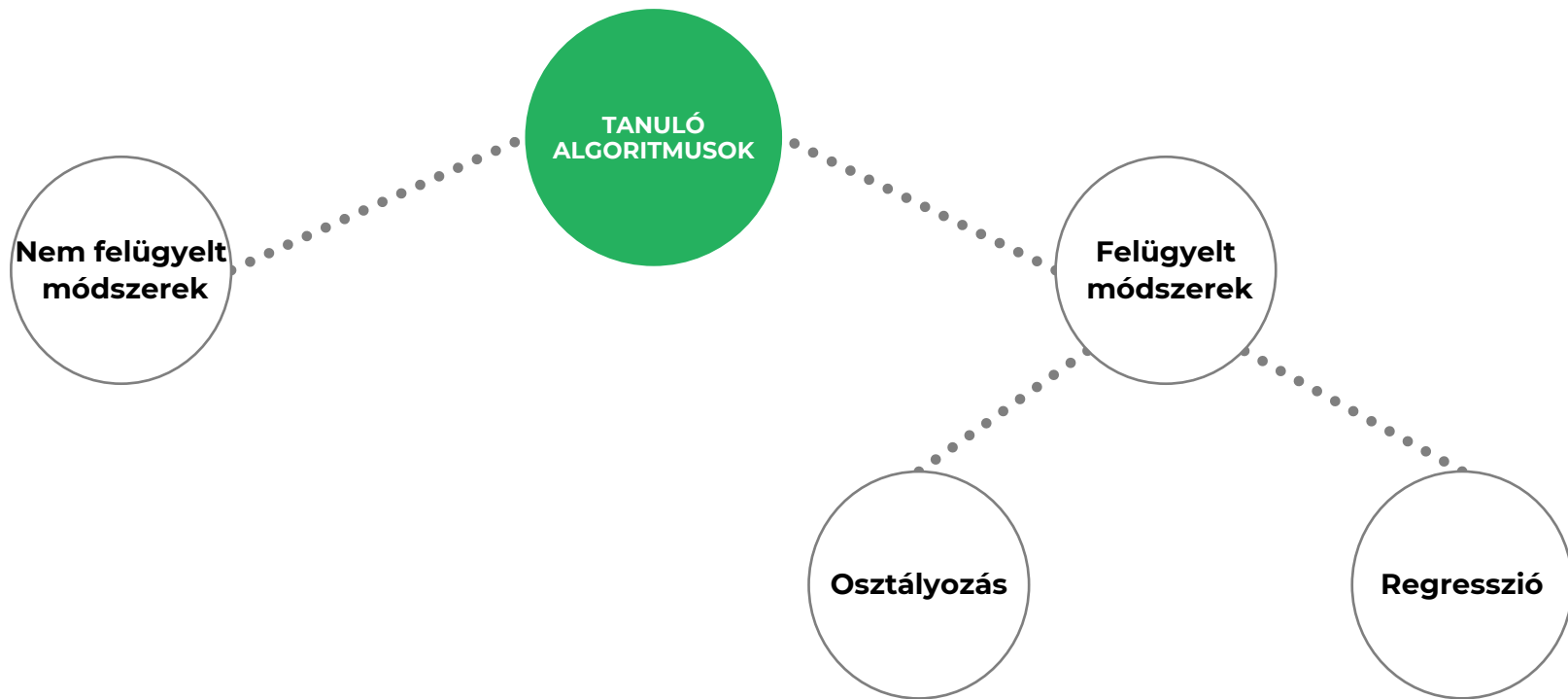
2024.10.14.

Jónás Dániel
data scientist

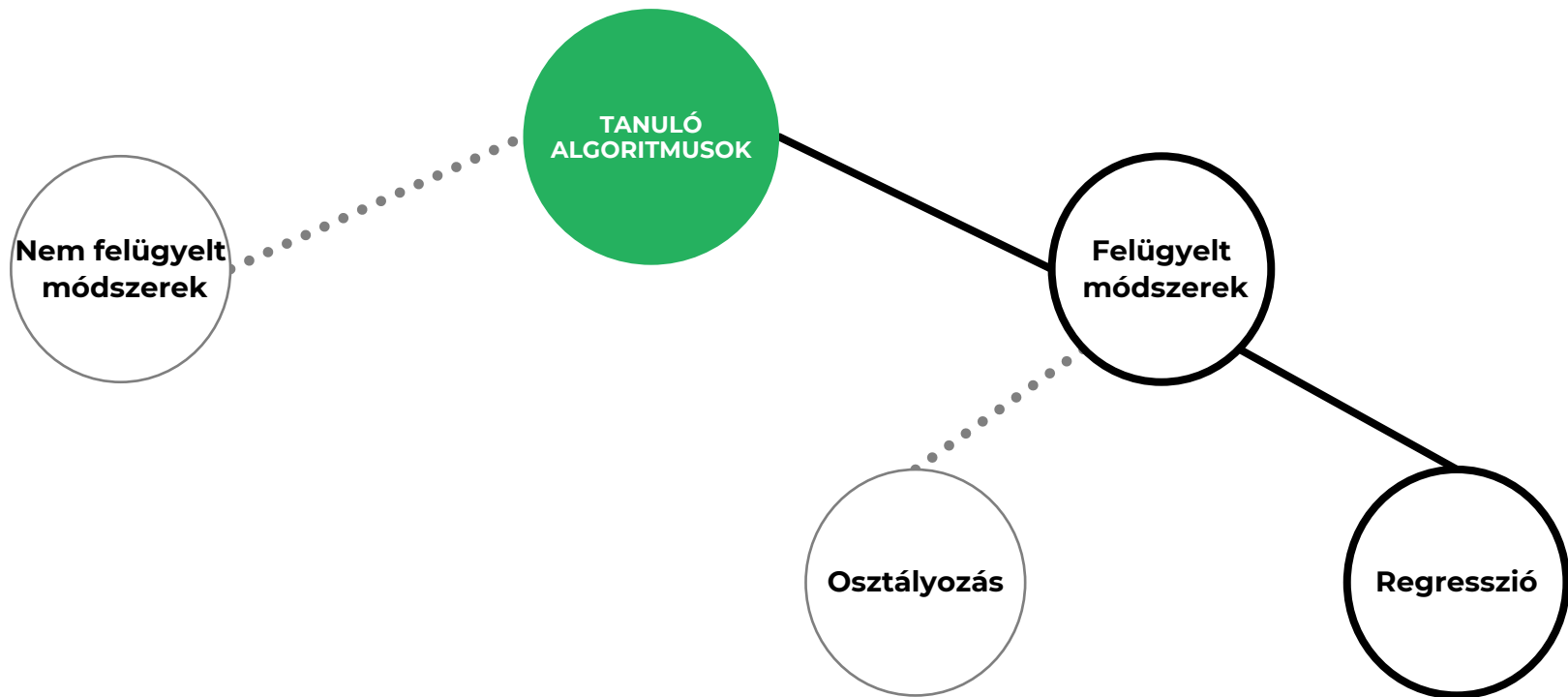


ADATELEMZÉSI PROBLÉMÁK CSOPORTOSÍTÁSA

MODELLEZÉSI ALAPOK



MODELLEZÉSI ALAPOK



A REGRESSZIÓ

Jellemzői

- Felügyelt tanulás
- Célváltozó: folytonos, numerikus változó
- Tipikus use-case: ár, szenzoradat, fogyasztás stb. előrejelzése
- Regressziós gépi tanulási modellek: lineáris, regularizált, döntési fa-alapú, stb..
- Kiértékelési metrikák: MAE, MSE, MAPE

Feladat

Van egy tanító adathalmazom
... ahol ismerem a tulajdonságokat és

Tulajdonságleíró attribútumok – **bemeneti változók**

Életkor	Foglalkozás	Lakásméret	Gyerekek	Autók száma
42	0432	82	3	1
88	0312	44	0	0
22	0101	32	0	1
38	4003	102	1	2

Feladat

Van egy tanító adathalmazom
... ahol ismerem a tulajdonságokat és

X				
X1	X2	X3	X4	X5

Tulajdonságleíró attribútumok – **bemeneti változók**

Életkor	Foglalkozás	Lakásméret	Gyerekek	Autók száma
42	0432	82	3	1
88	0312	44	0	0
22	0101	32	0	1
38	4003	102	1	2

Feladat

Van egy tanító adathalmazom
... ahol ismerem a tulajdonságokat és
... ahol ismerem a célváltozót

X					y
X1	X2	X3	X4	X5	

Tulajdonságleíró attribútumok – **bemeneti változók**

Célváltozó
Target
Címke

Életkor	Foglalkozás	Lakásméret	Gyerekek	Autók száma	Jövedelem
42	0432	82	3	1	280 000
88	0312	44	0	0	90 000
22	0101	32	0	1	180 000
38	4003	102	1	2	400 000

Feladat

És van egy új sor, ahol
... ismerem a tulajdonságokat

Életkor	Foglalkozás	Lakásméret	Gyerekek	Autók száma	Jövedelem
55	3257	62	1	3	????

X					y
X1	X2	X3	X4	X5	

Tulajdonságleíró attribútumok – **bemeneti változók**

Célváltozó

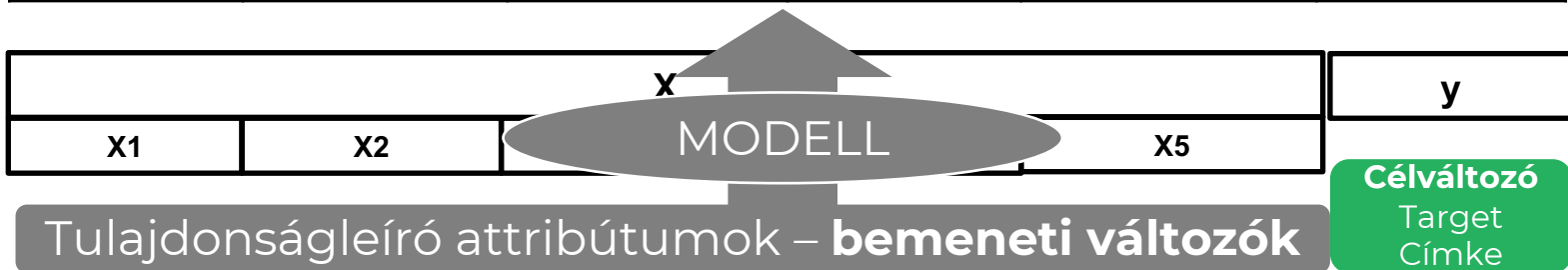
Target
Címke

Életkor	Foglalkozás	Lakásméret	Gyerekek	Autók száma	Jövedelem
42	0432	82	3	1	280 000
88	0312	44	0	0	90 000
22	0101	32	0	1	180 000
38	4003	102	1	2	400 000

Feladat

És van egy új sor, ahol
... ismerem a tulajdonságokat
... modell mondja meg mekkora a célváltozó

Életkor	Foglalkozás	Lakásméret	Gyerekek	Autók száma	Jövedelem
55	3257	62	1	3	????



Életkor	Foglalkozás	Lakásméret	Gyerekek	Autók száma	Jövedelem
42	0432	82	3	1	280 000
88	0312	44	0	0	90 000
22	0101	32	0	1	180 000
38	4003	102	1	2	400 000

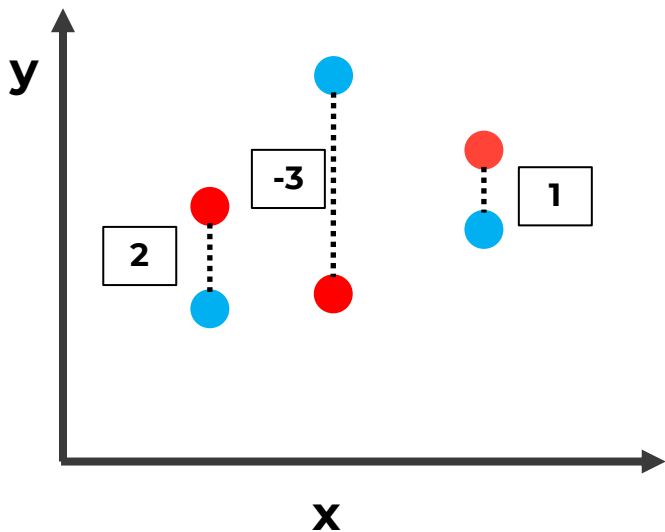


KIÉRTÉKELÉSI METRIKÁK

VISSZAMÉRÉSI FÜGGVÉNYEK

Mi lenne a legjobb tipp? Mi minimalizálja a hibát?
Attól függ, milyen hibát.

Hiba = valós érték – prediktált érték

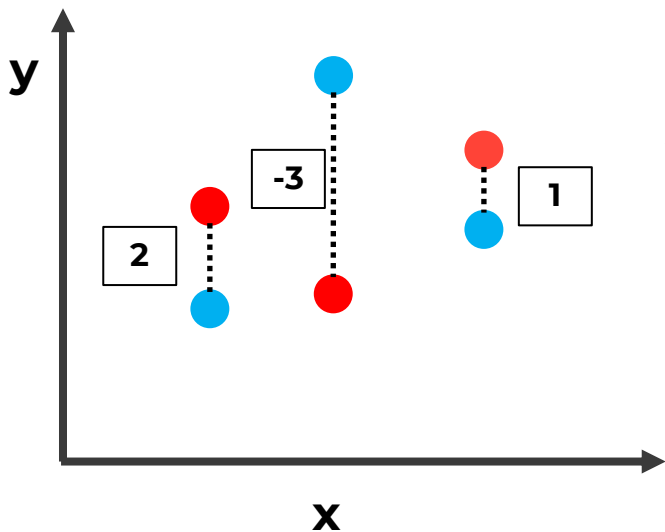


VISSZAMÉRÉSI FÜGGVÉNYEK

Mi lenne a legjobb tipp? Mi minimalizálja a hibát?
Attól függ, milyen hibát.

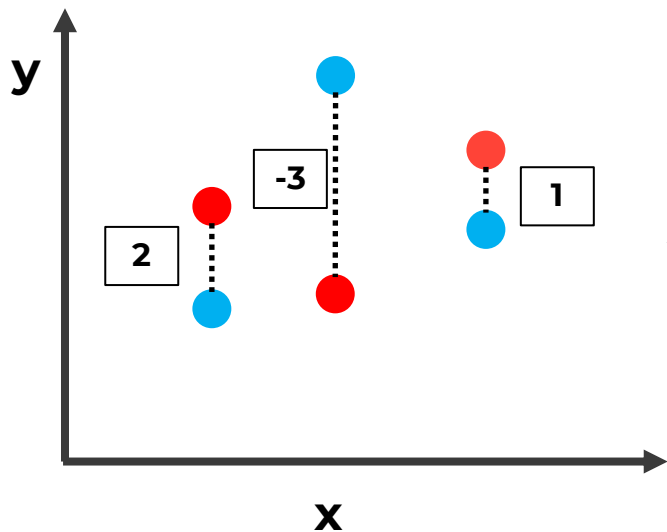
Hiba = valós érték – prediktált érték

$$\text{Átlagos Hiba} = (2 + -3 + 1) / 3 = 0$$



VISSZAMÉRÉSI FÜGGVÉNYEK

Mi lenne a legjobb tipp? Mi minimalizálja a hibát?
Attól függ, milyen hibát.



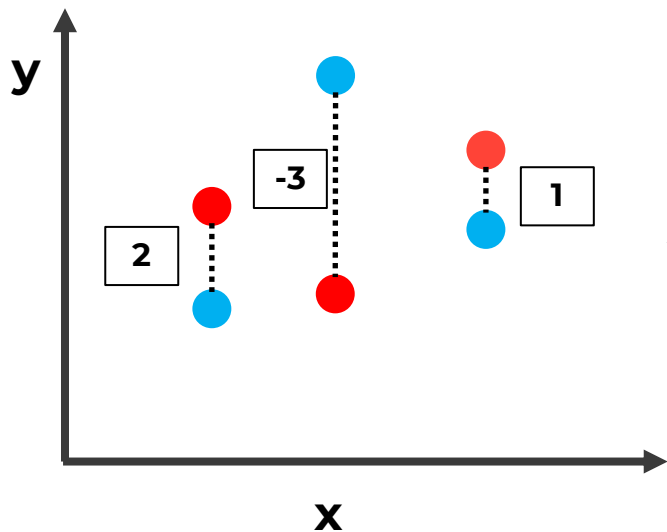
Hiba = valós érték – prediktált érték

$$\text{Átlagos Hiba} = (2 + -3 + 1) / 3 = 0$$

$$\text{Abszolút Átlagos Hiba} = (2 + 3 + 1) / 3 = 2$$

VISSZAMÉRÉSI FÜGGVÉNYEK

Mi lenne a legjobb tipp? Mi minimalizálja a hibát?
Attól függ, milyen hibát.



Hiba = valós érték – prediktált érték

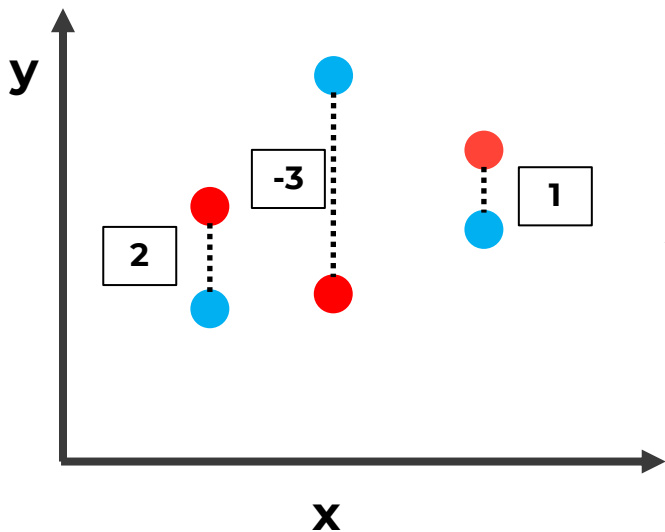
$$\text{Átlagos Hiba} = (2 + -3 + 1) / 3 = 0$$

$$\text{Abszolút Átlagos Hiba} = (2 + 3 + 1) / 3 = 2$$

$$\text{Négyzetes Átlagos Hiba} = (4 + 9 + 1) / 3 = 4.67$$

VISSZAMÉRÉSI FÜGGVÉNYEK

Mi lenne a legjobb tipp? Mi minimalizálja a hibát?
Attól függ, milyen hibát.



Hiba = valós érték – prediktált érték

$$\text{Átlagos Hiba} = (2 + -3 + 1) / 3 = 0$$

$$\text{Abszolút Átlagos Hiba} = (2 + 3 + 1) / 3 = 2$$

$$\text{Négyzetes Átlagos Hiba} = (4 + 9 + 1) / 3 = 4.67$$

$$\text{Abszolút Százalékos Hiba} = (0.3 + 1.3 + 0.2) / 3 = 0.6 = 60\%$$

KIÉRTÉKELÉSI METRIKÁK

MEAN ABSOLUTE ERROR (MAE)

- $|X_{\text{valós}} - X_{\text{pred}}|$
- $\text{MAE} = |10 - 12| = 2$
- A hiba mértéke a fontos, nem az iránya
- Könnyen értelmezhető

KIÉRTÉKELÉSI METRIKÁK

MEAN ABSOLUTE ERROR (MAE)

- $|X_{\text{valós}} - X_{\text{pred}}|$
- $\text{MAE} = |10 - 12| = 2$
- A hiba mértéke a fontos, nem az iránya
- Könnyen értelmezhető

MEAN SQUARED ERROR (MSE)

- $(X_{\text{valós}} - X_{\text{pred}})^2$
- $\text{MSE} = (10 - 12)^2 = 4$
- Könnyen optimalizálható
- A nagyobb eltérés nagyobb hibának számít

KIÉRTÉKELÉSI METRIKÁK

MEAN ABSOLUTE ERROR (MAE)

- $|X_{\text{valós}} - X_{\text{pred}}|$
- $\text{MAE} = |10 - 12| = 2$
- A hiba mértéke a fontos, nem az iránya
- Könnyen értelmezhető

MEAN SQUARED ERROR (MSE)

- $(X_{\text{valós}} - X_{\text{pred}})^2$
- $\text{MSE} = (10 - 12)^2 = 4$
- Könnyen optimalizálható
- A nagyobb eltérés nagyobb hibának számít

MEAN ABSOLUTE PERCENTAGE ERROR (MAPE)

- $\frac{|X_{\text{valós}} - X_{\text{pred}}|}{X_{\text{valós}}}$
- $\frac{|10 - 12|}{10} = 20\%$
- Általánosan értelmezhető
- A hiba mértéke függ az értékkészlet nagyságrendjétől

KIÉRTÉKELÉSI METRIKÁK

MEAN ABSOLUTE ERROR (MAE)

- Ismerni kell a nagyságrendet
- pl MAE = 100 jó vagy rossz?

KIÉRTÉKELÉSI METRIKÁK

MEAN ABSOLUTE ERROR (MAE)

- Ismerni kell a nagyságrendet
- pl MAE = 100 jó vagy rossz?

MEAN SQUARED ERROR (MSE)

- $MSE = 250.000.000.000$
- Jó vagy rossz?
- 50 milliós nagyságrendű célváltozóknál ez 1% körüli hiba
- $RMSE = \sqrt{MSE} = 500.000$

KIÉRTÉKELÉSI METRIKÁK

MEAN ABSOLUTE ERROR (MAE)

- Ismerni kell a nagyságrendet
- pl MAE = 100 jó vagy rossz?

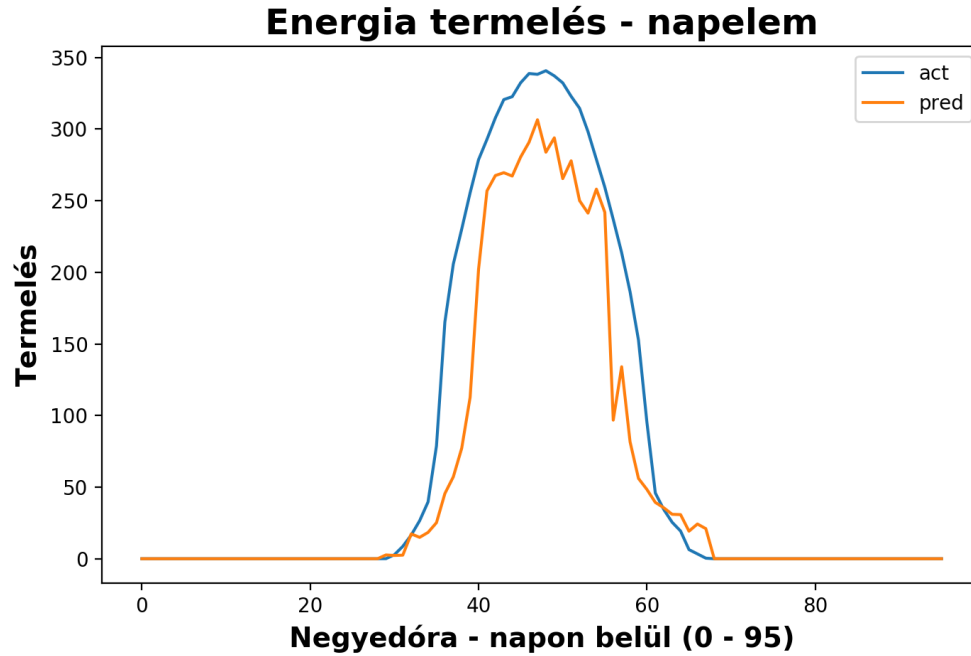
MEAN SQUARED ERROR (MSE)

- $MSE = 250.000.000.000$
- Jó vagy rossz?
- 50 milliós nagyságrendű célváltozóknál ez 1% körüli hiba
- $RMSE = \sqrt{MSE} = 500.000$

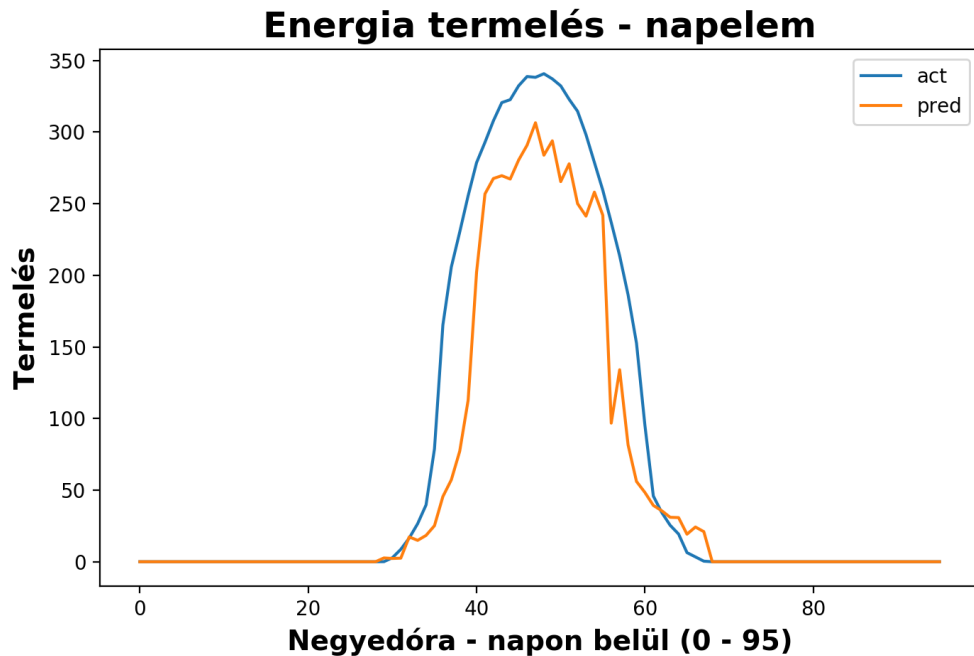
MEAN ABSOLUTE PERCENTAGE ERROR (MAPE)

- 0 értékű célváltozónál nem értelmezhető
- A hiba mértéke függ az értékkészlet nagyságrendjétől

MAPE

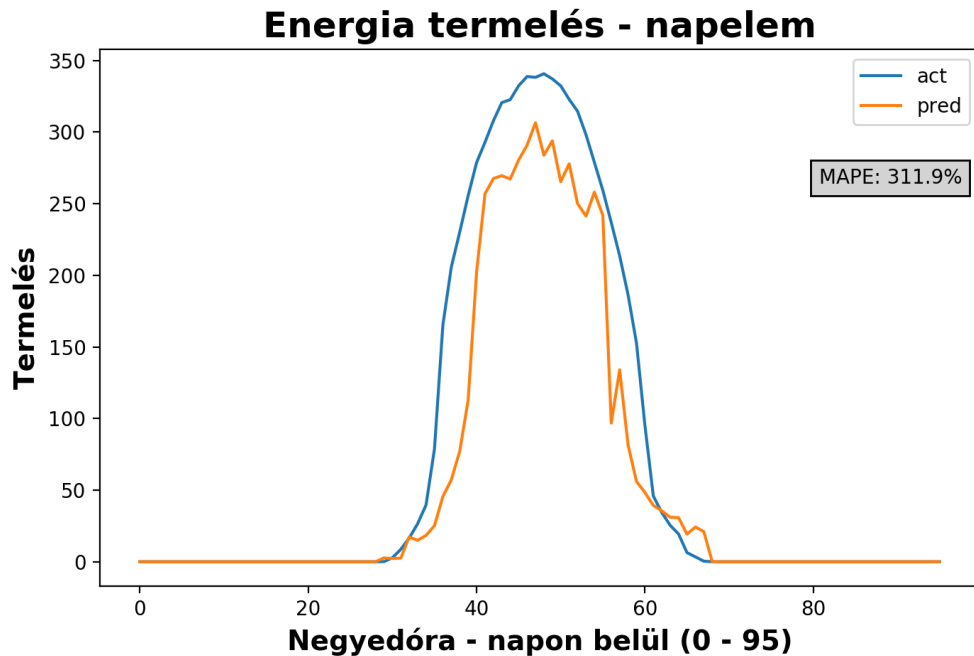


MAPE



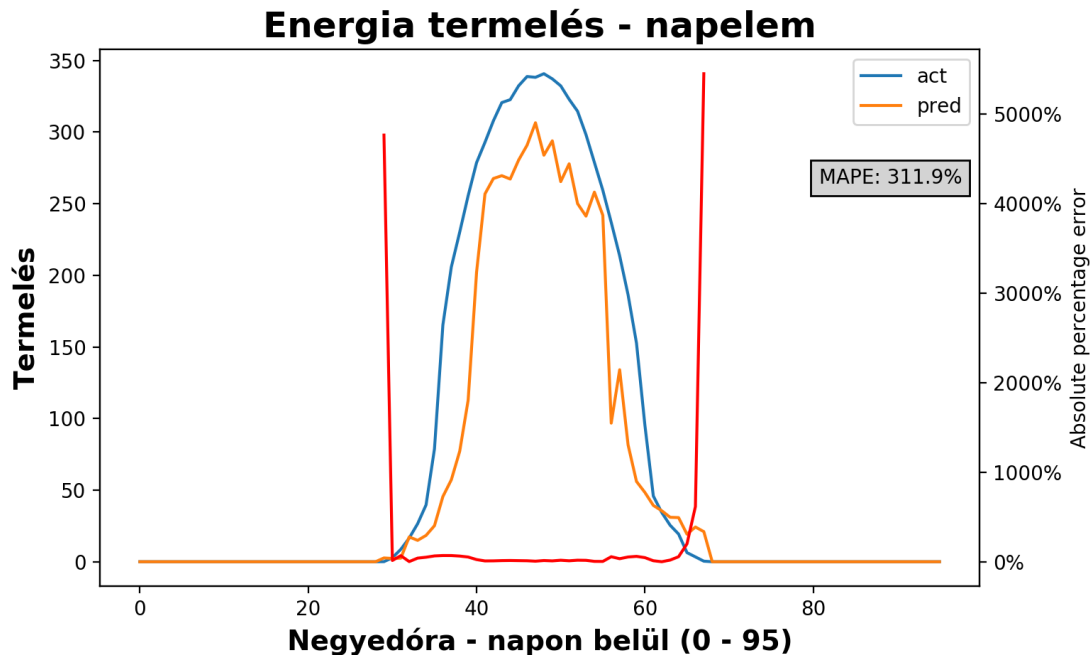
➤ Ha a célváltozó 0, akkor nem értelmezhető

MAPE



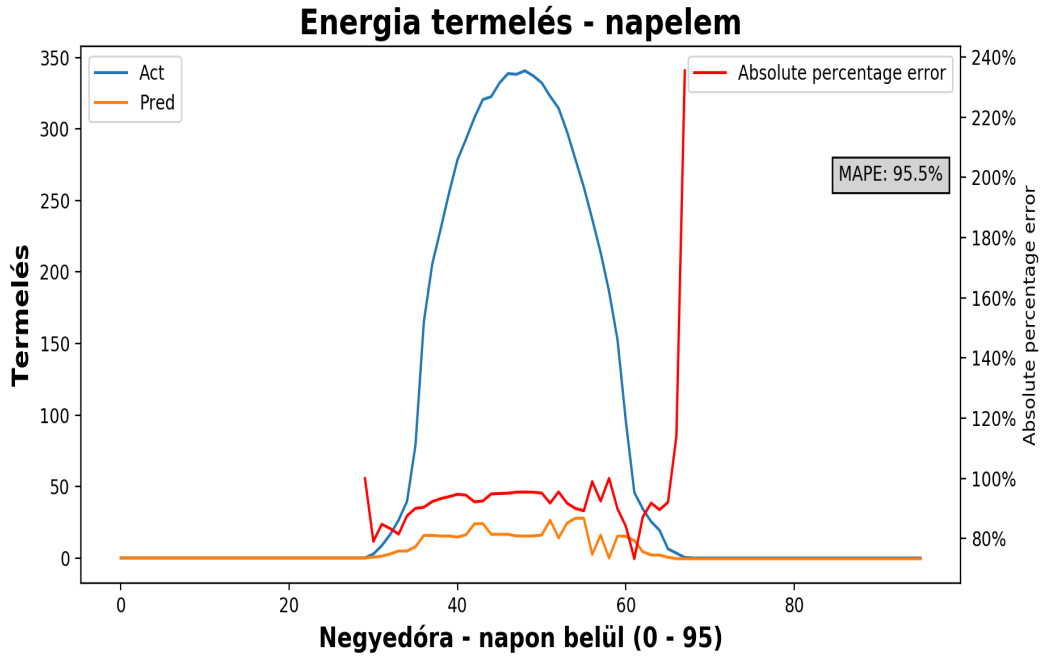
➤ Ha a célváltozó 0, akkor nem értelmezhető

MAPE

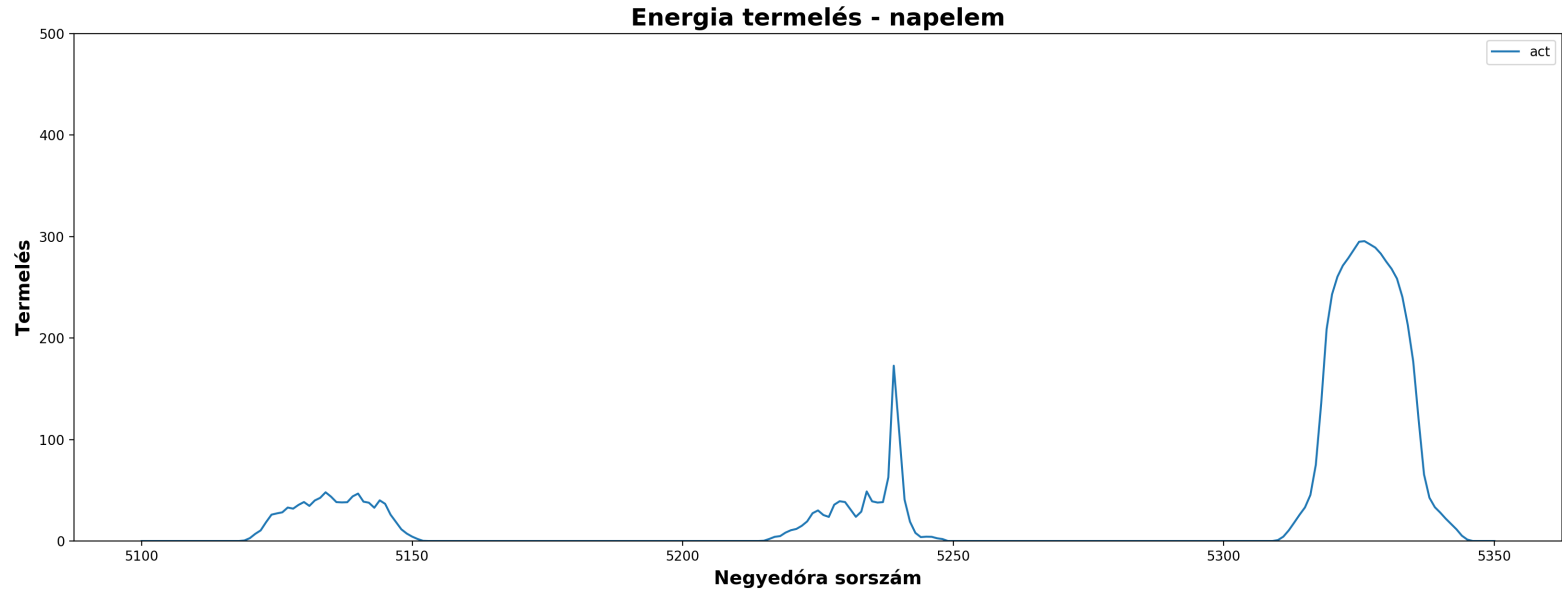


- Ha a célváltozó 0, akkor nem értelmezhető
- Ha az értékkészlet szórása nagy, a MAPE nem általánosan értelmezhető

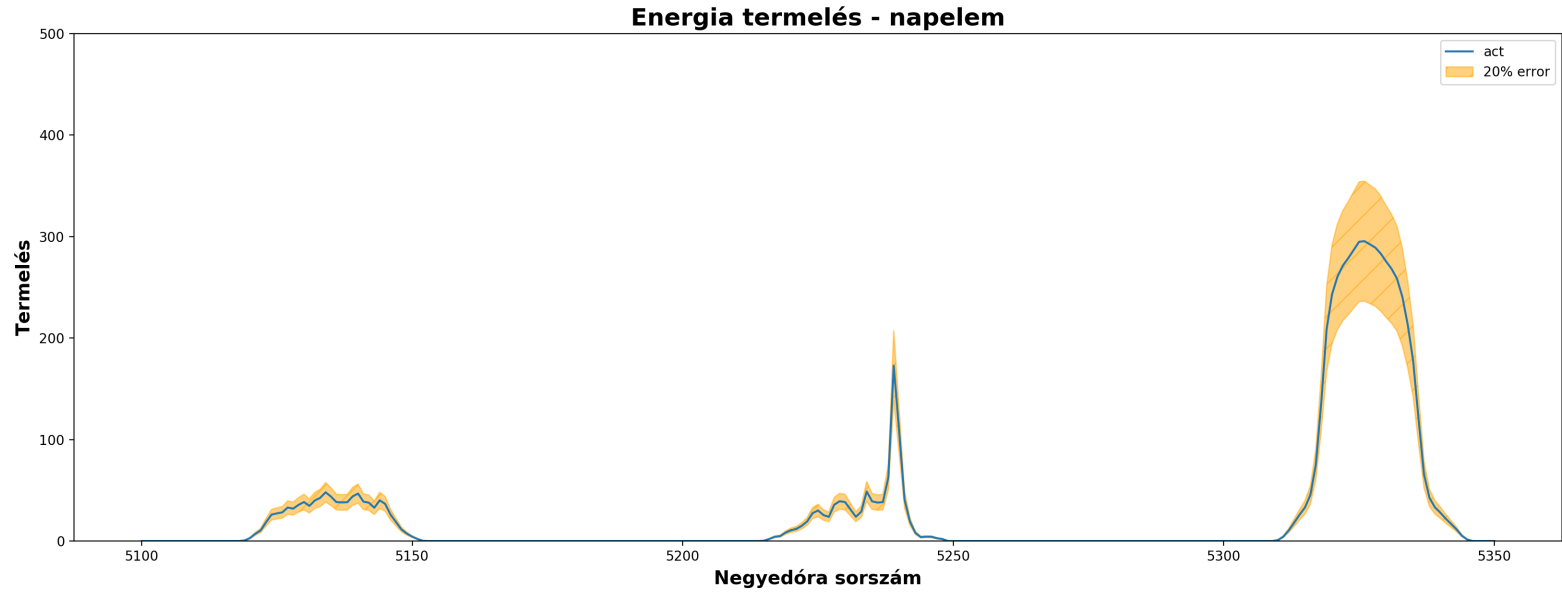
MAPE



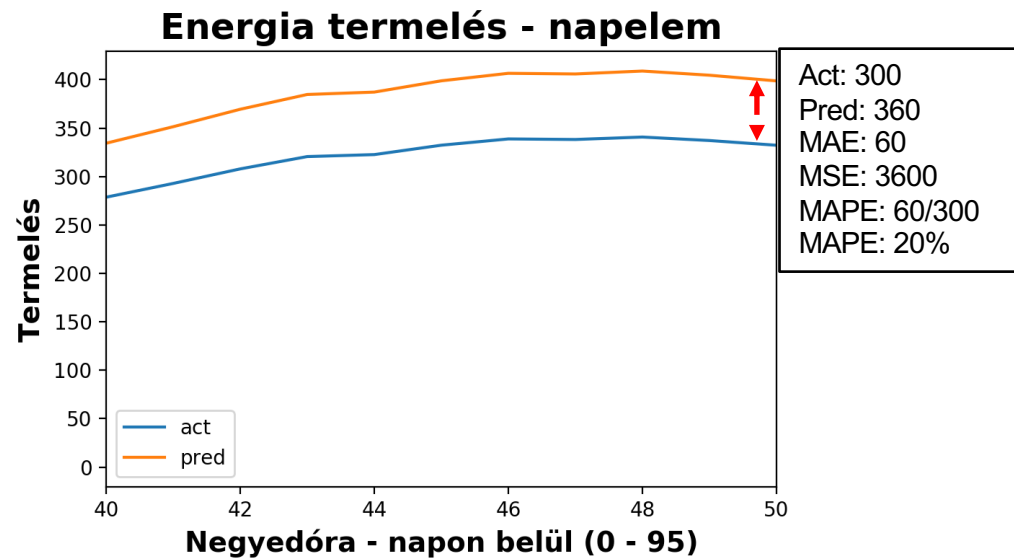
MAPE



MAPE

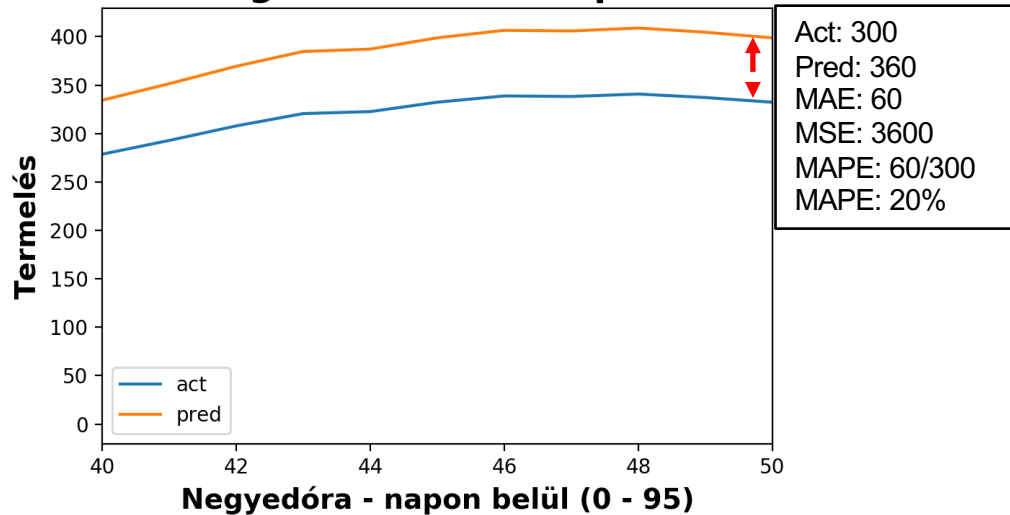


MAPE

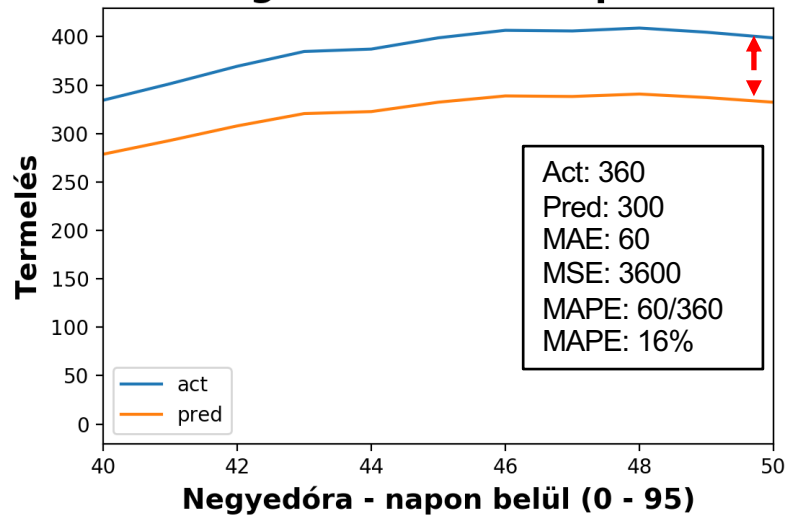


MAPE

Energia termelés - napelem



Energia termelés - napelem



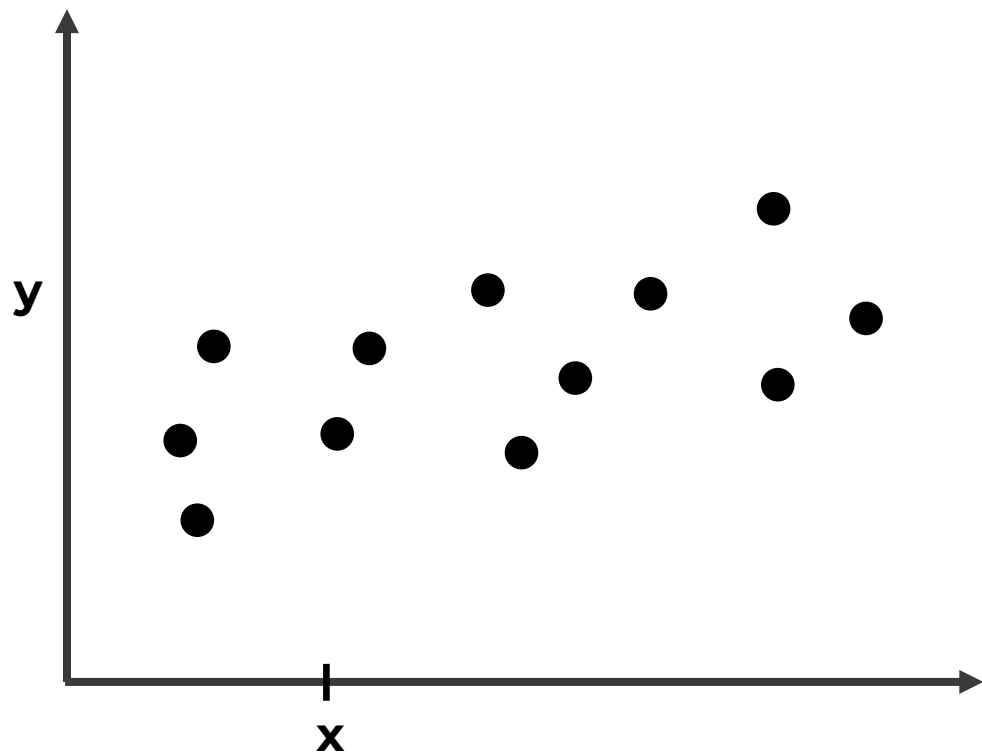
```
fifth_listing = normalized_listings.iloc[20][['accommodates', 'bathrooms']]
first_fifth_distance = distance.euclidean(first_listing, fifth_listing)

def predict_price_multivariate(new_listing_value, feature_columns):
    temp_df = norm_train_df
    temp_df['distance'] = distance.cdist(temp_df[feature_columns], [new_listing_value[feature_columns]])
    temp_df = temp_df.sort_values('distance')
    knn_5 = temp_df.price.iloc[:5]
    predicted_price = knn_5.mean()
    return(predicted_price)cols = ['accommodates', 'bathrooms']

norm_test_df['predicted_price'] = norm_test_df[cols].apply(predict_price_multivariate, feature_columns=cols, axis=1)
norm_test_df['squared_error'] = (norm_test_df['predicted_price'] - norm_test_df['price'])**2
mse = norm_test_df['squared_error'].mean()
mse = mse ** (1/2)
```

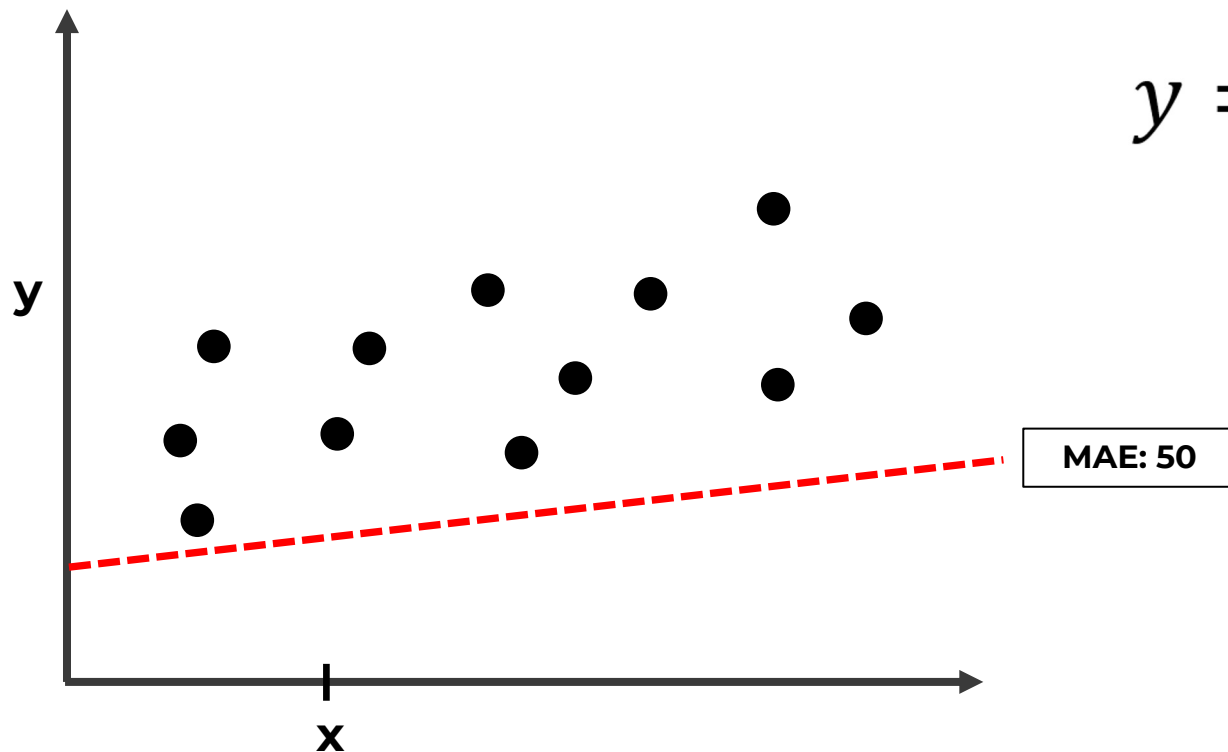
LINEÁRIS REGRESSZIÓ

Lineáris regresszió

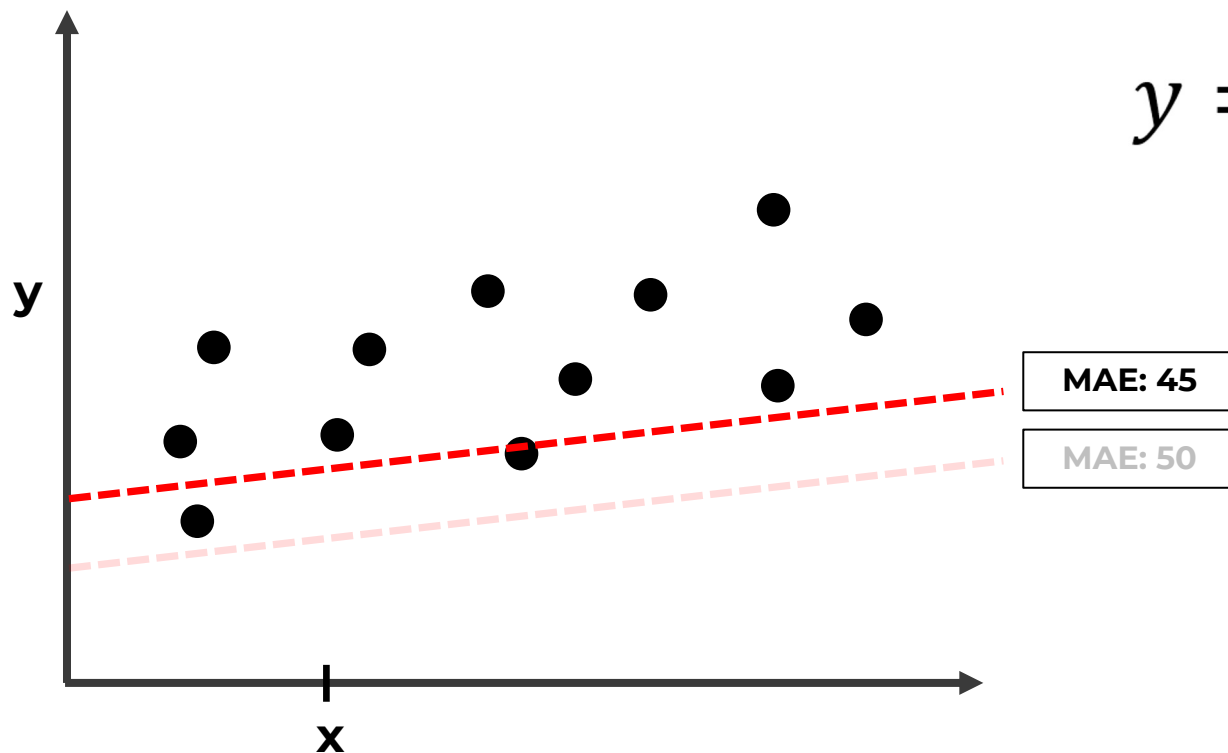


$$y = b_0 + b_1 * x$$

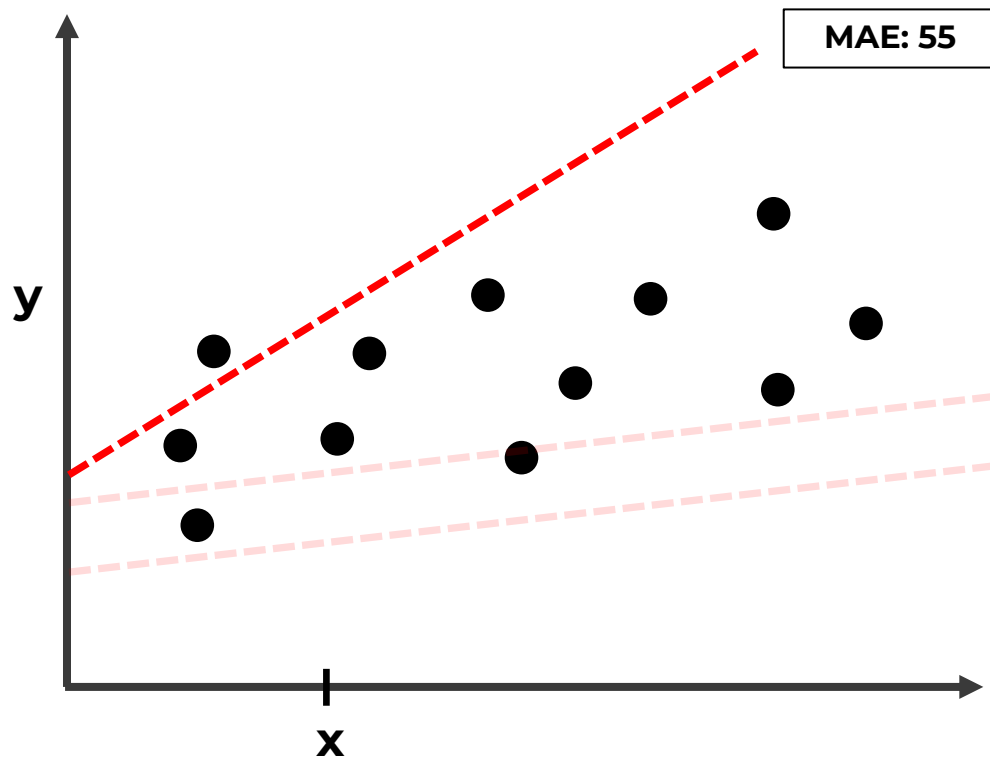
Lineáris regresszió



Lineáris regresszió

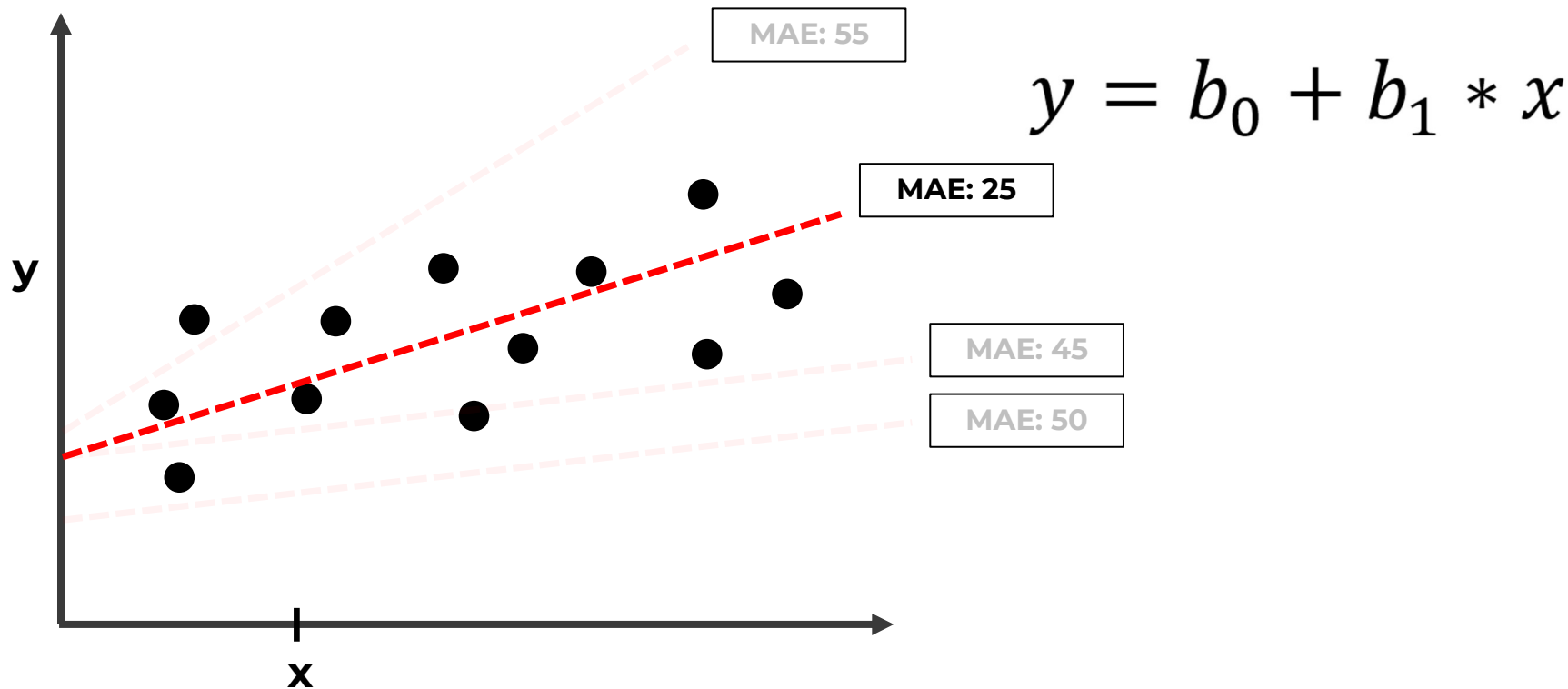


Lineáris regresszió



$$y = b_0 + b_1 * x$$

Lineáris regresszió



LINEÁRIS REGRESSZIÓ

<i>area</i>	<i>rooms</i>	<i>district</i>	<i>elevator</i>	<i>price</i>
120	4	3	1	90
50	2	3	0	45
82	3	8	1	59
46	1	13	0	40
70	3	5	0	71
65	2	3	1	54

X input features

y target

LINEÁRIS REGRESSZIÓ

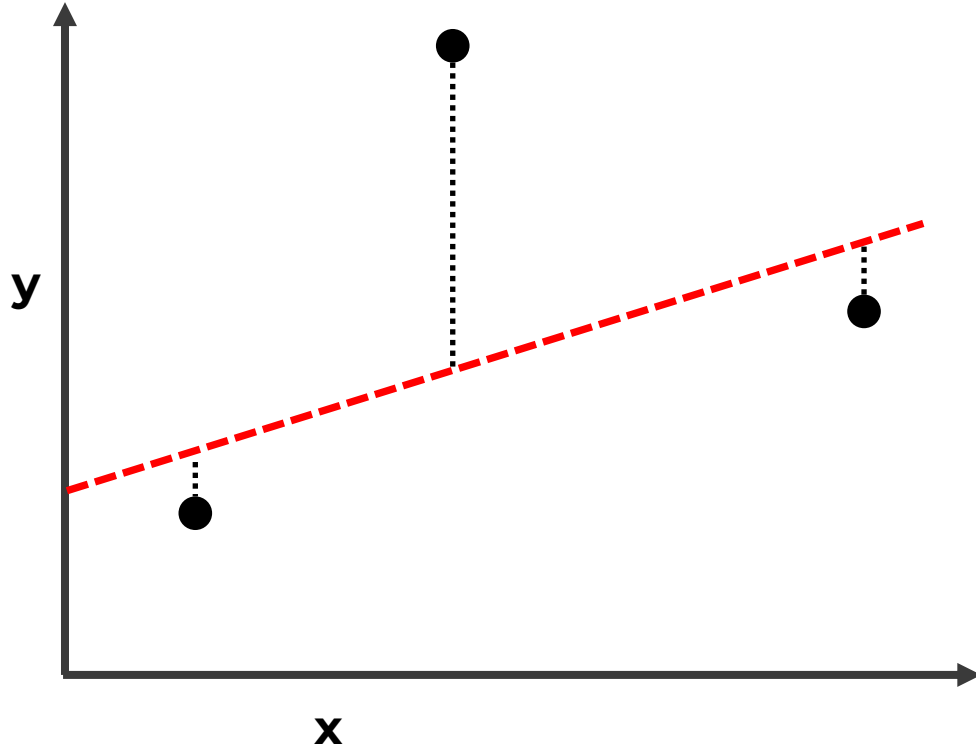
<i>area</i>	<i>rooms</i>	<i>district</i>	<i>elevator</i>	<i>price</i>
120	4	3	1	90
50	2	3	0	45
82	3	8	1	59
46	1	13	0	40
70	3	5	0	71
65	2	3	1	54

X input features

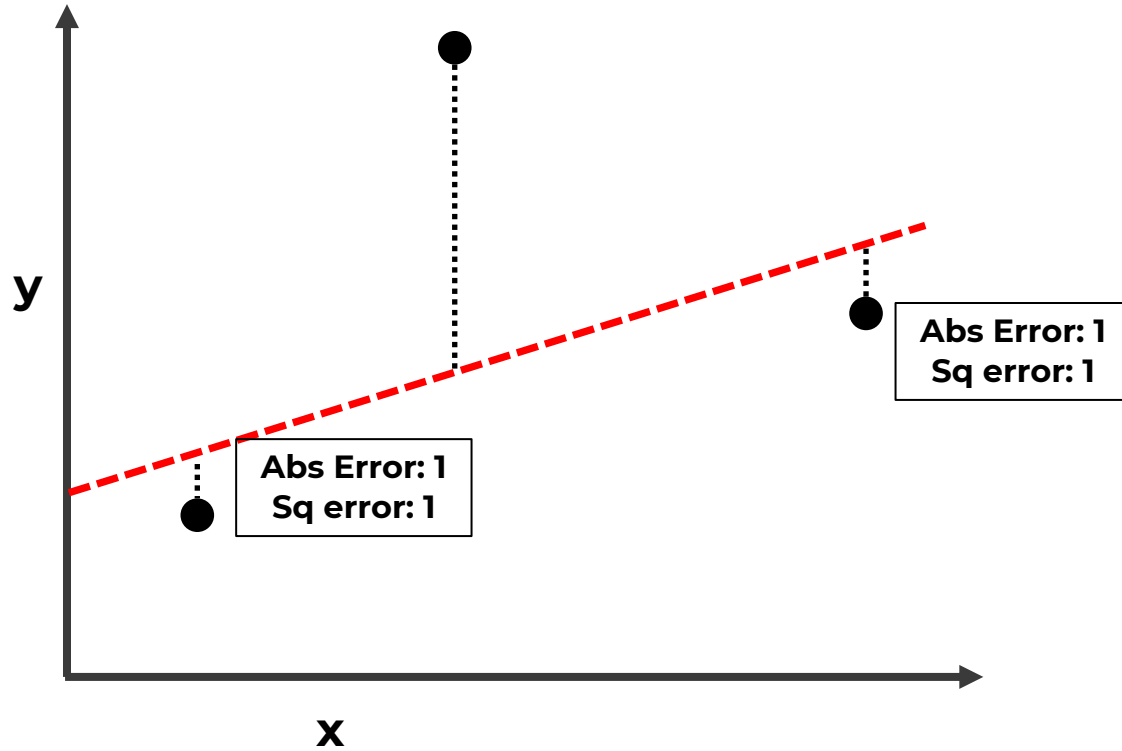
y target

$$y = b_0 + b_1 * x_1 + \dots + b_n * x_n$$

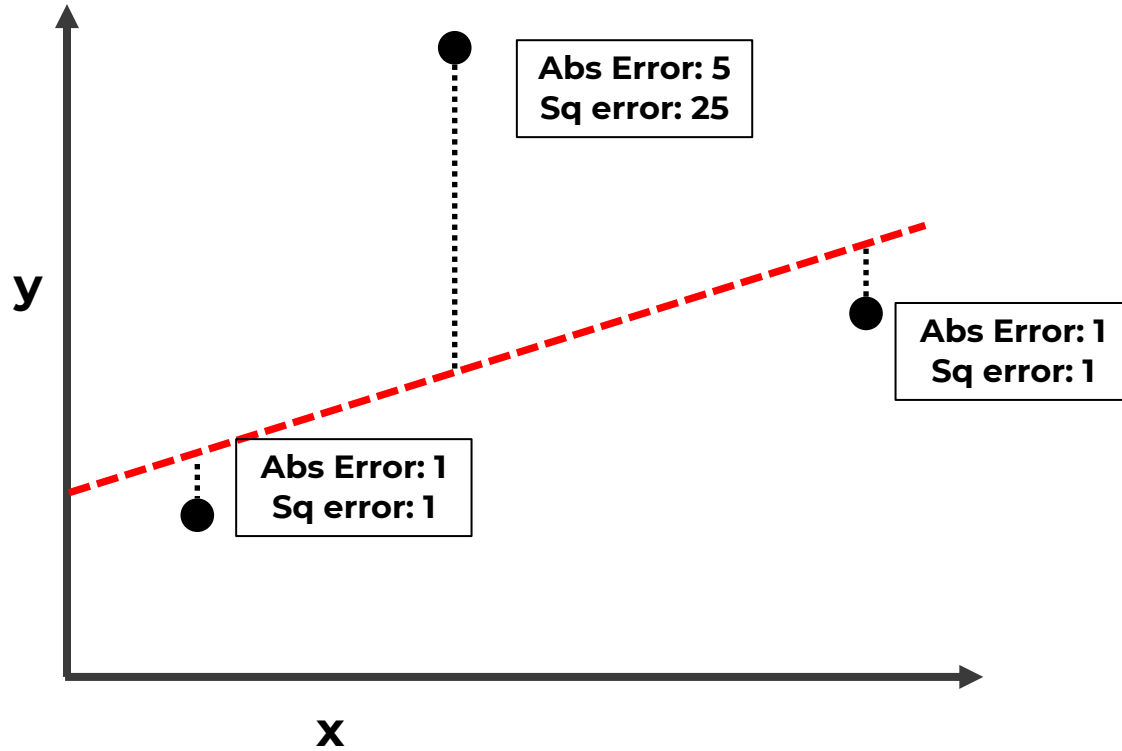
MSE



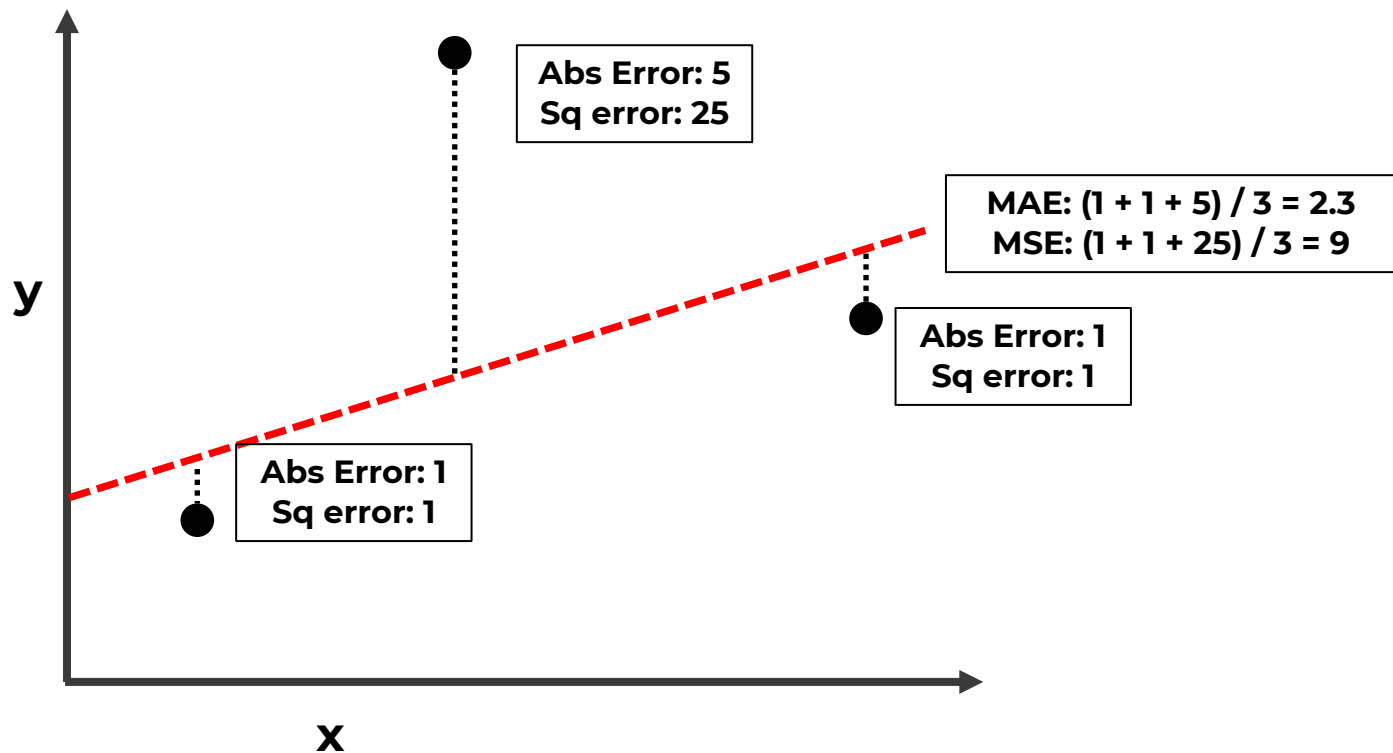
MSE



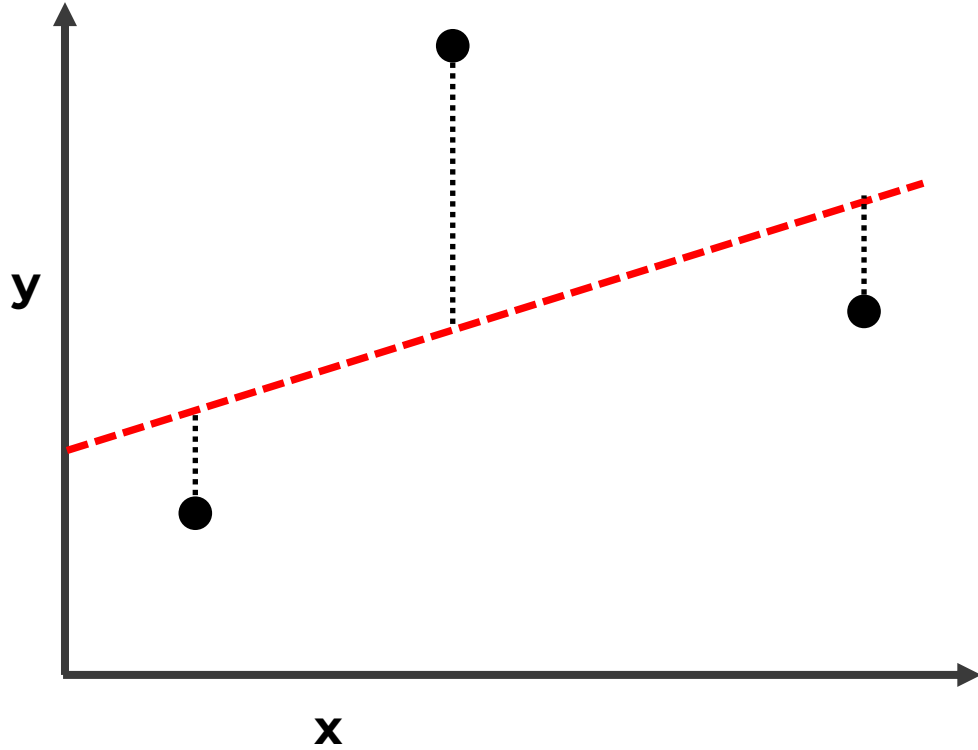
MSE



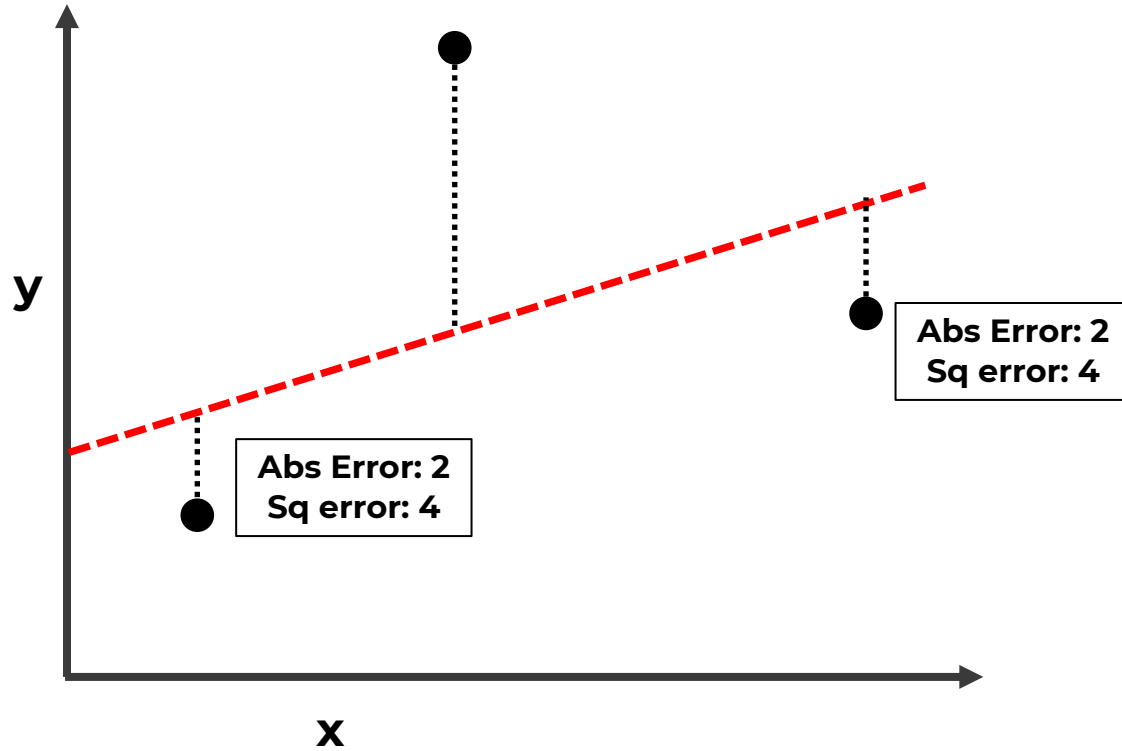
MSE



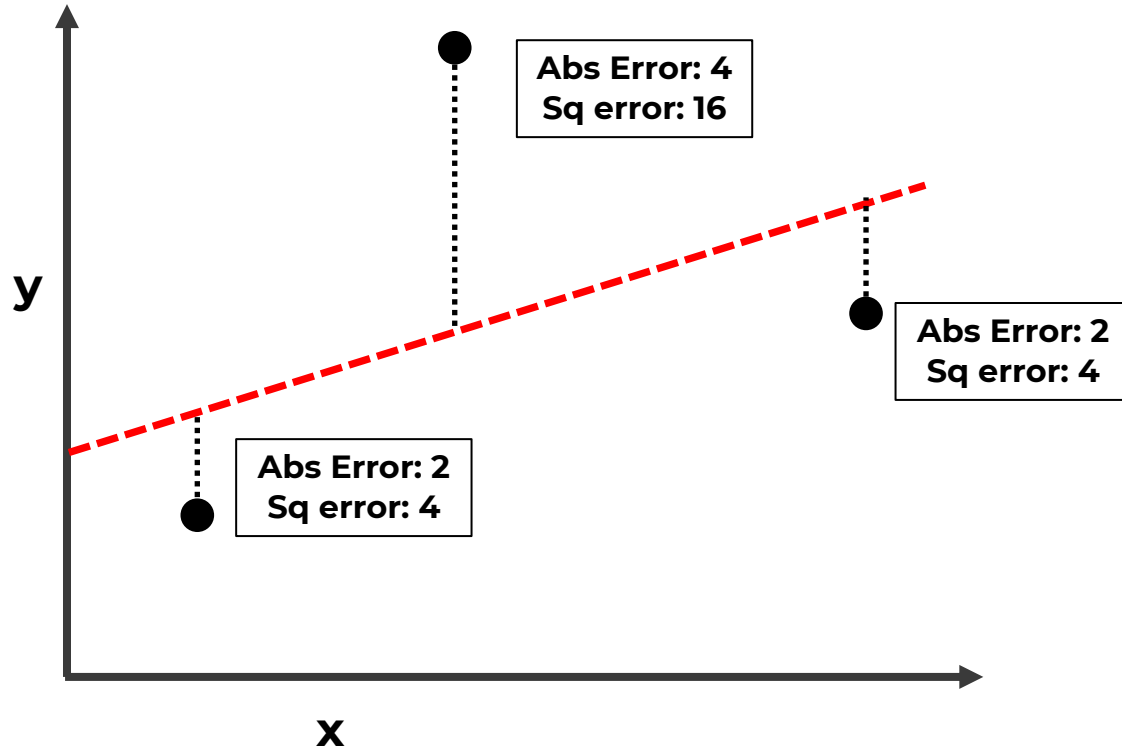
MSE



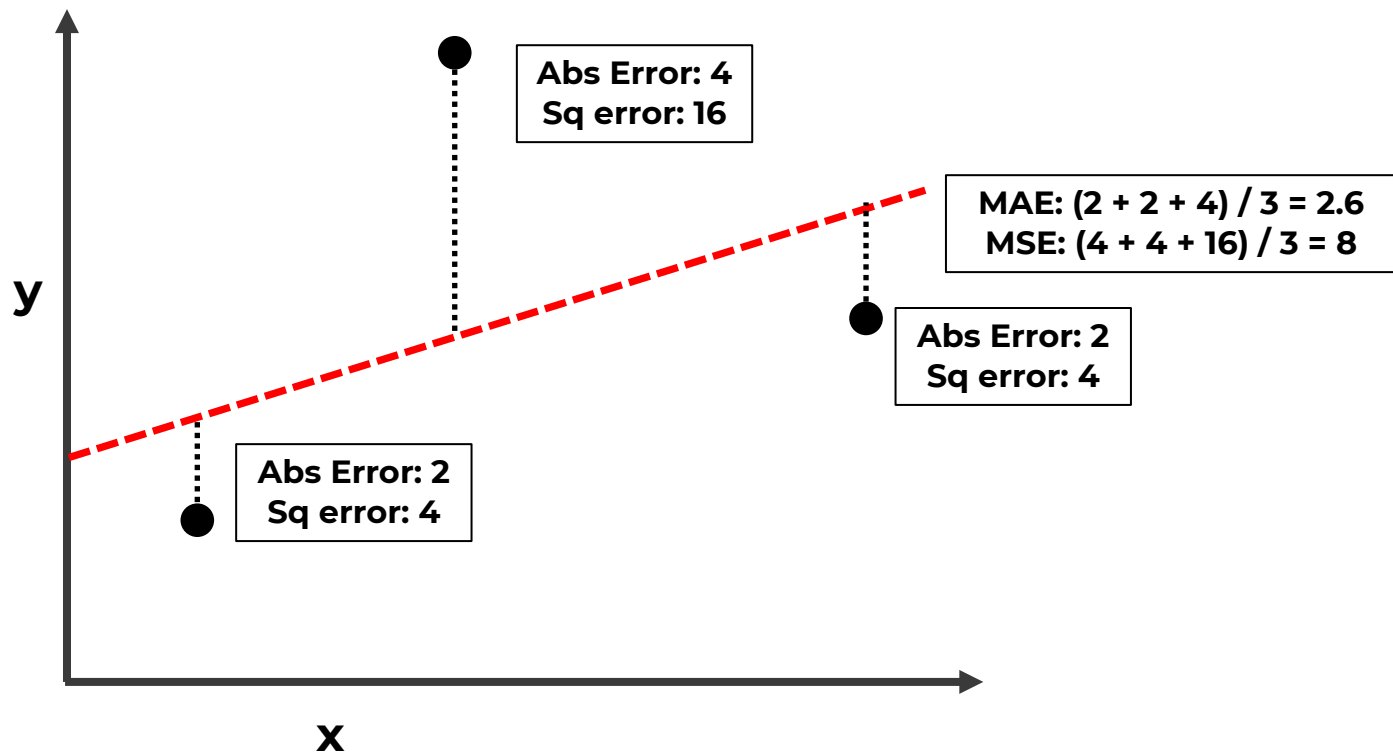
MSE



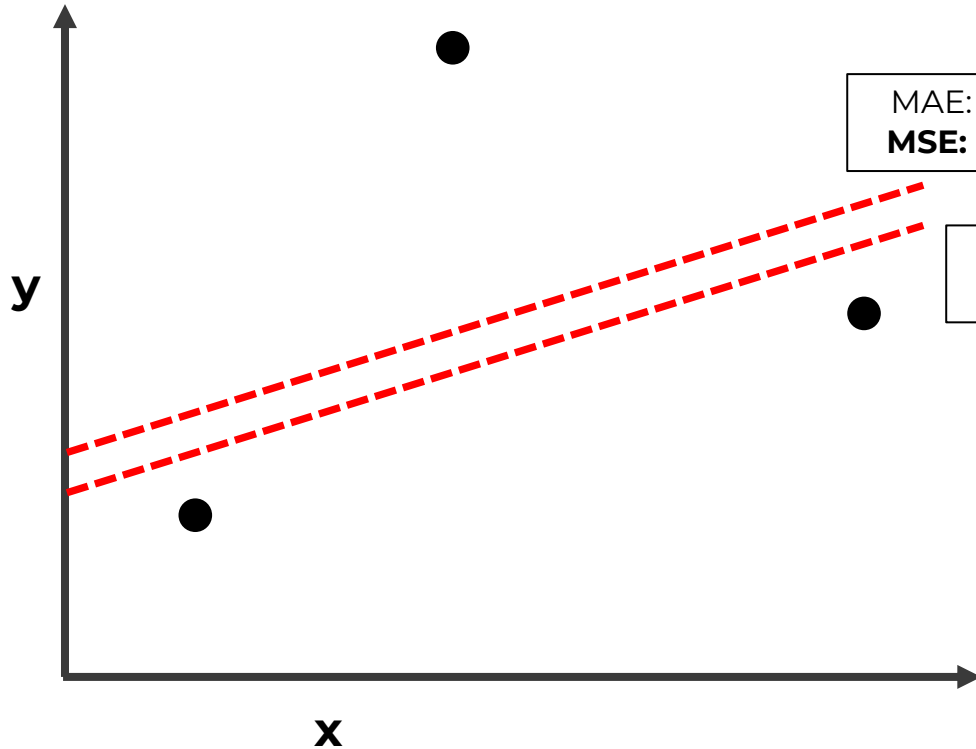
MSE



MSE



MSE



$$\text{MAE: } (2 + 2 + 4) / 3 = 2.6$$

$$\text{MSE: } (4 + 4 + 16) / 3 = 8$$

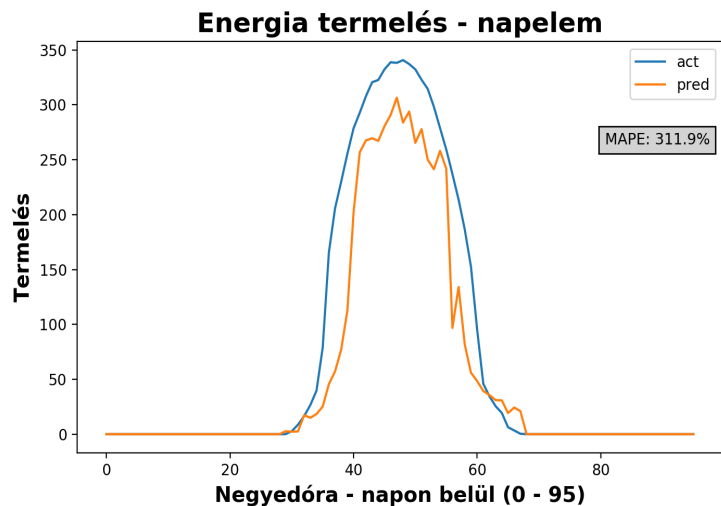
$$\text{MAE: } (1 + 1 + 5) / 3 = 2.3$$

$$\text{MSE: } (1 + 1 + 25) / 3 = 9$$

➤ Megoldásunk függ a célfüggvénytől

MAE vs. MAPE

Célfüggvény: MAE



Célfüggvény: MAPE

