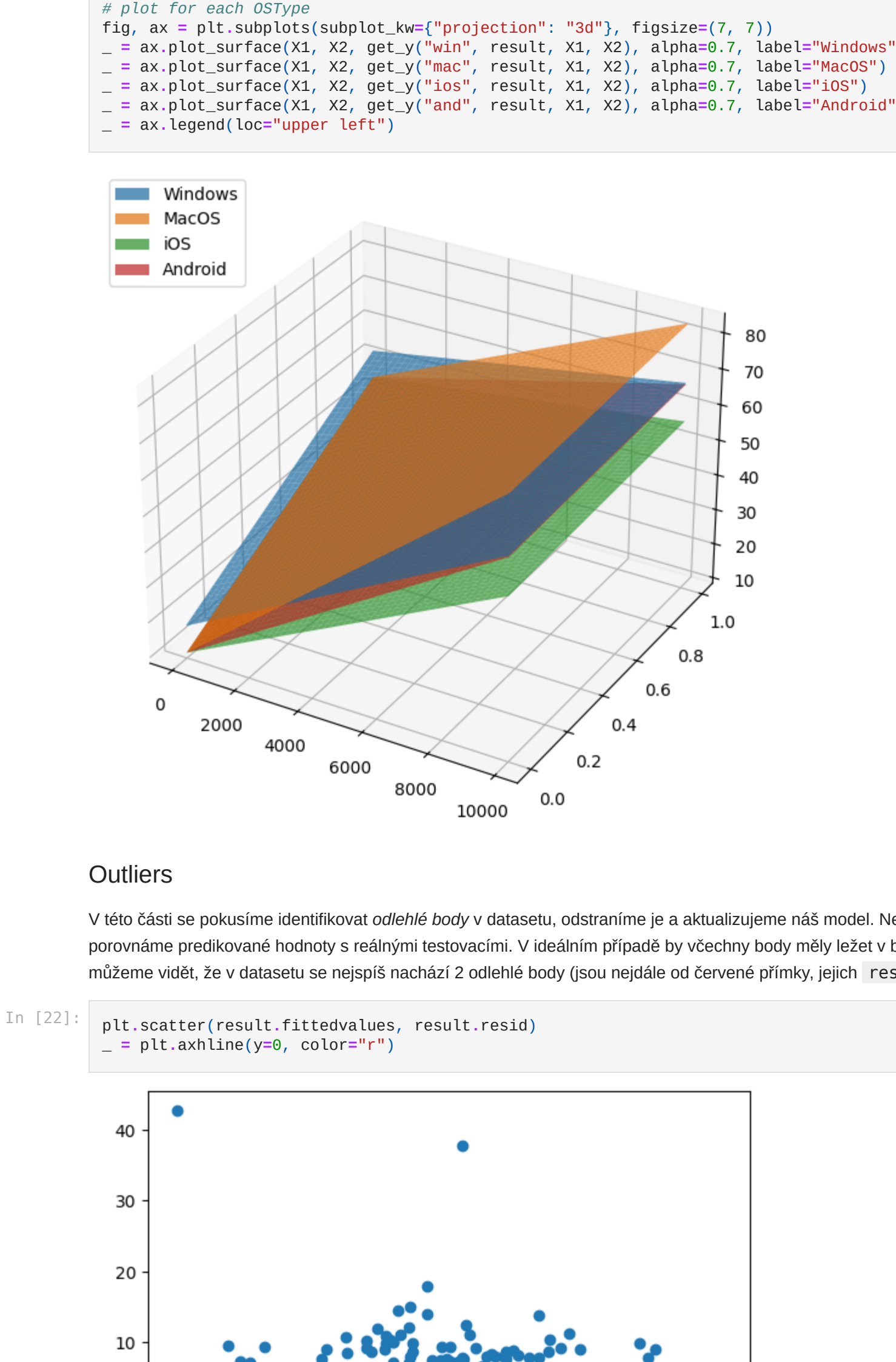
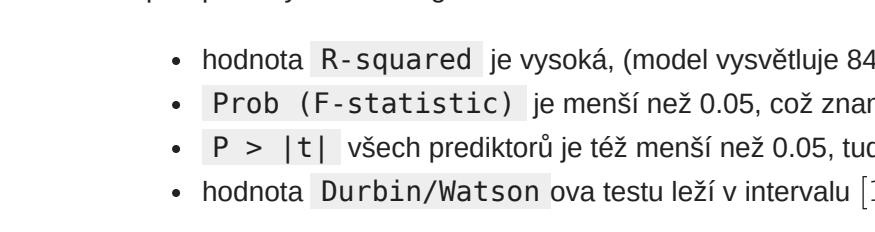


Můžeme si vykreslit graf pro jednotlivé operační systémy



## Outliers

V této části se pokusíme identifikovat odlehle body v datasetu, odstraníme je a aktualizujeme náš model. Nejprve si vykreslíme graf, kde provedeme predikční hodnoty s reálnými testovacími. V jedním případě by všechny body měly ležet v blízkosti červené přímky. Zde můžeme vidět, že v datasetu se nejspíš nachází 2 odlehle body (jsou nejedle od červené přímky, jejich `resid` je nad 30).



Pomocí funkce `outlier_test` testujeme odlehle body na našich pozorovaných datech. Jako výchozí metoda je použita **Bonferroniho korekce** na hladině spolehlivosti 95 %. Data, u kterých vyjde jako výsledek této korelace hodnota menší než 0.05 budeme brát za odlehlou hodnotu. Z testu získáme indexy těchto dat a z datasetu je odstraníme. Následně znovu otestujeme náš model.

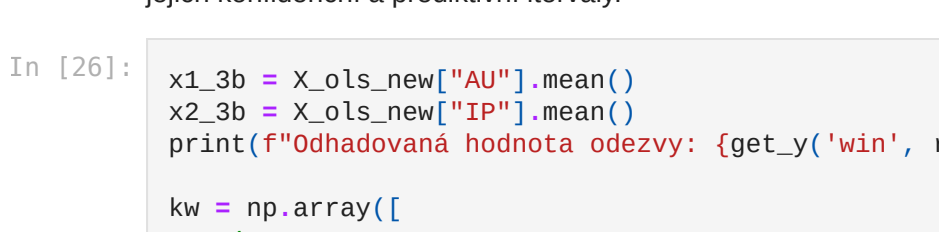
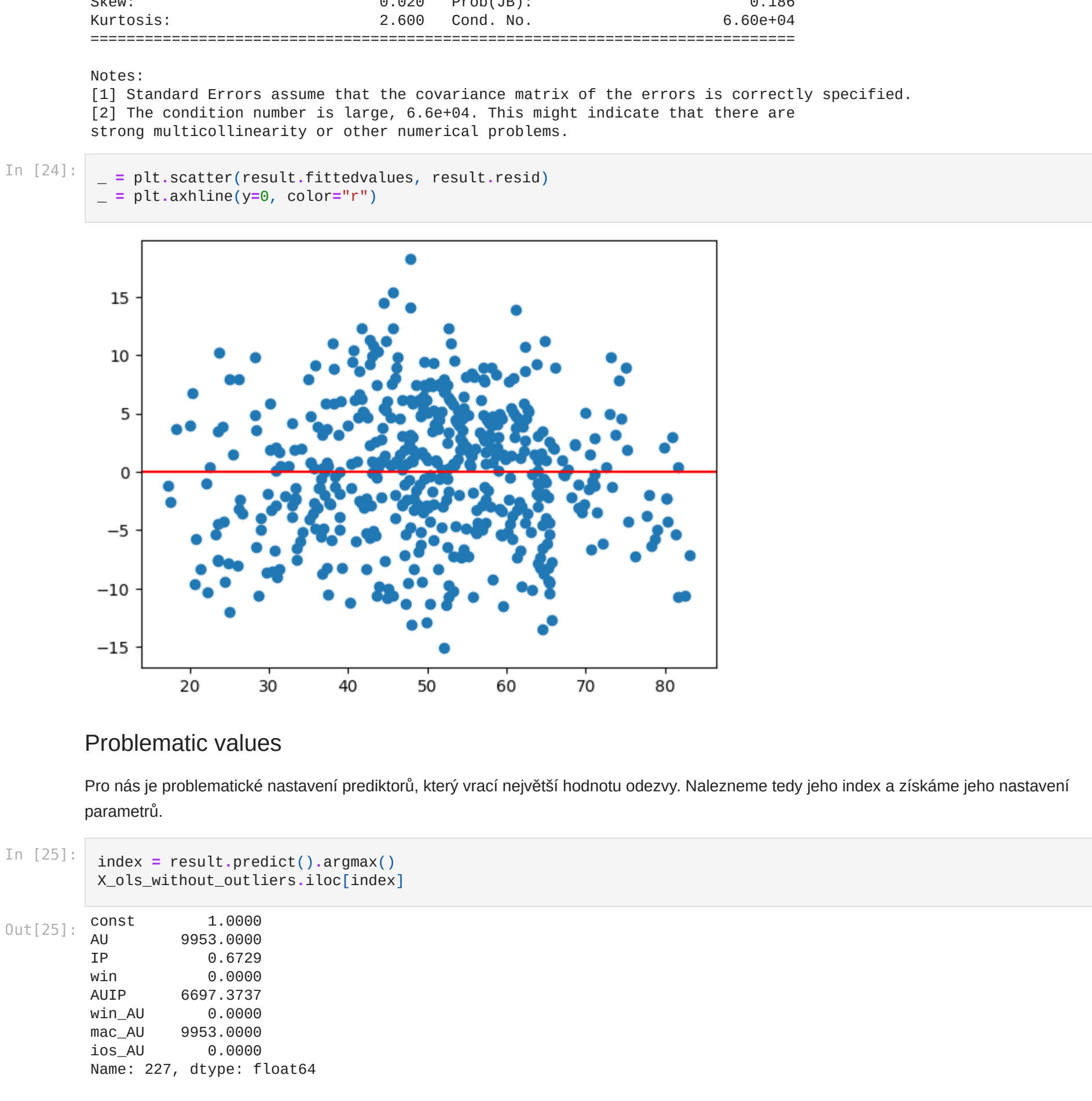
Nový model má lepší hodnotu `R-squared` (z 0.814 na 0.842). Upravený model můžeme popsat touto rovnicí:

$$y = 8.5942 + 0.0057x_{AU} + 36.6538x_{IP} + 8.1829C_{win} - 0.0035x_{AU}x_{IP} - 0.0008C_{win}x_{AU} + 0.0017C_{mac}x_{AU} - 0.0010C_{ios}x_{AU}$$

Nyní zkontrolujeme náš model. Pokud budeme analyzovat výpis informací o modelu, dospějeme k závěru, že náš model splňuje předpoklady lineární regrese.

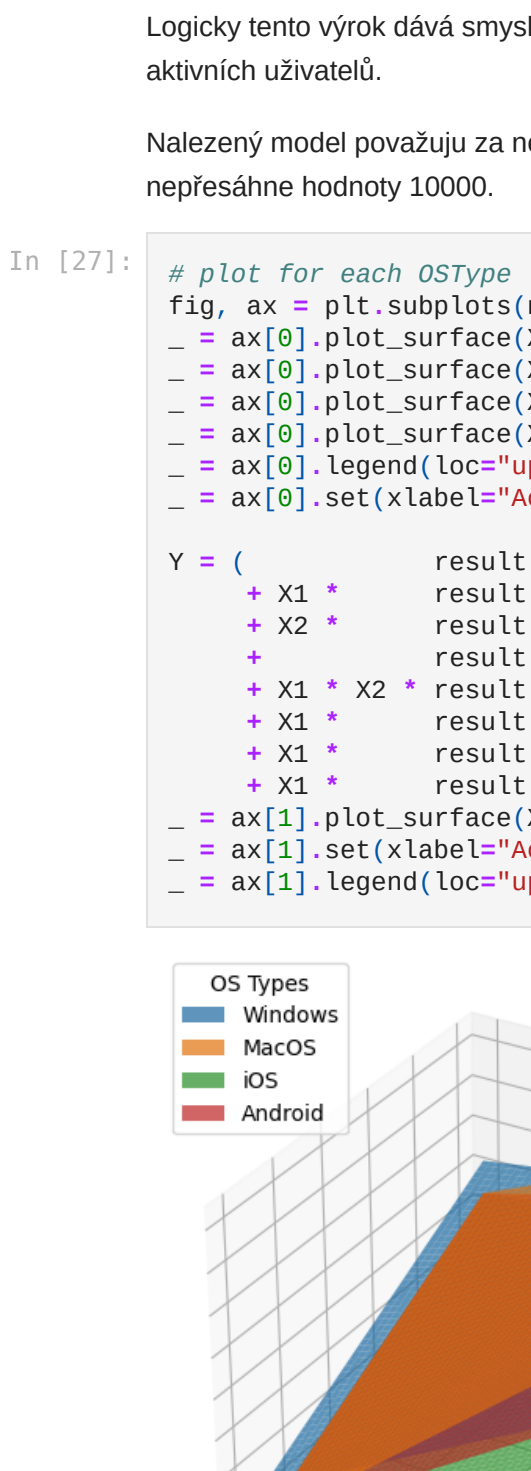
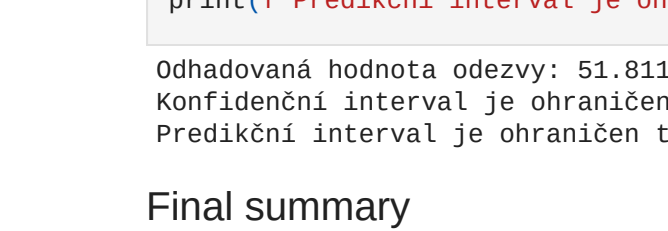
- hodnota `R-squared` je vysoká, (model vysvětluje 84.2 % změn hodnoty odezvy)
- `Prob (F-statistic)` je menší než 0.05, což znamená, že náš výběr prediktorů je statisticky signifikantní
- `P > |t|` všech prediktorů je větší než 0.05, tudíž každý prediktor v modelu má vliv na výpočet výsledné hodnoty odezvy
- hodnota `Durbin/Watson` ova testu leží v intervalu `[1, 2]`, tudíž neexistuje autokorelace mezi sousedními rezidui

Model je samozřejmě deformovaný, protože jsme neprovedli standardizaci. Po standardizaci by nám vyšel stejný model s různými koeficienty a významnostmi.



## Problematic values

Pro nás je problematické nastavení prediktorů, který vrací největší hodnotu odezvy. Nalezneme tedy jeho index a získáme jeho nastavení parametru.



## Predict a value

Nyní se pokusíme odhadnout hodnotu odezvy pro uživatele `Windows`, při průměrném nastavení ostatních parametruů. Následně vypíšeme i jejich konfidenční a predikční intervaly.



## Final summary

Jak již bylo řečeno v předchozích buňkách, výpis informací modelů nám ukazuje, že máme dobrý model. Podíváme-li se na grafy pro jednotlivé operační systémy, můžeme z něj vyčíst, že s rostoucím počtem **aktivních uživatelů** roste i predikovaná **hodnota odezvy**. Logicky tento výrok dává smysl, nicméně ne vždy to musí platit. Hodnota `InteractingPct` má větší vliv na odezvu jen při malém počtu aktivních uživatelů.

Nalezený model považují za nejspíš vhodným a můžeme podle něj predikovat hodnoty odezvy, pokud počet aktivních uživatelů nepřesáhne hodnoty 10000.

