

Indeks	Punkty	Uwagi - sprawozdanie	Punkty (prezentacja)	Uwagi - prezentacja
258544	5	<p>Jeśli mamy dane w złej skali, to łatwo wiele obserwacji uznać za outliery. Duży plus za znalezienie błędu w danych!</p>	3	<p>"Opis danych - ok Dokładny opis z czego korzystamy - ok.</p> <p>Język matematyczno-statystyczny:</p> <p>Szczegóły podział na 3 części - niepotrzebne. Przekształcenie danych (symetryzacja danych) Skośność, Kurtoza (na rysunku) Przycięcie zmiennych</p> <p>Uwaga: Można dać ładne etykiety na przekształconej skali</p> <p>Uproszczenie zmiennych - bardzo dobre. Ale czy należy to tłumaczyć osobom, które tego nie rozumieją?</p> <p>Bardzo dużo szczegółów technicznego dochodzenia do rozwiązania. A nas interesuje:</p> <ul style="list-style-type: none"> <li>* jak wygląda model, który został wybrany</li> <li>* jak należy go rozumieć</li> <li>* na ile da się przewidywać liczbę udostępnień</li> </ul> <p>Bardzo dobry wykres 'roznica predykcji i rzeczywistych' vs 'indeksy w sortowaniu' Jak poprawić -&gt; dobre pomysły.</p>
258486, 258	5	<p>Przy wypisywaniu tabel warto użyć pakietu xtable lub knitr::kable</p>	5	<p>"Cel projektu -&gt; w miarę jasno opisane</p> <p>Opis danych: krótko jakie są grupy zmiennych</p> <p>5 Eksploracja: identyfikacja potencjalnych problemów z danymi.</p>
258478	4	<p>Dobre uzasadnienie czemu bierzemy skalę logarytmiczną Co do wyboru ostatecznego modelu: co prawda <math>R^2</math> jest niskie w obu przypadkach 0.08 i 0.12, ale to jest jednak różnica prawie 50%! Brak diagnostyki modelu</p>	4	<p>"Początek dobry: nie ma nic super, ale mamy prosty model, który trochę tłumaczy.</p> <p>Bardzo dobre wytłumaczenie co się dzieje w modelu. Brak opisu danych.</p> <p>Opis - kilka hitów, ale większość średniaków. Przejście do innej skali (logarytm to zgrzyt, ale malutki)</p> <p>Mini minus - warningi w chunkach warning=FALSE, message=FALSE</p> <p>Jeśli mamy średnie dla dwóch grup - można sprawdzić czy warto w ogóle raportować</p>

258554	4	<p>Podział zbioru na treningowy, walidacyjny i testowy powinien nastąpić PRZED eksploracją danych!          Usunięcie obserwacji z bardzo dużą liczbą udostępnień jest ok. Usunięcie tych z małą liczbą może być niebezpieczne. Warto sprawdzić jak model zachowuje się dla takich obserwacji          Usunięcie obserwacji 29088 ok          Jeśli chodzi o usunięcie obserwacji z dużą liczbą linków - jest to nieco pochopne. Dlaczego one nam przeszkadzają? Może lepiej zastosować transformację? Połączenie poziomów zmiennej channel ok</p>		<p>"Początek: Ok. Co robimy, co możemy dostać.          Opis danych - bardzo dobrze.          Analiza udostępnień - co za tym stoi - ciekawe.          Dużo niepotrzebnych zmiennych - będziemy się chcieli ich pozbyć.</p> <p>Czyszczenie danych: obserwacje zbyt zróżnicowane (może mogłaby pomóc transformacja?),          usuwane są obserwacje odstające (uwaga, czy to nie są przypadkiem te najbardziej popularne?)</p> <p>5 Troszkę przesada z łączeniem grup. Wchodzimy w szczegóły - w jakim celu.</p>
258512, 258	4.5	<p>Podział zbioru na treningowy, walidacyjny i testowy powinien nastąpić PRZED eksploracją danych!          Bardzo fajny wykres shares vs all          Usunięcie obserwacji ze względu na błędne wartości n_tokens_title bardzo ok          Właściwie zamiast korelacji między zmiennymi objaśniającymi powinniśmy mierzyć Variance Inflation Factor (VIF), który mówi w jakim stopniu zmienna jest tłumaczona liniowo przez wszystkie pozostałe          Co do modelu model1_lm_log - widać niespełnienie założeń. I widać też, że powodem tak naprawdę jest kilka obserwacji o bardzo dużej liczbie udostępnień          Bardzo fajny pomysł z porównaniem artykułów parami</p>		<p>"Początek ok. Cel, co robimy i po co.</p> <p>Opis danych - z czym mamy do czynienia, grupy zmiennych.</p> <p>Jak wygląda zmienna objasniata - mamy bardzo skrajne różnice w udostępnieniach.</p> <p>Większość zmiennych nie ma związku - pokazujemy to co ma wpływ na liczbę udostępnień.</p> <p>Tytuły -&gt; średnia długość tytułu daje szansę na dużo sharów</p> <p>Fajna idea odnośnie porównań między dwoma artykułami. Ale przy posortowaniu mamy problem wielokrotnego testowania.</p> <p>Wnioski też ok -&gt; zamiast przewidywać liczbę sharów czy jest popularny czy nie.</p> <p>5</p>