

作业说明：

分类任务， 根据一个人的年龄、职业、学历等信息预测此人的收入情况（大于 50000 美元或小于 50000 美元）， 因此为 2 分类问题。

数据集说明：

数据存放于 data 文件夹中， 共 6 个文件：

train.csv： 原格式的训练数据， 54256 行数据， 代表 54256 个角色

第一列为 id, 最后一列为预测结果, 即收入大于 50000 美元或小于 50000 美元；

其余各列为用于分类的属性值， 例如： 年龄、职业等；

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
1	id	age	class of work	detailed occupation	education	wage per hour	enroll in education	marital status	major industry	major occupation	race	hispanic or latino	sex	member of union	reason for full or part-time	capital gain	capital loss	dividends	tax filer status	re	
2	0	33	Private	34	26	Masters c	0	Not in un	Married-i	Finance ii	Adm supp	Asian or f	All other	Female	Not in un	Not in un	Full-time	0	0	0	Joint bott
3	1	63	Private	7	22	Some col	0	Not in un	Never m	Manufact	Adm supp	White	All other	Female	Not in un	Not in un	Full-time	0	0	0	Single
4	2	71	Not in un	0	0	7th and E	0	Not in un	Married-i	Not in un	Not in un	White	All other	Male	Not in un	Not in un	Not in la	0	0	0	Joint bott
5	3	43	Local gov	43	10	Bachelors	0	Not in un	Married-i	Educator	Professio	White	All other	Female	Yes	Not in un	Full-time	0	0	0	Joint bott
6	4	57	Local gov	40	32	Some col	0	Not in un	Widowed	Entertain	Other ser	Amer Ind	All other	Female	Not in un	Not in un	Full-time	0	0	0	Head of f
7	5	42	Private	16	4	Masters c	0	Not in un	Married-i	Manufact	Professio	Asian or f	All other	Male	Not in un	Not in un	Children	0	1902	165	Joint bott
8	6	16	Not in un	0	0	11th grac	0	High sch	Never m	Not in un	Not in un	White	Central o	Male	Not in un	Not in un	Children	0	0	0	Single
9	7	16	Private	33	19	10th grac	0	High sch	Never m	Retail tra	Sales	Other	Other Sp	Female	Not in un	Not in un	Full-time	0	0	0	Nonfiler
10	8	20	Private	5	36	High sch	0	Not in un	Never m	Manufact	Machine	White	All other	Male	Not in un	Job loser	Unemplo	0	0	0	Single
11	9	38	Not in un	0	0	Bachelors	0	Not in un	Married-i	Not in un	Not in un	White	All other	Female	Not in un	Not in un	Children	0	0	125	Joint bott
12	10	34	Not in un	0	0	Some col	0	Not in un	Separate	Not in un	Not in un	White	All other	Female	Not in un	Not in un	Not in la	0	0	0	Nonfiler
13	11	39	Self-emp	37	5	Bachelors	0	Not in un	Married-i	Business	Professio	White	All other	Male	Not in un	Not in un	Full-time	0	0	5000	Joint bott
14	12	62	Not in un	0	0	7th and E	0	Not in un	Married-i	Not in un	Not in un	White	Mexican-	Male	Not in un	Not in un	Not in la	0	0	0	Nonfiler
15	13	25	Not in un	0	0	High sch	0	Not in un	Married-i	Not in un	Not in un	White	All other	Male	Not in un	Not in un	Children	0	0	0	Joint bott
16	14	30	Private	24	26	Bachelors	0	Not in un	Married-i	Manufact	Adm supp	White	All other	Female	No	Not in un	Children	0	0	150	Joint bott
17	15	35	State gov	45	11	Doctorate	0	Not in un	Married-i	Other prc	Professio	White	All other	Female	Not in un	Not in un	Children	0	0	0	Joint bott
18	16	40	Private	30	38	High sch	0	Not in un	Separate	Commun	Transport	Black	All other	Female	Not in un	Not in un	Children	0	0	175	Head of f
19	17	34	Not in un	0	0	Bachelors	0	Not in un	Married-i	Not in un	Not in un	White	All other	Female	Not in un	Not in un	Children	0	0	0	Joint bott
20	18	80	Not in un	0	0	Some col	0	Not in un	Married-i	Not in un	Not in un	White	All other	Female	Not in un	Not in un	Children	0	0	1000	Joint bott
21	19	31	Private	3	35	High sch	0	Not in un	Never m	Mining	Precision	White	All other	Male	Not in un	Not in un	Full-time	0	0	0	Single
22	20	28	Local gov	47	28	Some col	0	Not in un	Married-i	Public ad	Protective	White	All other	Male	Not in un	Not in un	Full-time	0	0	0	Joint bott
23	21	39	Private	24	2	Masters c	0	Not in un	Married-i	Manufact	Executive	White	All other	Female	Not in un	Not in un	Children	0	2415	1000	Joint bott
24	22	2	Not in un	0	0	Children	0	Not in un	Never m	Not in un	Not in un	White	Mexican	Male	Not in un	Not in un	Children	0	0	0	Nonfiler
25	23	67	Not in un	0	0	Masters c	0	Not in un	Married-i	Not in un	Not in un	White	All other	Female	Not in un	Not in un	Not in la	0	0	304	Joint bott
26	24	15	Private	33	29	7th and E	0	Not in un	Never m	Retail tra	Other ser	White	All other	Female	Not in un	Re-entra	Children	0	0	0	Single
27	25	60	Private	45	19	Masters c	0	Not in un	Married-i	Other prc	Professio	White	All other	Male	Not in un	Not in un	Full-time	0	0	1000	Joint bott

test\_no\_label.csv： 原格式的测试数据， 与 tain.csv 类似， 没有给出预测结果， 需

要模型输出；

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V
1	id	age	class of work	education	education	education	education	enroll in	enroll in	enroll in	enroll in	enroll in	enroll in	enroll in	enroll in	enroll in	enroll in	enroll in	enroll in	enroll in	enroll in	enroll in
2	0	37	Private	42	30	Associate	0	Not in un	Married	Medical	Other ser	White	Hispanic o	sex	member c	reason for full or part	capital ga	capital los	dividends	tax filer st	region of	
3	1	48	Private	31	35	High sch	0	Not in un	Married	Utilities a	Precision	White	All other	Female	Not in un	Not in un	Full-time	0	0	0	Joint bot	Not in un
4	2	68	Not in un	0	0	High sch	0	Not in un	Married	Not in un	Not in un	White	All other	Female	Not in un	Not in un	Children	0	0	0	Single	Not in un
5	3	74	Private	12	36	11th grad	700	Not in un	Married	Manufact	Machine	White	All other	Male	No	Not in un	Full-time	0	0	0	Joint bot	Not in un
6	4	18	Federal g	50	14	Some col	0	Not in un	Never ms	Public ad	Technicia	White	All other	Male	Not in un	Re-entrai	Unemplo	0	0	0	Nonfiler	Not in un
7	5	46	Private	27	4	High sch	0	Not in un	Married	Manufact	Professio	White	All other	Male	Not in un	Not in un	Full-time	0	0	100	Joint bot	Not in un
8	6	17	Not in un	0	0	9th grade	0	High sch	Never ms	Not in un	Not in un	White	All other	Male	Not in un	Not in un	Children	0	0	0	Nonfiler	Midwest
9	7	38	Self-emp	45	17	Bachelors	0	Not in un	Widowed	Not in un	Not in un	White	All other	Female	Not in un	Not in un	Children	15024	0	0	Joint bot	Not in un
10	8	90	Not in un	0	0	Bachelors	0	Not in un	Never ms	Construct	Precision	White	All other	Male	Not in un	Not in un	Children	0	0	0	Nonfiler	Northeas
11	9	36	Private	4	34	High sch	0	Not in un	Never ms	Construct	Precision	White	All other	Male	Not in un	Not in un	Children	0	0	0	Single	Not in un
12	10	6	Not in un	0	0	Children	0	Not in un	Never ms	Not in un	Not in un	Other	Puerto Ri	Female	Not in un	Not in un	Children	0	0	0	Nonfiler	Not in un
13	11	16	Private	4	34	High sch	0	Not in un	Never ms	Construct	Precision	White	Central o	Male	Not in un	Not in un	PT for eco	0	0	0	Single	Not in un
14	12	37	Self-emp	35	17	Bachelors	0	Not in un	Married	Finance ii	Sales	White	All other	Male	Not in un	Not in un	Children	0	2415	0	Joint bot	Not in un
15	13	51	Private	21	36	7th and 8	0	Not in un	Separate	Manufact	Machine	Black	All other	Female	Not in un	Not in un	Children	0	0	0	Single	Not in un
16	14	36	Private	42	30	11th grad	650	Not in un	Never ms	Medical	Other ser	White	All other	Female	No	Not in un	Full-time	0	0	0	Head of f	Not in un
17	15	54	Private	33	2	High sch	0	Not in un	Married	Retail tra	Executive	White	All other	Male	No	Not in un	Full-time	0	0	0	Joint bot	Not in un
18	16	64	Private	34	24	High sch	0	Not in un	Married	Finance ii	Adm supp	Black	All other	Female	Not in un	Not in un	Full-time	0	0	0	Joint one	Not in un
19	17	36	Not in un	0	0	Some col	0	Not in un	Married	Not in un	Not in un	Asian or I	All other	Male	Not in un	Not in un	Children	0	0	0	Single	Not in un
20	18	47	Not in un	0	0	1st, 2nd 3	0	Not in un	Married	Not in un	Not in un	White	Central o	Female	Not in un	Not in un	Not in la	0	0	0	Joint bot	Not in un
21	19	39	Private	37	12	Some col	0	Not in un	Married	Business	Professio	White	All other	Male	No	Not in un	Full-time	0	0	0	Joint bot	Not in un
22	20	45	Private	30	2	Some col	0	Not in un	Married	Commun	Executive	White	All other	Male	Not in un	Not in un	Full-time	0	0	275	Joint bot	Not in un
23	21	35	Private	43	10	12th grad	0	Not in un	Divorced	Educator	Professio	White	All other	Female	Not in un	Not in un	Full-time	0	0	0	Head of f	Not in un
24	22	47	Private	30	33	High sch	1950	Not in un	Married	Commun	Precision	White	All other	Male	Yes	Not in un	Children	0	0	644	Joint bot	Not in un
25	23	54	Private	32	16	Some col	0	Not in un	Married	Wholesale	Sales	White	All other	Female	No	Not in un	Full-time	0	0	5	Joint bot	Not in un
26	24	46	Private	19	35	Bachelors	0	Not in un	Married	Manufact	Precision	White	All other	Male	Not in un	Not in un	Children	0	0	0	Joint bot	Not in un

sample\_submission.csv: 输出样例

两列数据，第一列为 id，第二列为预测结果

	A	B	C
1	id	label	
2	1	0	
3	2	0	
4	3	0	
5	4	0	
6	5	0	
7	6	0	
8	7	0	
9	8	0	
10	9	0	
11	10	0	
12	11	0	
13	12	0	
14	13	0	
15	14	0	
16	15	0	
17	16	0	
18	17	0	
19	18	0	
20	19	0	
21	20	0	
22	21	0	
23	22	0	
24	23	0	
25	24	0	
26	25	0	
27	26	0	



## 输入：

直接使用处理过的数据即可

训练数据：X\_train、Y\_train

测试数据：X\_test

## 输出：

输出一个.csv 文件，参照 sample\_submission.csv 的格式，第一列为测试数据中所有人的 id，第二列为收入预测结果。

## 示例程序：

1. Logistic\_regression，逻辑回归方法，运行后会创建 output\_logistic.csv，内容为对测试数据的预测；具体细节见代码及注释。

其余输出内容：

训练测试数据参数：

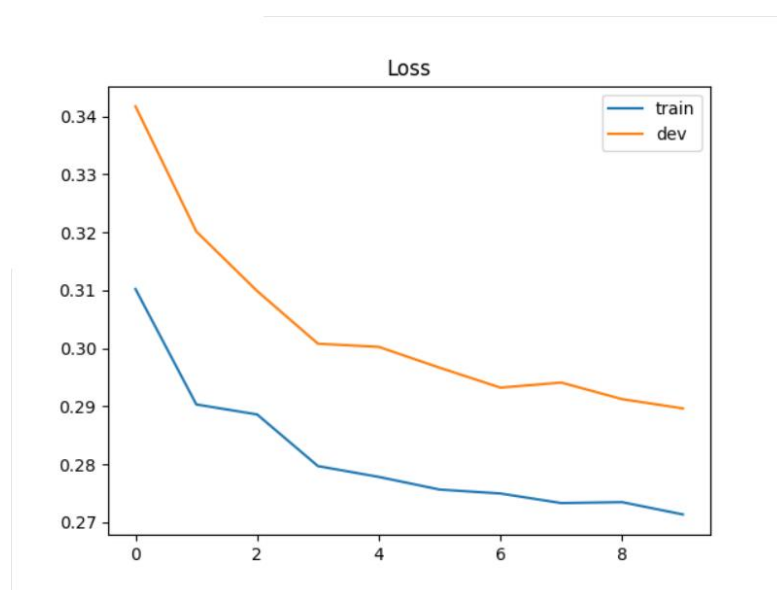
```
Size of training set: 48830
Size of development set: 5426
Size of testing set: 27622
Dimension of data: 510
```

训练集、验证集效果：

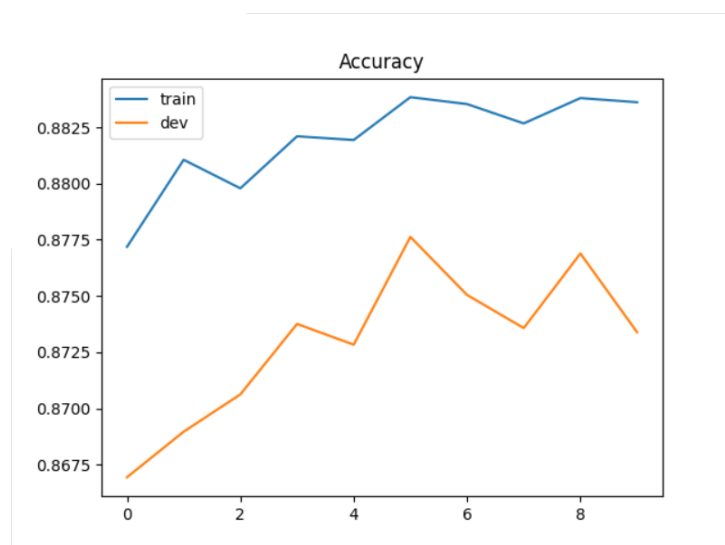
```
Training loss: 0.271355435246406
Development loss: 0.2896359675026287
Training accuracy: 0.8836166291214418
Development accuracy: 0.8733873940287504
```

训练效果可视化:

损失函数变化图:



准确率变化图:



数据前 10 项特征对应的权重:

```
Other Rel <18 never married RP of subfamily -1.4195759775765406
Child 18+ ever marr Not in a subfamily -1.2958572076664743
Unemployed full-time 1.171255828588591
Other Rel <18 ever marr RP of subfamily -1.167791807296237
Italy -1.093458143800618
Vietnam -1.0630365633146408
num persons worked for employer 0.9389922773566495
1 0.8226614922117187
```

2. Generative\_model 方法和 Logistic regression 方法类似, 不同之处在于 Generative model 可以直接计算出  $w$  和  $b$  的最佳解, 而 Logistic regression 是将  $w$  和  $b$  进行初始化, 通过迭代训练来更新  $w$  和  $b$ , 代码除了求解  $w$  和  $b$  地方不一样, 其他地方类似。