

# Tytuł raportu

## Badanie zależności liniowej oceny czekolady przez krytyków w zależności od jej składu

### Wstęp

Wybrałem ten temat, ponieważ:

- bardzo lubię czekoladę
- chciałem dowiedzieć się, zawartość jakich składników ma największy wpływ na ocenę czekolady
- chciałem sprawdzić, czy ocena czekolady jest funkcją (liniową) składu

### Opis danych - struktura zbiorów, opis zmiennych, pochodzenie

Zbiór danych zawiera następujące zmienne:

- ref - unikalny dla każdej firmy numer
- company - nazwa firmy
- company\_location - lokalizacja firmy
- review\_date - data recenzji czekolady
- country\_of\_bean\_origin - kraj pochodzenia ziaren kakao
- specific\_bean\_origin\_or\_bar\_name - dokładniejsza lokalizacja pochodzenia ziaren kakao
- cocoa\_percent - zawartość procentowa kakao
- rating - ocena
- counts\_of\_ingredients - ilość składników
- beans - czy czekolada zawiera ziarna kakao?
- cocoa\_butter - czy czekolada zawiera masło kakaowe?
- vanilla - czy czekolada zawiera wanilię?
- lecithin - czy czekolada zawiera lecytynę?
- salt - czy czekolada zawiera sól
- sugar - czy czekolada zawiera cukier
- sweetener\_without\_sugar - czy czekolada zawiera słodzik
- first\_taste - smak czekolady przy pierwszej próbie
- second\_taste - smak czekolady przy drugiej próbie
- third\_taste - smak czekolady przy trzeciej próbie
- fourth\_taste - smak czekolady przy czwartej próbie

Dane pochodzą z następującej strony: <https://www.kaggle.com/soroushghaderi/chocolate-bar-2020>

# Opis procesu przygotowywania danych do analizy - kolejne kroki

W pierwszej kolejności odrzucam zmienne:

- ref
- company
- company\_location
- review\_date
- country\_of\_bean\_origin
- specific\_bean\_origin\_or\_bar\_name
- first\_taste
- second\_taste
- third\_taste
- fourth\_taste

Celem tego zabiegu jest uniezależnienie modelu od zmiennych niezaliczających się do składu czekolady. Zmienne takie, jak nazwa firmy, oraz data recenzji nie powinny mieć wpływu na wynik.

**Następnie sprawdzam ilość unikalnych wartości, w celu wyeliminowania zbędnych zmiennych.**

```
cocoa_percent: [ 42.  46.  50.  53.  55.  56.  57.  58.  60.  60.5  61.  62.
 63.  64.  65.  66.  67.  68.  69.  70.  71.  71.5  72.  72.5
 73.  73.5  74.  75.  76.  77.  78.  79.  80.  81.  82.  83.
 84.  85.  86.  87.  88.  89.  90.  91.  99. 100. ]
rating: [1.  1.5  1.75 2.  2.25 2.5  2.6  2.75 3.  3.25 3.5  3.75 4. ]
counts_of_ingredients: [1 2 3 4 5 6]
beans: ['have_bean']
cocoa_butter: ['have_cocoa_butter' 'have_not_cocoa_butter']
vanilla: ['have_not_vanilla' 'have_vanilla']
lecithin: ['have_lecithin' 'have_not_lecithin']
salt: ['have_not_salt' 'have_salt']
sugar: ['have_not_sugar' 'have_sugar']
sweetener_without_sugar: ['have_not_sweetener_without_sugar'
'have_sweetener_without_sugar']
```

Na podstawie powyższych wyników postanowiłem, że: - należy usunąć zmienną beans, ponieważ posiada tylko jedną unikalną wartość - zmienne coco\_butter, vanilla, lecithin, ssalt, sugar i sweetener\_without\_sugar to tak na prawdę zmienne logiczne i powinny zostać na takowe zamienione - zmienna cocoa\_percent może zostać znormalizowana

## Analiza danych - przyjęte założenia, krótki opis metod i obranej metodologii analizy

Zbadana została korelacja liniowa między zmiennymi - zmienna counts\_of\_ingredients wykazała dużą korelację ze zmiennymi vanilla, lecithin i cocoa\_butter, co może oznaczać, że uwzględnia ona zawartość tych składników - zmienne sugar, oraz sweetener\_without\_sugar

wykazują dużą (ok. -0.96) korelację ujemną

Postanowiłem w związku z tym usunąć zmienną `counts_of_ingredients`, oraz zostawić na razie zmienne zawierające informacje o cukrze, oraz słodziku, ponieważ zastąpienie cukru za pomocą słodzika może mieć wpływ na ocenę, ponieważ słodziki często smakują inaczej (np. stewia).

## **Modelowanie danych - przyjęte założenia, krótki opis metod i obranej metodologii budowania modeli**

- Wykorzystany został model regresji liniowej.
- dane po przygotowaniu wykazały słabe rezultaty
- w następnej kolejności wystąpiła próba sprawdzenia, czy zaokrąglenie w dół wartości oceny polepszy model
- zmienne zawartości cukru i słodzika były mocno (ujemnie) skorelowane, w związku z czym wyznaczyłem na ich podstawie zmienną mówiącą, czy czekolada zawiera słodzik lub cukier.
- ostatecznie usunięta zmienna ta została usunięta, by sprawdzić jaki wpływ będzie to miało na model

## **Rezultaty, wnioski i ich dyskusja**

- Składniki nie wykazały znacznej korelacji ze zmienną *rating*, najwyższy wynik (wartość bezwzględna) uzyskała zmienna *vanilla* (ok. -0.1).
- Model regresji liniowej nie był w stanie strafnie przewidzieć oceny krytyka na podstawie składu czekolady.
- Zmiany w modelu pogarszały jedynie jego skuteczność.