ENCODER (stringlarni raqamlarga o'tkazish)

- One Hot Encoding > One Hot Encoding va Get dummies
- Label Encoding
- Target Encoding
- Frequence Encoding

Encoder turlari va ishlatilishi

1. Label Encoding (Ordinal Encoding - tartibli)(Alphabet)

- Qanday ishlaydi: Har bir kategoriya raqam bilan almashtiriladi (0, 1, 2...).
- Qachon ishlatiladi: Kategoriyalar oʻzaro tartibli boʻlsa (masalan, low, medium, high).
- Misol:
- arduino
- Copy code
- Kategoriya: ['Low', 'Medium', 'High']
- Kodlangan: [0, 1, 2]

Kamchilik:

 Modellar kategoriyalarni tartibli deb qabul qiladi, lekin ular tartibsiz boʻlsa, notoʻgʻri natijalar boʻlishi mumkin.

2. One-Hot Encoding (Nominal -taribsiz)

- Qanday ishlaydi: Har bir kategoriya uchun alohida ustun yaratiladi, va shu ustunda 1 yoki 0 yoziladi.
- Qachon ishlatiladi:

Tartibsiz kategoriyalar uchun.

- Misol:
- CSS
- Copy code
- ['Apple', 'Banana', 'Cherry'] →
- Apple Banana Cherry
 - 1 0 0 0 1 0 0 0 1

Kamchilik:

 Juda koʻp kategoriya boʻlsa, ustunlar soni oshib ketadi (sparcity muammosi).

3. Target Encoding (maqsadli nishon)(ikkinchi nomi MEAN encoding)

OUTPUT ga bevosita bog'liq bo'lgan encoding.

- Qanday ishlaydi: Har bir kategoriya uchun maqsadli nishon (target) qiymatlarning oʻrtacha yoki boshqa statistikasi hisoblanadi.
- Qachon ishlatiladi:
 - Yugori oʻlchovli kategoriyali ma'lumotlar uchun.

Misol:

Kategoriya: ['A', 'B', 'A', 'B', 'C']
Target: [100, 200, 150, 250, 300]

1. Target Encoding

Formula:

 $\label{eq:Target Encoding} \text{Target Encoding}(k) = \frac{\text{Kategoriya k uchun target yig'indisi}}{\text{Kategoriya k uchun namunalar soni}}$

Hisob-kitob:

A:

Target Encoding(A) =
$$\frac{100 + 150}{2} = \frac{250}{2} = 125$$

B:

$$\operatorname{Target} \, \operatorname{Encoding}(B) = \frac{200 + 250}{2} = \frac{450}{2} = 225$$

C:

Target Encoding
$$(C) = \frac{300}{1} = 300$$

Natija:

yaml

Copy code

A: 125, B: 225, C: 300

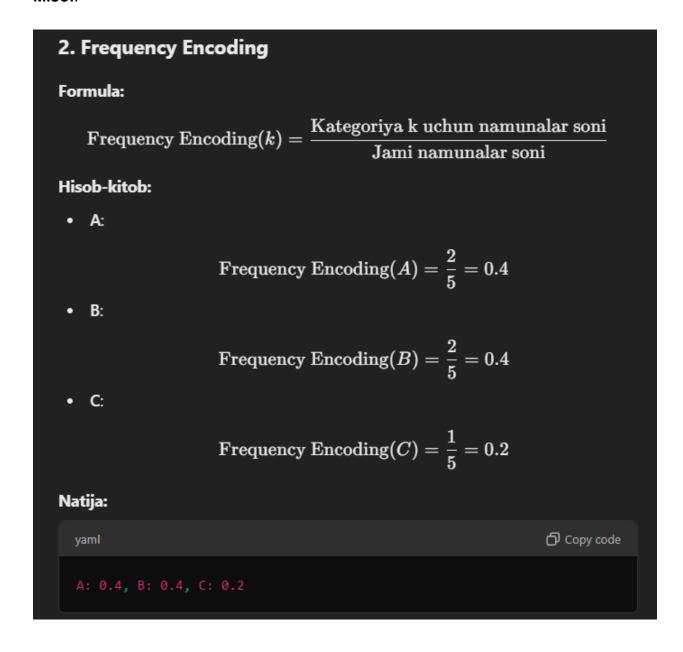
Kamchilik:

Overfitting xavfi bor, chunki target ma'lumotlardan foydalanyapti.

4. Frequency Encoding (Nominal -taribsiz)

Objectlar nechi marta ishlayotganini tekshirib o'shani o'rtachasi olinadi.

- Qanday ishlaydi: Har bir kategoriya ma'lumot to'plamida necha marta uchraganiga qarab kodlanadi.
- Qachon ishlatiladi:
 - Juda katta kategoriyali ma'lumotlar uchun samarali.
- Misol:



- Kamchilik:
- Ushbu kodlash kategoriyalar orasidagi oʻzaro bogʻliqlikni toʻliq ifodalamaydi.

Qanday encoderni tanlash kerak?

- 1. Tartibli kategoriya: Label Encoding yoki Target Encoding.
- 2. Tartibsiz kategoriya: One-Hot Encoding.
- 3. Katta kategoriyalar soni: Frequency yoki Target Encoding.

Xulosa

- Label va One-Hot Encoding asosiy klassik usullar.
- Target va Frequency Encoding murakkab va katta kategoriyalar uchun mos.

One Hot Encoder o'zi 2 methodga bo'linadi. (OneHotEncoder va get_dummies) OneHotEncoder vs pd.get_dummies

Xususiyat	OneHotEncoder (sklearn)	pd.get_dummies (pandas)
Kutubxona	scikit-learn	pandas
Interfeys	Machine Learning pipeline'lariga mos	DataFrame bilan ishlash uchun qulay
Natija turi	NumPy array yoki sparse matrix	Pandas DataFrame
NaN qiymatlar	Qoʻllab-quvvatlanmaydi	NaN qiymatlarni avtomatik boshqaradi
Model bilan ishlash	Sklearn modellariga mos	Toʻgʻridan-toʻgʻri modelda ishlatilmaydi
Bir nechta ustunlar	Qoʻllab-quvvatlanadi	Qoʻllab-quvvatlanadi
Kategoriyalarni boshqarish	Oʻqitish va sinov ma'lumotlari bilan bir xil kategoriya hosil qiladi	Kategoriyalarni avtomatik aniqlaydi
Qoʻllanish joyi ML pipeline va modellar bilan birga ishlatish uchun		Oddiy tahlillar va visualization uchun qulay

Qachon qaysi biridan foydalanish kerak?

1. OneHotEncoder:

- o Machine Learning pipeline'larida, scikit-learn modellarida ishlatiladi.
- NumPy array yoki sparse matrixni yaratadi, bu ML modellar bilan samarali ishlaydi.

2. pd.get_dummies:

- o Pandas DataFrame bilan ishlayotganingizda.
- Oddiy tahlil yoki vizualizatsiya uchun ma'lumotlarni One-Hot Encoding qilishda.

Ordinal va Nominal kategoriyalar uchun encodinglar jadvali

Kategoriy a turi	Tavsifi	Mos encodingla r	Nega ishlatiladi?	Misol
Ordinal Category	Tartibli: Kategoriyala r orasida tartib mavjud.	- Label Encoding- Ordinal Encoding	Tartibni saqlab qolish uchun.	['Low', 'Medium', 'High'] → [0, 1, 2]
				One-Hot: ['Red', 'Green', 'Blue']
		- One-Hot	Bogʻliqlikni	\rightarrow
	Tartibsiz:	Encoding-	oldini olish yoki	[[1,0,0],[0,1,0],[0,0,1]] Frequenc
Nominal	Kategoriyala	Frequency	koʻp	y : ['Red', 'Green', 'Blue', 'Red']
Category	r orasida	Encoding-	kategoriyalarni	\rightarrow {'Red': 2, 'Green': 1, 'Blue':
	tartib yo'q.	Target	optimallashtiris	1} Target : ['Red', 'Green', 'Blue']
		Encoding	h uchun.	→ {'Red': 100, 'Green': 200,
				'Blue': 150}

One Hot Encoder bilan Label Encoder tanlash

Qaysi Encoderni tanlashdan oldin Cardinality ni aniqlab olamiz. Cardinality bu bitta feature (xususiyat=ustun) ni ichidagi classlar soni.

Masalan rang degan feature bor va uning ichida 5 xil rang (class)lar bor. Cardinalitysi 5 ga teng.

Agar cardinality = > 5 dan bo'lsa Label Encoding Agar cardinality = < 5 bo'lsa One Hot Encodingdan foydalaniladi.