

Previsão de Notas de Review (Olist)

Introdução à Ciência de Dados — NES

Tiago Cavalcante Trindade

Julho de 2025

- Prever a **review_score** (1-5 estrelas) antes do cliente avaliar.
- Base de dados pública da **Olist**.
- **Features**
 - Numéricas: tempo de entrega, valor do carrinho, frete, nº de itens, pagamentos, dia da semana, mês.
 - Texto: título + mensagem da review → vetor FastText spaCy (300 dim) → SVD (120 dim).
- Balanceamento por undersampling; divisão temporal 80/20 (treino/teste).

Modelos Avaliados

- ① Regressão Linear (arredondada & cortada)
- ② Regressão Logística Multiclasse
- ③ Regressão Logística Ordinal

Implementados com `scikit-learn` (+ `mord` para o modelo ordinal).

Matrizes de Confusão

Linear

True \ Pred	1	2	3	4	5
1	41	256	156	39	0
2	19	222	195	107	2
3	12	109	195	235	12
4	1	20	157	425	36
5	1	9	138	492	62

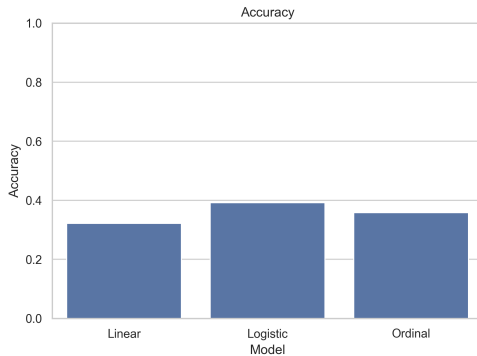
Logistic

True \ Pred	1	2	3	4	5
1	231	152	47	34	28
2	144	180	76	69	76
3	55	110	85	155	158
4	11	33	59	196	340
5	7	18	59	158	460

Ordinal

True \ Pred	1	2	3	4	5
1	123	235	78	55	1
2	67	226	122	125	5
3	24	120	131	268	20
4	6	30	97	437	69
5	1	15	87	464	135

Acurácia por Modelo



Modelo	Acurácia	MAE	RMSE
Regressão Linear (arred.)	0.321	0.917	1.105
Regr. Logística Multiclasse	0.391	0.882	1.263
Regr. Logística Ordinal	0.358	0.823	1.115

- Logística Multiclasse: melhor acerto exato (aproximadamente 40%).
- Logística Ordinal: menor MAE e RMSE — erra geralmente só ± 1 estrela.
- Regressão Linear é apenas um baseline simples.

Escolha

- Se o KPI exige acerto exato da nota \Rightarrow **Logística Multiclasse.**
- Se o mais importante é limitar a distância do erro \Rightarrow **Logística Ordinal.**