Homework 3 - Al

1. Danh sách nhóm 17

- a. Nguyễn Thái An 20020278
- b. Nguyễn Đức Anh 20020074
- c. Lê Tuấn Anh 20020286

2. Giới thiệu bài toán

- Trong bài này, bài toán nhóm sẽ thực hiện là bài toán dự đoán nhãn thu nhập của người
 lao động trên tập dữ liệu của cuộc thi UET Hackathon 2022 Data Science.
- Trong đó, các nhãn của dữ liệu được liệt kê trong bảng sau.

Nhãn	Ý nghĩa
7	Rất cao
6	Trung bình cao
5	Cao
4	Trung bình
3	Thấp
2	Trung bình thấp
1	Rất thấp

- Ở bài này, có 6 file dữ liệu chia thành 3 thể loại (info, work và label) cho hai tập (train và test). Nhưng trong bài này, nhóm sẽ chỉ sử dụng tập train để huấn luyện và kiểm thử.
- Ở tập info, có những trường sau đây:

Tên cột	Ý nghĩa
birthYear	Ngày sinh
gender	Giới tính
address	Địa chỉ
id_bh	Mã cá nhân

- Ở tập work, có những trường sau đây:

Tên cột	Ý nghĩa
id_bh	Mã cá nhân
id_management	Mã đơn vị quản lý
id_office	Mã văn phòng
company_type	Loại hình công ty
job/role	Vị trí, chức vụ của cá nhân
from_date	Công việc bắt đầu từ ngày
	nào
to_date	Công việc kết thúc vào ngày
	nào
employee_lv	Cấp độ của chức vụ

address	Địa chỉ
---------	---------

Ở tập label, có những trường như sau:

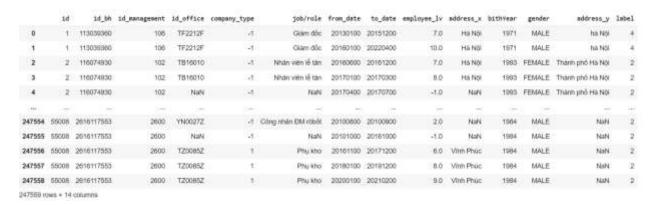
Tên cột	Ý nghĩa
id_bh	Mã cá nhân
label	Nhãn của thu nhập

Trong đó nhãn label là cột mà nhóm sẽ đi dự đoán.

3. Xử lý dữ liệu

a. Xử lý các file liên quan đến train đầu tiên

- Ta thấy ở cả 3 tập dữ liệu, đều có cùng một cột Mã khách hang, vì vậy ta cần ghép 3 tập dữ liệu lại với nhau để tiện cho quá trình xử lý.
- Sau khi ghép, ta sẽ có một tập dữ liệu lớn hơn với 247559 dòng và 14 cột



- Sau đó, gộp chung các thông tin address x và address y thành 1 cột chung là address

	ist	inf, belo	id_management	id office	company type	jels/rale	from date	to date	employes lv	address x	hithYear	gender	address_y	label	address
	. 1	112029360	106	TF2212F	-1	Giám átic	20130100	20151200	7.0	HAND	1971	MAGE	há hiệu		Hà Ngi
1	. 19	115039360	106	1102107	-4.	Glám đốc	20160100	20220400	10.0	Hallo	1971	MALE	149.7406	. 4	H9.1906
2	2	116074690	102	7816010	-1	Nhân viên liệ Sin	20160600	20161200	7.0	He No	1999	FOMALE	Thanh phố Ha Nói	2	Ha No
3	2	116074000	102	1816010	-1	Nhan viers iế tân	20170100	20170300	4.0	Harvor	1990	PEMALE	Thinh phố Ha Nói	2	H8.7600
4	2	116074930	102	NaNi	-1	Nati	20170400	20170700	-1.0	howe	1993	PEMALE	Thanh µnd Ha Noi	- 1	Thann phố Ha hội
-				-			-	-	-	-	-	-		-	-
247554	55008	2616117552	2600	V9400272	-4	Công nhân ĐM Hìbởi	20100600	20100900	2.0	New	1964	MALE	NiN	- 2	him
247555	55008	3616117550	2633	NaN	- 4	Nahi	20101000	20161000	1.0	heate	1984	MALE	run	- 1	Nati
247556	35008	2616117333	2600	7200852	1	Physides	20161100	30171200	6.0	Vinn Phos	1984	MALE	NeN	2	Vinn Phos
247557	55008	2618117552	2600	T20065Z		Phy kho	20190100	20191200	8.0	Ven Phúc	1954	MALE	Natio	2	Vinn Phúc
247558	55008	2916117553	2600	T200812	1	Physina	20000100	20210200	9.0	Vinn Phoc	1984	NACE	Nati	1	Vinh Phos

 Lấy dữ liệu năm chính xác từ hàng to_date, ta lấy 4 số đầu, bỏ 4 số sau bởi nó không có ý nghĩa. Từ 4 số đầu tiên, ta ra được năm

	id	id_bh	ld_management	ld_office	company_type	job/role	from_date	to_date	employee_lv	address_x	bithYear	gender	address_y	label	address	lastest
	1	113039360	106	TF2212F	-1	(iriam dőc	20130100	20151200	7.0	Hà No	1971	MALE	na Nói	4	Ha No	2015
1	- 1	113039560	106	TF2212F	-1	Gián đốc	20160100	20220400	10.0	HR NO	1971	MALE	hé No	4	Ha No	2022
2	2	116074930	102	TB16010	-1	Nhân viên lễ sán	2016/06/00	20161200	7.0	Ha No:	1993	FEMALE	Thanh phó Hà Nói	ž	Ha No	2016
3	- 2	116074930	102	TB16010	of	Ahán vián lễ tân	20170100	20170300	8.0	H8 NO	1993	FEMALE	Thanh phó Ha Noi	- 1	HA NO	2017
4	2	116074930	102	NaN	- 4	NaN	20170400	20170700	-1.0	NaN	1993	FEMALE	Thành phố Hà Nói	12	Thank phó Há Nói	2017
-	-	-		-		-	-	-		-	-	-	-			
247554	55008	2616117552	2600	VN0027Z	-1	Công nhân ĐM róbốt	20100800	20100900	2.0	NaNi	1984	MALE	NaN	2	NaN	2010
247555	55008	2616117553	3600	Nan	-11	Nati	20101000	20161000	-1.0	Note	1984	MALE	NaN	- 2	NaN	2016
247556	55008	2616117553	2600	7Z0085Z	31	Phy ino	20161100	20171200	6.0	Vinh Phác	1984	MALE	Nahy	ĘΣ	Vinh Phúc	2017
247557	55008	2616117553	1600	720085Z		Phy kho	20180100	20191200	8.0	Vinn Phoc	1984	MALE	NAME.	2	Vinh Phúc	2019
247558	59008	2616117553	3600	TZ0085Z		Phu kho	20200100	20210200	9.0	Vinh Phic	1984	MALE	hate	2	Vinh Phúc	2021

Ta tiến hành kiểm tra tỉ lệ bị thiếu từ dữ liệu.

id	0.000000
id bh	0.000000
id_management	0,000000
id_office	0.049152
company_type	0.000000
job/role	0.135148
from_date	0,000000
to_date	0.000000
employee_lv	0.000020
address_x	0.258649
bithYear	0,000000
gender	0.000000
address y	0.422740
label	0.000000
dtype: float64	

- Ta thấy được cột id_office bị thiếu khoảng 4%, cột job/role bị thiếu 13.5%, cột address_x bị thiếu 25%, cột address_y bị thiếu 42%. Các cột còn lại thì không bị thiếu.
- Trong 4 cột bị thiếu này, ta thấy cột job/role có ý nghĩa rất lớn đến việc xem xét thu nhập, vì điều kiện để biết được thu nhập của một người là phải xem chức vụ của họ là gì. Nên đối với những dòng bị thiếu job/role, thì ta tiến hành loại bỏ để có thể đạt được ý nghĩa cao hơn cho dữ liệu.

- Sau khi loại bỏ những hàng bị thiếu job/role, ta còn lại 214102 dòng và 14 cột.

	10	1d_bh	1d_management	id_office	company_type	job/role	from_date	to_date	employee_lv	address_x	bithyear	gender	address_y	label	
0	1	113039360	106	TF2212F	-11	Glám đốc	20130100	20151200	7.0	Hs No	1971	MALE	ha No	4	
1	1	113038960	106	TF2212F	-4	Glám đốc	20180100	20220400	10,0	Ha No	1971	MALE	ha No	4	
2	2	116074030	102	TB16010	- 1	Nhân viên lễ tân	20160600	20161200	7.0	Ha No	1993	FEMALE	Thành phố Hà Nội	2	
3	2	110074930	102	TB19010	3 30	Nhân viên lễ tân	20170100	20170300	8,0	Ha No	1983	FEMALE	Thanh phố Hà Nội	2	
5	2	116074930	102	TB16010	- 4	Nhān viện Sales Admin	20170800	20191200	8.0	Ha No	1993	FEMALE	Thanh phố Hà Nội	2	
100	-	900							-		-	-		- 17	
247553	55008	2616117555	2000	YN0027Z	-4	Han X3	20100700	20100700	2.0	MaN	1984	MALE	14014	2	
247554	55008	2616117563	2600	YN0027Z	- 4	Công rihân ĐM rộbắt	20100800	20100900	2.0	NuN	1984	MALE	NiN	2	
247556	55008	2616117553	2000	TZ0085Z	1 1	Phy kho	20161100	20171200	6.0	Virin Phoc	1984	MALE	NoN	2	
247557	65008	2616117553	2900	TZ0085Z		Phy kho	20180100	20191200	8.0	Vinh Phúc	1984	MALE	NoN	2	
247558	55000	2616117553	2600	T20065Z	1	Phy kho	20200100	20210200	9.0	Winh Phúc	1964	MALE	NaN	2	
21410210	ows = 14	columns													

Sau đó, ta tiếp tục lọc các giá trị bị trùng nhau. Ta lọc dựa trên 3 tiêu chí là 'id_bh',
 'to_date', 'from_date'. Ta dùng code như sau:

```
df_filter = df.sort_values(
    by=["id_bh", 'to_date", 'from_date"],
    ascending=[True, False, False]
).drop_duplicates(
    subset=['id_bh'],
    keep='first',
    ignore_index=True
)
df_filter
```

- ⇒ Ta chỉ còn 27490 dòng và 14 cột
- Sau đó, ta tiến hành tính số tuổi của các nhân viên trong danh sách, tạo them 1 cột là 'age' cho dễ nhìn. Từ đó, ta bắt đầu lọc những người bị delay thời gian dựa trên thời gian tham gia, sau đó có dữ liệu để dễ dàng làm việc hơn

	h	id,bh	id_management	id_office	company_type	job/role	from_date	to_date	employee_lv	address,x	bithYear	gender	address_y	label	address	limitest	age	delay
0	468	100000725	100	HW01180	6	CVC.TP	20210600	20220400	17.0	HN	1963	MALE	Ha Noi	5	HN	2023	50	0
a	2900	100007067	2400	TA0002A		Nhân viên bản xông dầu	20190100	20220400	9.0	TP Bắc Giang	1971	MALE	8ác Gwng	- 4	TP Bắc Giang	2022	51	0
2	450	100007555	109	7(3361)	140	Nhân siên	20211000	20220400	9.0	Viet Nam	1970	MALE	NaN	2	Việt Nam	2022	52	0
3	293	9 100008102	100	HW0013Z	+1	Chuyên viên chinh, Úy wên thường trực Ủy ban	20200000	20225400	18.0	Ha Noi	1970	MALE	нN	5	Ha Noi	2022	52	0
4	896	7 100098777	714	HN05360	÷1	PCT UBND Noyen	28201000	20220400	17.0	HA NÓI	1965	MALE	NaN	- 6	HÀNŌI	2023	57	0
1	1 6		12	1 2		- 2	100	1 2	-	-	-	50			- 1	- 2		
17485	3205	9719632107	129	7502775	. 4	Ký sự công nghệ thóng tin	20210700	20211200	14.0	NaN	1995	MALE	NaN	- 6	NoN	2021	27	t

 Tiếp theo, t dùng hàm fillna để điền vào các vị trí bị trống của office và address (trong bảng cũ hiện là NaN)

```
df_filter['id_office'].fillna('Unknow', inplace=True)
df_filter['address_x'].fillna('Việt Nam', inplace=True)
df_filter['address_y'].fillna('Việt Nam', inplace=True)
df_filter
```

⇒ Đây là kết quả

	id	id_bh	id_management	id_affice	company_type	job/role	from_date	to_date	employee_lv	address_x	bithYear	gender	address_y	label	address	lastest	age	delay
0	4660	100000725	100	HW01180	1.9	CVC, TF	20210600	20220400	17,0	999	1963	MALE	Ha Nor	5	HN	3022	58	D
1	29905	100007067	2400	TADODZA	1	Mhán viên bán xông tiệu	20190100	20220400	9.0	TP Bác Glang	1971	MALE	8.6c Glang	4	TP Bác Giarig	2022	51	0
2	4503	100007555	109	1133611	9	Nhán viên	20211000	20220400	9.0	Viet Nam	1970	MALE	Việt Nam	2	Viet Nam	2022	52	0
3	2939	100008102	100	HW0013Z	э	Chuyển viên chính, Ủy viên Thường trực Ủy ban	20200800	20225400	18,0	Ha Noi	1970	MALE	HN	3	Ha Nos	.2022	52	0
ä	8967	100008777	114	HN05360	31	PCT UBNO huyên	20201000	20225400	17.0	HÀ NÓI	1965	MALE	Viot Nam	6	HÁ NÓI	2022	57	0
100	-	-			-		1	- 1		- 4				100			-	
27485	32655	9719632107	129	1502775	-1	Ký nơ công nghệ thông tin	28210700	20211200	14.0	Việt Nam	1995	MALE	Việt Nam	6	Net	2021	27	or
27486	32685	9809227380	4009	7111934	э	Quitn'ly ben hong	20210900	20220400	7.0	Tinn Nghệ An	1984	MALE	Viet Num	2	Tinh Nghè An	2022	30	0

- Sử dụng hàm thuộc Catboost, ta chia các dữ liệu đã ghi nhận được ra thành 4 file
 - b. Xử lý dữ liệu liên quan đến test (làm tương tự 3a)
 - Sau khi hoàn tất việc xử lý dữ liệu, ta gộp kết quả với các file đã được chia, sử dụng hàm để lưu các dữ liệu thành 1 file csv

4. Xây dựng mô hình

- Bài toán sử dụng mô hình Catboost, mô hình này triển khai cây đối xứng (symmetric trees) giúp giảm thời gian dự đoán và nó cũng có độ sâu của cây nông hơn theo mặc định (sáu)
- CatBoost tận dụng các hoán vị ngẫu nhiên
- 5. Kết quả

