

# SUNS – Zadanie 3:

## Súborové učenie

---

Petra Kirschová

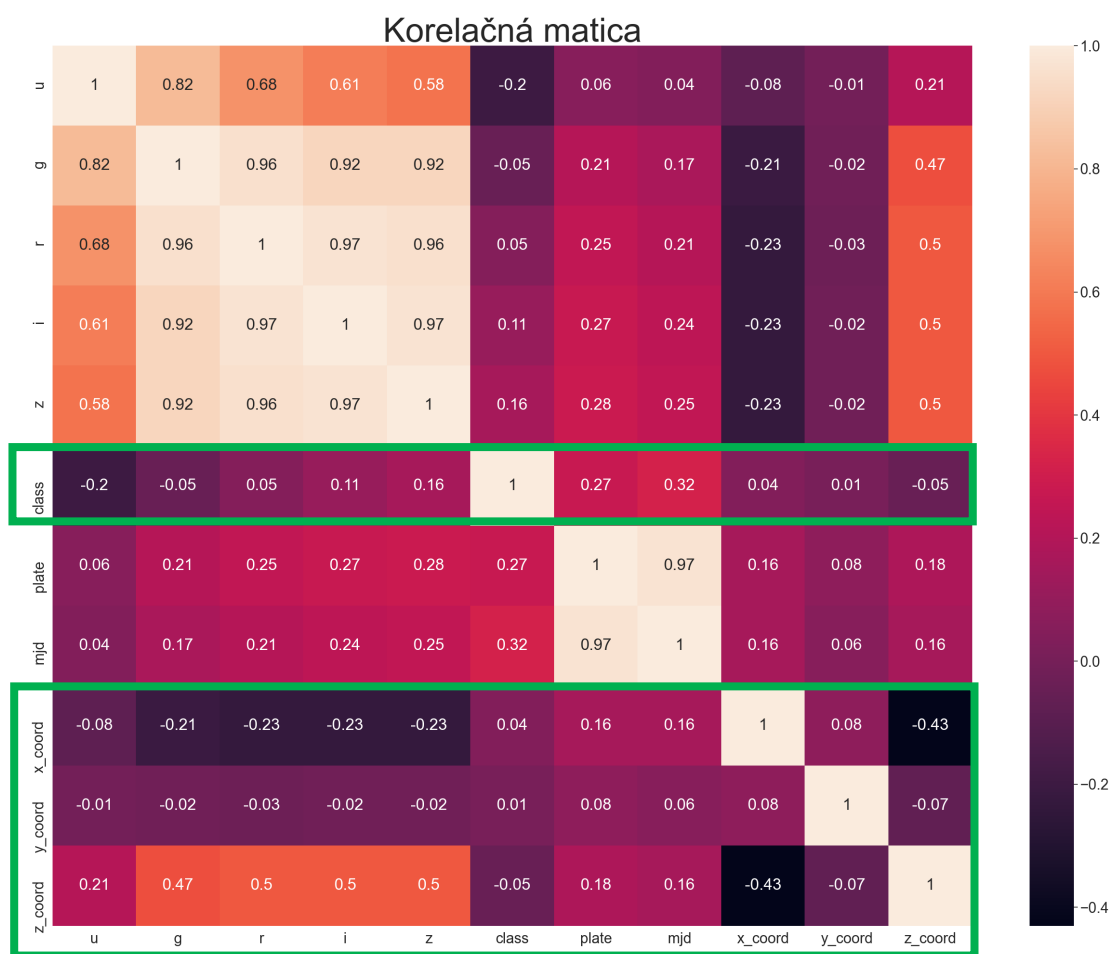
# 1 ANALÝZA A SPRACOVANIE DÁT

Na načítanie a prípravu dát som použila knižnicu `pandas`. Boli k dispozícii zvlášť trénovacie (11740 záznamov) a testovacie dáta (2935 záznamov).

Datasety neobsahovali žiadne null hodnoty a tiež nebolo potrebné mazať outliers, keďže strojové učenie nebolo ovplyvnené prítomnosťou outlierov.

Z datasetov som odstránila stĺpce, ktoré predstavujú identifikáciu záznamov: *objid*, *specobjid*, *run*, *rerun*, *camcol* a *field*.

Po odstránení stĺpcov vyzerala korelačná matica takto:



Na korelačnej matici som sledovala stĺpec *class* a súradnice *x\_coord*, *y\_coord* a *z\_coord*.

Stĺpec *class* má pomerne nízke korelácie s ostatnými ukazovateľmi, no napriek tomu klasifikátor dosahoval veľmi dobré výsledky, keď som na trénovanie použila celý dataset (okrem *class*).

Zo súradníc *x\_coord*, *y\_coord* a *z\_coord* majú *z\_coord* vysoké korelácie takmer so všetkými ostatnými parametrami, hlavne s meraniami *u*, *g*, *r*, *i*, *z*.

## 2 KLASIFIKÁCIA

Zo súborových klasifikátorov som si vybrala random forest klasifikátor a ako ďalší som zvolila MLP klasifikátor. Na ich tréovanie som použila knižnicu `sklearn`, na vykreslenie grafov `seaborn` a `matplotlib` a na vykreslenie a exportovanie stromu `graphviz`.

### 2.1 RANDOM FOREST KLASIFIKÁTOR

Pri random foreste nebolo potrebné dáta normalizovať ani konvertovať na číselné reprezentácie.

Vstupné parametre tréovania sú stĺpce: *u*, *g*, *r*, *i*, *z*, *plate*, *mjd*, *x\_coord*, *y\_coord* a *z\_coord*.

Výstupným parametrom je stĺpec *class*, ktorý obsahuje 3 rôzne triedy objektov: STAR, GALAXY a QSO.

Pre klasifikátor som nastavila:

```
classifier = RandomForestClassifier(criterion='entropy',  
                                  max_leaf_nodes=40, min_samples_leaf=10)
```

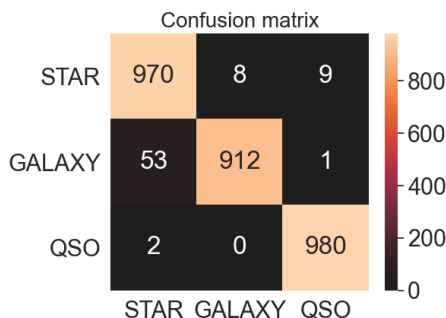
- **criterion='entropy'** – miera správnosti/úspešnosti rozdeľovania dát
- **max\_leaf\_nodes=40** - maximálny počet listov v stromoch
- **min\_samples\_leaf=10** minimálny počet vzoriek v liste

#### 2.1.1 Výsledky

Pri tréovaní random forest klasifikátora som dosiahla veľmi dobré výsledky, ktoré sa pohybovali medzi 97-98%. Na classification report sú uvedené úspešnosti utriedenia do jednotlivých kategórií a tiež celková presnosť klasifikácie.

***** RANDOM FOREST CLASSIFIER *****				
	precision	recall	f1-score	support
GALAXY	0.946	0.983	0.964	987
QSO	0.991	0.944	0.967	966
STAR	0.998	0.998	0.994	982
accuracy			0.975	2935
macro avg	0.976	0.975	0.975	2935
weighted avg	0.976	0.975	0.975	2935

Na confusion matrix je tiež vidno, že tréovanie prebehlo veľmi úspešne. Na diagonále sa nachádzajú vzorky, ktoré boli správne zatriedené do jednotlivých kategórií. Nesprávne zatriedených vzoriek je v porovnaní s veľkosťou datasetu zanedbateľné množstvo.



## 2.2 SLABÝ KLASIFIKÁTOR

Slabým klasifikátorom v random forest modeli je 1 strom. V modeli som nechala defaultný počet slabých klasifikátorov – 100, t.j. celkovo sa vyhodnocuje 100 rozhodovacích stromov.

Jeden uzol stromu obsahuje:

- *Rozhodovaciú podmienku*
- *Entropy* – miera rozdielnosti dát
- *Samples* - počet vzoriek
- *Value* – koľko vzoriek prislúcha jednotlivým triedam
- *Class* – výsledok zatriedenia v danom uzli stromu

Strom na základe rozhodovacej podmienky rozdelí dáta na 2 časti – vľavo pôjdu tie, pre ktoré podmienka platí a vpravo ostatné. Tento krok sa opakuje, kým nie sú splnené zastavovacie podmienky tréovania.

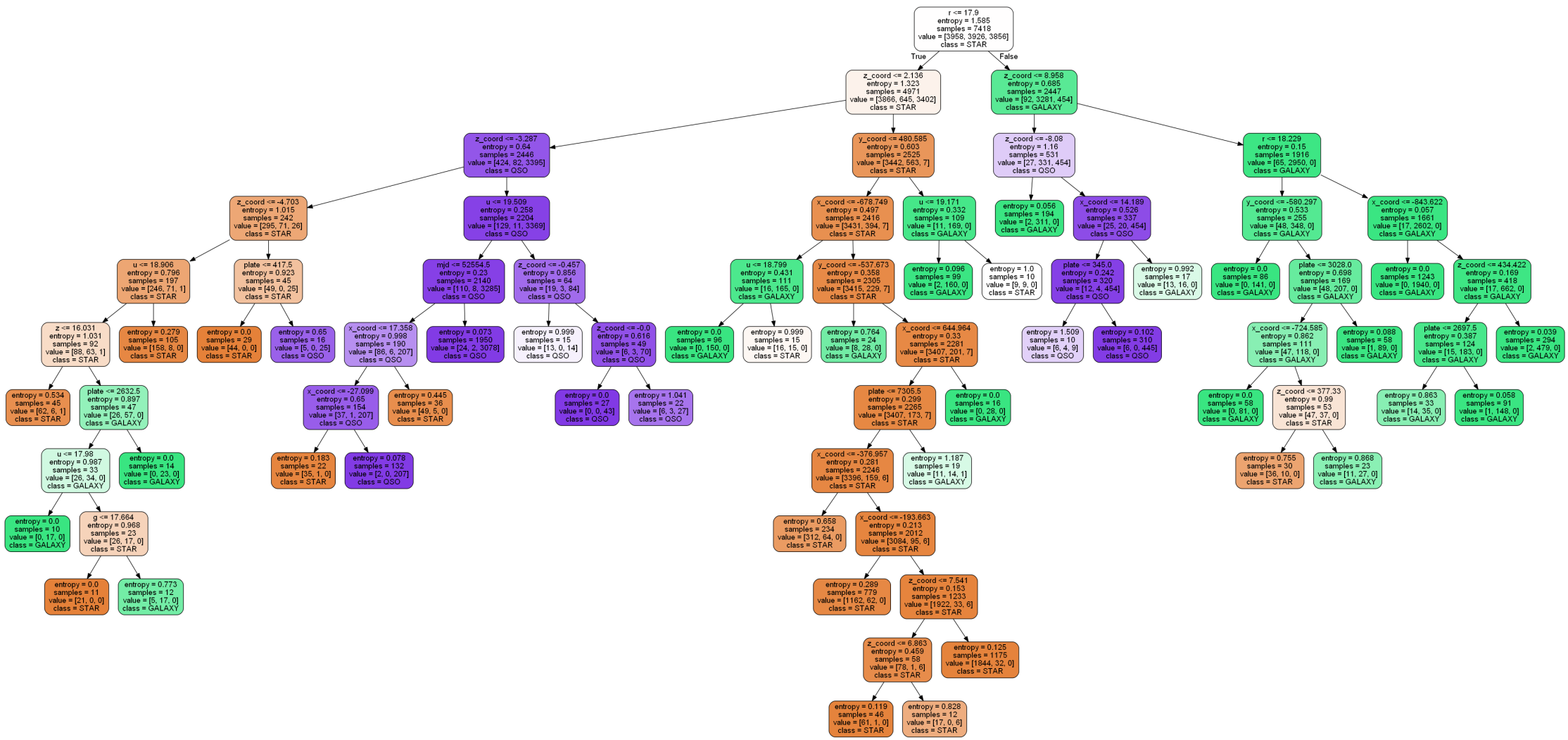
Entropia vyjadruje, ako veľmi sú dáta v danej vzorke navzájom rozdielne. V roote stromu je entropia vysoká, pretože všetky triedy sú takmer rovnako zastúpené. Klesaním nižšie v strome sa entropia znižuje a v niektorých listoch dosiahla aj nulu.



Výslednou triedou v danom uzle stromu je tá trieda, ktorá má vo vzorke najväčšie zastúpenie (najväčšie číslo v poli *value*).

Vykreslený strom je výsledkom pruningu, kvôli tomu, že nenastavenie zastavovacej podmienky viedlo k pretrénovaniu modelu. Pôvodný strom bol tiež príliš veľký a koncové uzly obsahovali málo dát na to, aby boli smerodajné. Pri veľkom strome vo väčšine vetiev nastala situácia, kedy sa ďalším delením stromu nezískala žiadna výpovedná hodnota, keďže sa vzorky v oboch častiach datasetu zaradili do rovnakej triedy.

Po pruningu sa v strome nachádza tiež zopár takých vetiev, ktorých rozdelením dostaneme rovnaký výsledok, no vzhľadom na dosiahnutú úspešnosť je tento model relatívne dobre rozložený. Sprísňovanie zastavovacích podmienok sa potom odzrkadľovalo na celkovej úspešnosti klasifikátora.



## 2.3 MLP KLASIFIKÁTOR

Vstup a výstup MLP klasifikátora bol rovnaký, ako pri random foreste, vstupné dáta sa na rozdiel od random forestu museli normalizovať. Výstupnú množinu, t.j. stĺpec *class* som nechala bez úprav v pôvodnom tvare – tréning fungoval, aj keď som vo výstupnej množine nechala triedy vo formáte string.

Nastavenie klasifikátora:

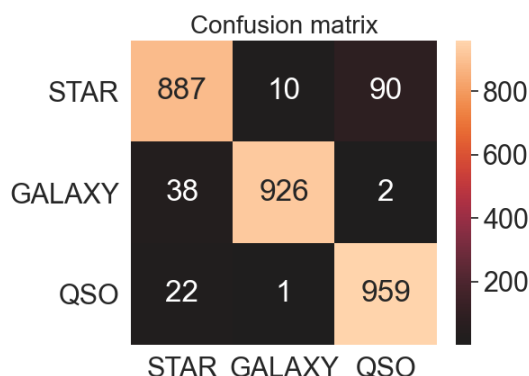
```
classifier = MLPClassifier(max_iter=600, alpha=0.01,  
                           hidden_layer_sizes=(20,))
```

- aktivačná funkcia=**relu** (default)
- solver=**adam** (default)
- **max\_iter=600** - maximálny počet iterácií
- **alpha= 0.01** – regularizačná penalta
- **hidden\_layer\_sizes=(20,)** - 1 skrytá vrstva s 20 neurónmi

### 2.3.1 Výsledky

MLP klasifikátor dosiahol trochu nižšie úspešnosti, ako random forest, ale stále sú veľmi vysoké. Celkovo sa presnosti MLP klasifikátora pohybovali okolo 93-94%.

***** MLP CLASSIFIER *****				
	precision	recall	f1-score	support
GALAXY	0.937	0.899	0.917	987
QSO	0.988	0.959	0.973	966
STAR	0.912	0.977	0.943	982
accuracy			0.944	2935
macro avg	0.946	0.945	0.945	2935
weighted avg	0.946	0.944	0.944	2935



## 3 REGRESIA

Spomedzi súborových klasifikátorov som na regresiu použila random forest regresor a ako ďalší regresor som si vybrala k-nearest neighbors regresor. Regresory som implementovala pomocou knižnice `sklearn` a na vykreslenie grafov som použila `seaborn`, `matplotlib` a `plotly`.

### 3.1 RANDOM FOREST REGRESOR

Pri random forest regresore nebolo potrebné dáta normalizovať, ale vstupný stĺpec `class` bolo treba konvertovať na číselné reprezentácie tried.

Vstupné parametre trénovania sú stĺpce: `u`, `g`, `r`, `i`, `z`, `plate`, `mjd` a `class`.

Výstupnými parametrami sú stĺpce so súradnicami `x_coord`, `y_coord` a `z_coord`, ktoré určujú polohu sledovaného objektu vo vesmíre.

Nastavenie regresora:

```
regressor = RandomForestRegressor(criterion="mse", max_depth=15)
```

- **criterion="mse"** – miera správnosti rozdelenia dát = mean squared error
- **max\_depth=15** – maximálna hĺbka stromu, zastavovacia podmienka

#### 3.1.1 Výsledky

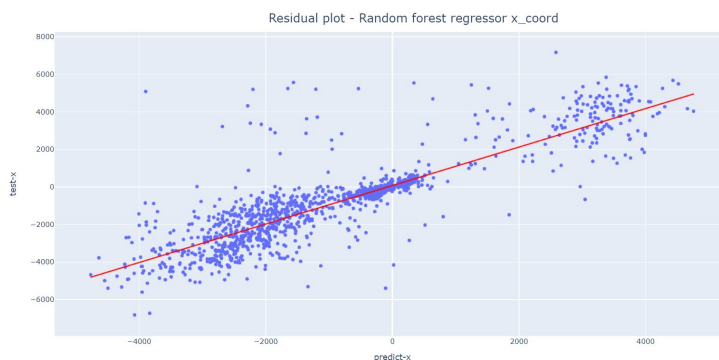
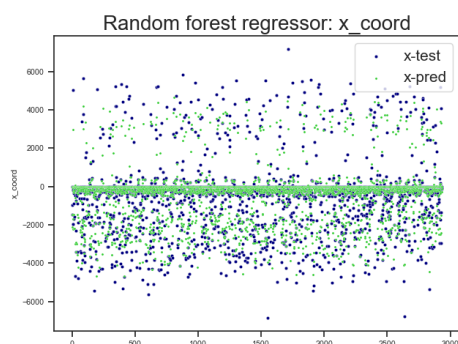
Random forest regresor dosahoval dosť dobré výsledky, dosiahol R2 score pre celý výstup okolo 73% a mean absolute error cca 280. Pre jednotlivé súradnice bola pre `x_coord` najvyššia chyba a pre `z_coord` najnižšia. R2 score bolo najvyššie pre `x_coord` a najnižšie pre `y_coord`.

```
***** RANDOM FOREST REGRESSOR *****
[x_coord,y_coord,z_coord]: MAE = 280.0764288871531 R2 = 0.7265198924894873

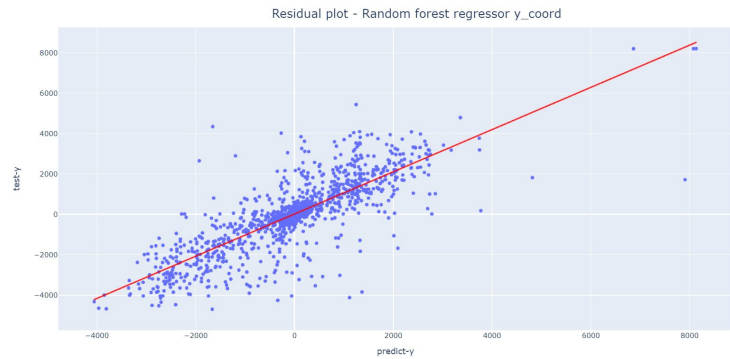
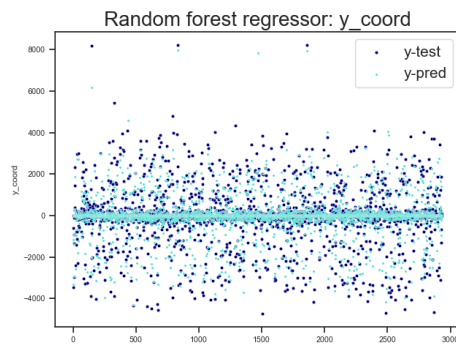
x_coord: MAE= 324.1265951496431 R2= 0.7641190247821393
y_coord: MAE= 291.8999896621916 R2= 0.6537591774454663
z_coord: MAE= 224.20270184962345 R2= 0.7616814752408566
```

V 2D a 3D grafoch som porovnala výstup testovacích a predikovaných dát. Na grafoch je vidno, že regresor priblížil výstupné hodnoty relatívne dobre, pozície predpovedaných bodov sú veľmi blízko k reálnym bodom. Na 2D grafoch sú porovnané predpovedané a očakávané výstupy (v závislosti od id alebo pri residual plote priamo závislosť medzi predikovaným a očakávaným výstupom) pre každú súradnicu zvlášť. 3D graf obsahuje body `[x_coord, y_coord, z_coord]` v trojrozmernom priestore.

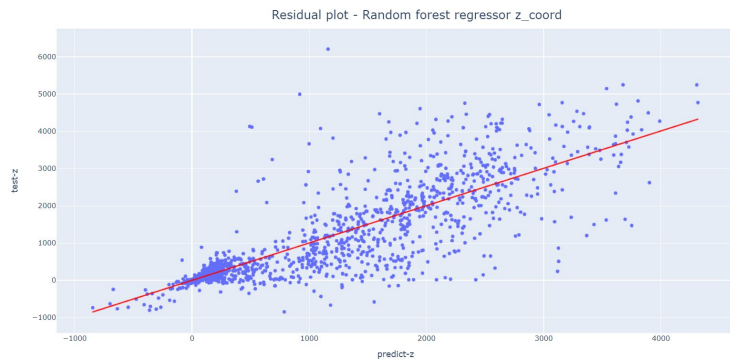
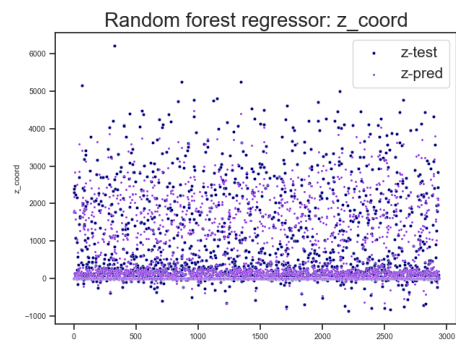
#### x\_coord: porovnanie očakávaných hodnôt s predpovedanými + residual plot



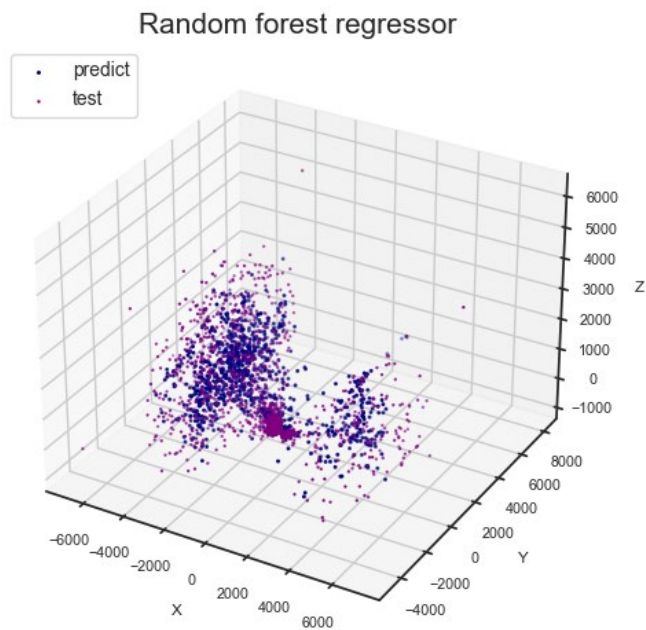
### y\_coord: porovnanie očakávaných hodnôt s predpovedanými + residual plot



### z\_coord: porovnanie očakávaných hodnôt s predpovedanými + residual plot



### [x\_coord, y\_coord, z\_coord]: porovnanie očakávaných hodnôt s predpovedanými





## 3.2 K-NEAREST NEIGHBORS REGRESOR

K-nearest neighbors regresor predpovedá výsledné hodnoty na základe toho, aké hodnoty dosahujú iné dáta v jeho okolí (neighbors). Dáta tiež nebolo treba normalizovať, ale vstupný stĺpec *class* sa konvertoval na číselné reprezentácie tried.

Nastavenie regresora:

```
regressor = KNeighborsRegressor(n_neighbors=8, weights='distance')
```

- **n\_neighbors=8** – počet „susedov“, podľa ktorých sa vyhodnocuje výsledok
- **weights='distance'** – váhy susedov sa určujú podľa ich vzdialenosti k sledovanému bodu, t.j. body, ktoré sú bližšie viac ovplyvnia výsledok, ako tie, ktoré sú ďalej.

### 3.2.1 Výsledky

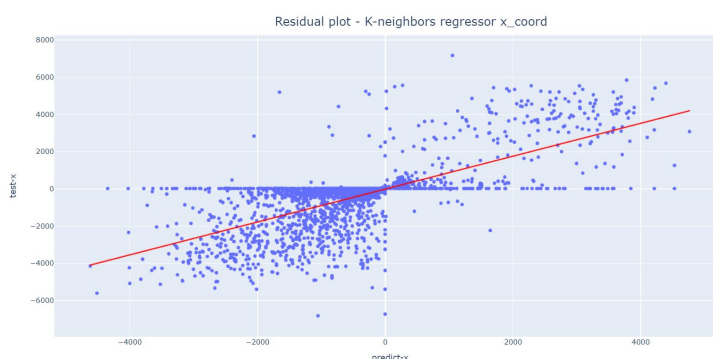
Výsledky KNN regresora boli horšie, ako pri random forest, R2 bolo celkovo okolo 0.5, MAE cca 490, pre jednotlivé súradnice sa chyby pohybovali približne medzi 400-700, R2 score medzi 0.4-0.6.

```
***** K-NEIGHBORS REGRESSOR *****
[x_coord,y_coord,z_coord]: MAE = 491.929882806797 R2 = 0.5040752458655859

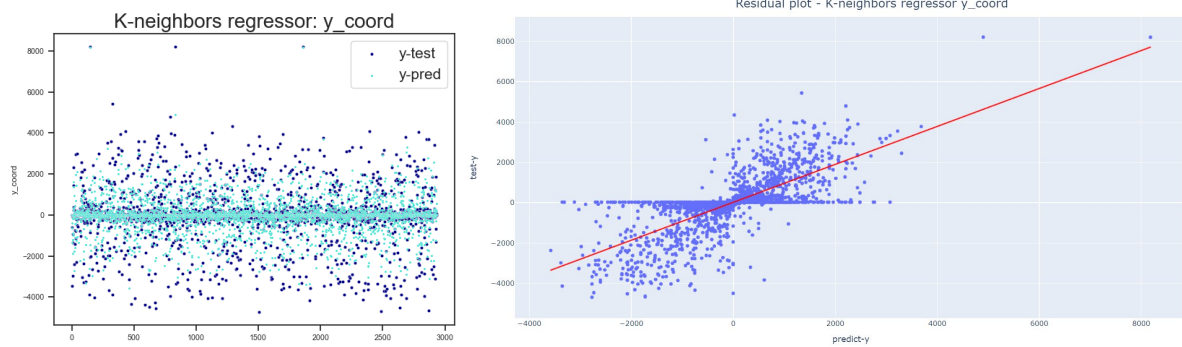
x_coord: MAE= 656.3278901376939 R2= 0.499784264419148
y_coord: MAE= 414.38610080972256 R2= 0.5721132709708907
z_coord: MAE= 405.07565747297576 R2= 0.4403282021839513
```

Grafy, na ktorých sú vykreslené pôvodné a predpovedané výstupy sú na prvý pohľad veľmi podobné s grafmi random forest regresora, rozdiely sú pozorovateľné na residual grafoch. Pri tomto regresore sú predpovedané body rozptýlené viac, ako pri random foreste a najviac sa ich zoskupuje pozdĺž nuly.

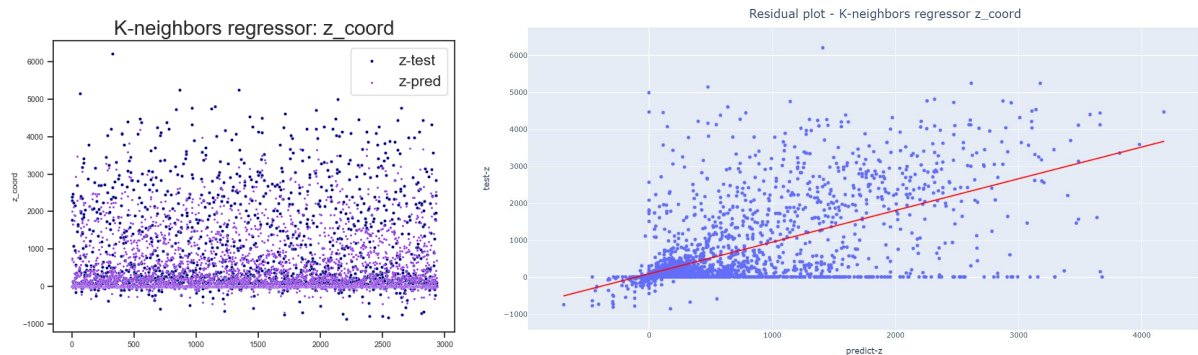
#### x\_coord: porovnanie očakávaných hodnôt s predpovedanými + residual plot



### y\_coord: porovnanie očakávaných hodnôt s predpovedanými + residual plot



### z\_coord: porovnanie očakávaných hodnôt s predpovedanými + residual plot



### [x\_coord, y\_coord, z\_coord]: porovnanie očakávaných hodnôt s predpovedanými

