

به نام خدا  
دانشگاه صنعتی امیرکبیر  
(پلی تکنیک تهران)  
دانشکده مهندسی کامپیوتر



دانشگاه صنعتی امیرکبیر  
(پلی تکنیک تهران)

تمرین دوم درس یادگیری ماشین

گزارش سوال‌های پیاده‌سازی

استاد درس: دکتر احسان ناظر فرد

دانشجو: فاطمه غلامزاده

۹۹۱۳۱۰۰۳

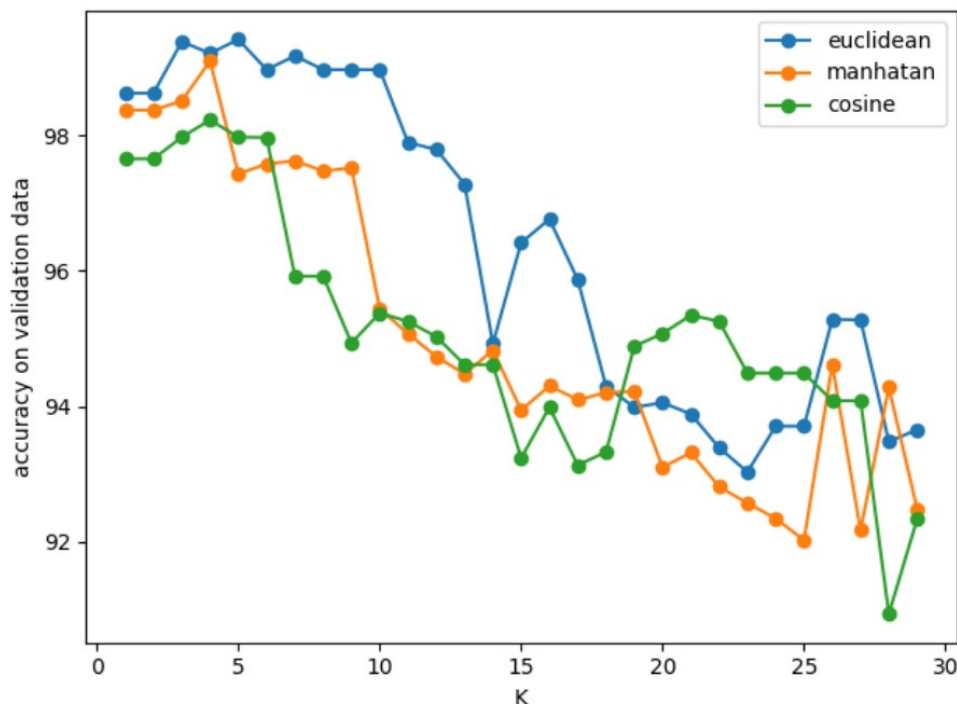
نیم سال دوم ۱۳۹۹-۱۴۰۰

## سوال‌های پیاده‌سازی

### سوال اول:

**الف)** دیتا به سه دسته آموزش، ارزیابی و تست که هر کدام به ترتیب سهم ۶۰ و ۲۰ و ۲۰ درصدی دارند تقسیم شده است.

برای یافتن بهترین مدل، آموزش و اندازه گیری دقت مدل بر روی داده های ارزیابی انجام گرفت بدین صورت که برای مقادیر  $k$  از ۱ تا ۳۰ و برای سه فاصله اقلیدسی، منهتن و کسینوسی، دقت اندازه گیری شد. نمودار دقت برای  $k$ ها و توابع فاصله مختلف رسم شده است :



بهترین مدل، مدلی است که از فاصله اقلیدسی و با  $k=5$  استفاده می‌کند که دقت آن برابر ۹۹,۴۱٪ شد. بنابراین آموزش مدل را با این پارامترها انجام داده و دقت را برای داده‌های تست اندازه گیری می‌کنیم. میزان دقت بر روی داده‌های آموزش، ارزیابی و تست در جدول زیر آورده شده است:

Model parameters	Train accuracy	Validation accuracy	Test accuracy
K=4 & Euclidean	۹۸,۴۲	۹۹,۴۱	۹۸,۳۳

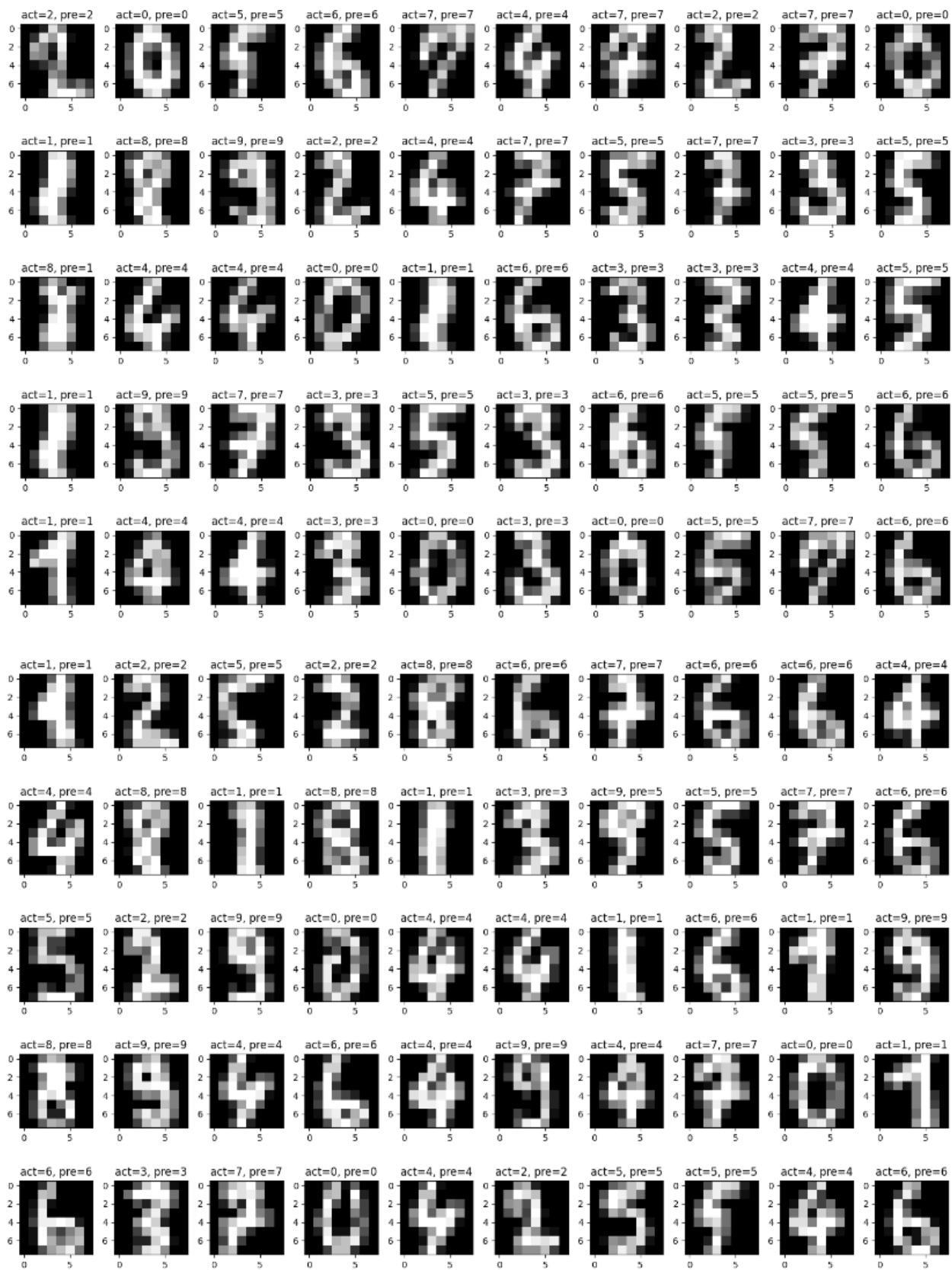
میزان خطای آموزش، ارزیابی و تست:

Model parameters	Train Error	Validation Error	Test Error
K=4 & Euclidean	1.58%	0.59%	1.67%

ماتریس درهم ریختگی:

		Predicted Class									
True Class		0	1	2	3	4	5	6	7	8	9
	0	33	0	0	0	0	0	0	0	0	0
	1	0	28	0	0	0	0	0	0	0	0
	2	0	0	33	0	0	0	0	0	0	0
	3	0	0	0	34	0	0	0	0	0	0
	4	0	0	0	0	46	0	0	0	0	0
	5	0	0	0	0	0	45	1	0	0	1
	6	0	0	0	0	0	0	35	0	0	0
	7	0	0	0	0	0	0	0	33	0	1
	8	0	1	0	0	0	0	0	0	29	0
	9	0	0	0	0	1	1	0	0	0	38

مقدار واقعی و مقدار پیش بینی شده برای ۱۰۰ رقم که به طور تصادفی از داده های تست انتخاب شده اند رسم شده است:



(ب)

در این قسمت از کتابخانه‌ی sklearn و KNeighborsClassifier استفاده شده است. بهترین مدل برای  $k=3$  به دست آمد که میزان دقت آن روی داده های ارزیابی ۹۹,۱۷٪ شد.

میزان دقت :

Model parameters	Train accuracy	Validation accuracy	Test accuracy
K=3 & Euclidean	۹۸,۹۷	۹۹,۱۷	۹۸,۶۱

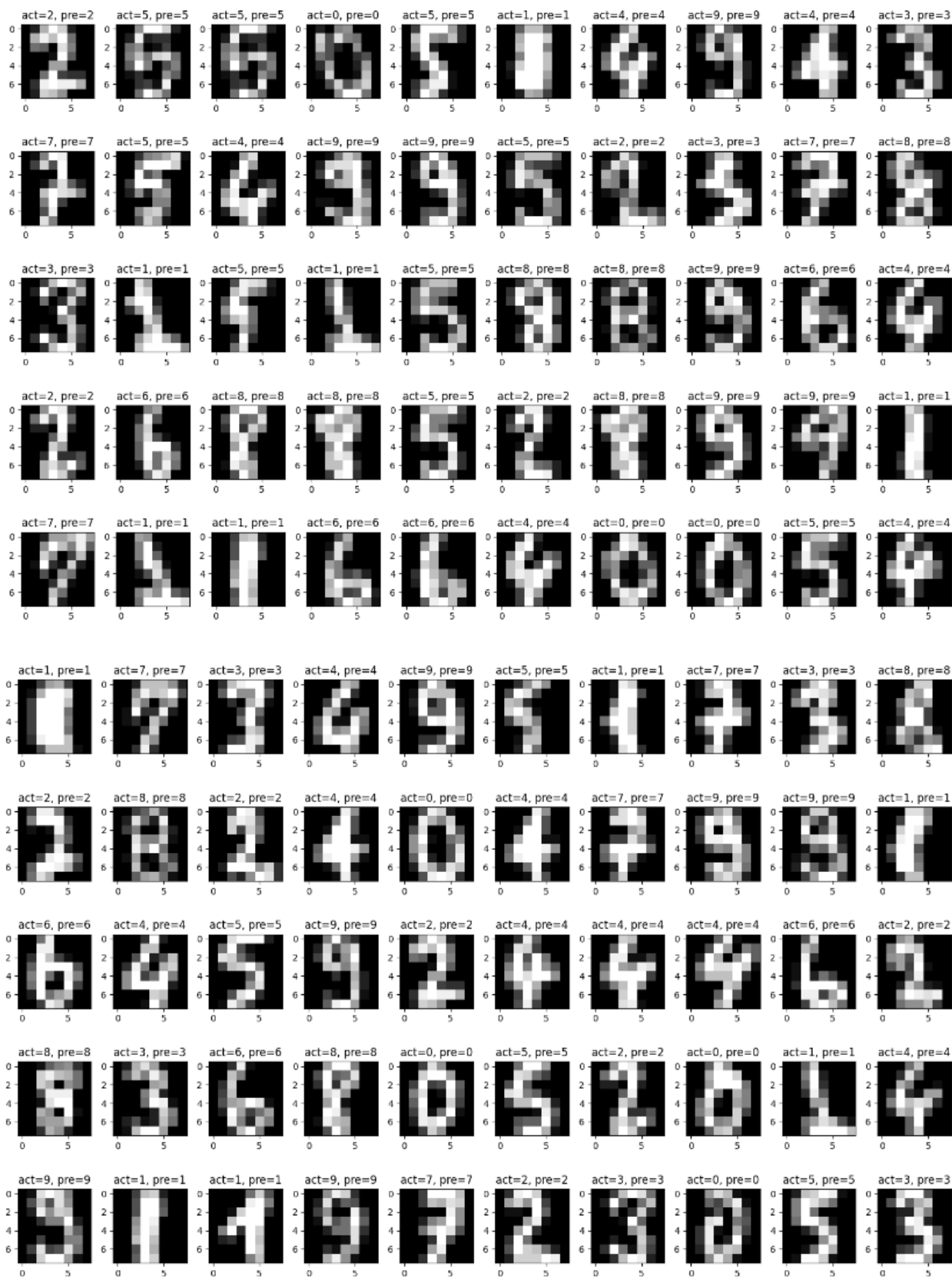
میزان خطای آموزش، ارزیابی و تست:

Model parameters	Train Error	Validation Error	Test Error
K=3 & Euclidean	1.03%	0.87%	1.39%

ماتریس در هم ریختگی:

		Predicted Class									
True Class		0	1	2	3	4	5	6	7	8	9
	0	33	0	0	0	0	0	0	0	0	0
	1	0	28	0	0	0	0	0	0	0	0
	2	0	0	33	0	0	0	0	0	0	0
	3	0	0	0	34	0	0	0	0	0	0
	4	0	0	0	0	46	0	0	0	0	0
	5	0	0	0	0	0	46	0	0	0	1
	6	0	0	0	0	0	0	35	0	0	0
	7	0	0	0	0	0	0	0	33	0	1
	8	0	1	0	0	0	0	0	0	29	0
	9	0	0	0	0	1	1	0	0	0	38

مقدار واقعی و مقدار پیش بینی شده برای ۱۰۰ رقم که به طور تصادفی از داده های تست انتخاب شده اند رسم شده است:



## مقایسه نتایج قسمت الف و ب:

همانطور که مشاهده می‌شود با استفاده از کتابخانه آماده خطای آموزش و تست مقدار کمی کاهش پیدا کرده است. خطای آموزش از 1.58% به 1.03% و خطای تست از 1.67% به 1.39% کاهش یافته است. با مقایسه ماتریس‌های درهم‌ریختگی هم می‌توان مشاهده کرد که با استفاده از کتابخانه آماده، تعداد داده‌های درست دسته‌بندی شده فقط ۱ مورد افزایش یافته است.

## سوال دوم:

### الف)

**پیش پردازش:** ابتدا داده‌ها را شافل می‌کنیم چون کلاس‌ها به ترتیب در داده‌ها ظاهر شده‌اند و این موجب کاهش دقت دسته‌بندی می‌شود. هم‌چنین برخی از ستون‌های داده‌ها را استاندارد می‌کنیم.

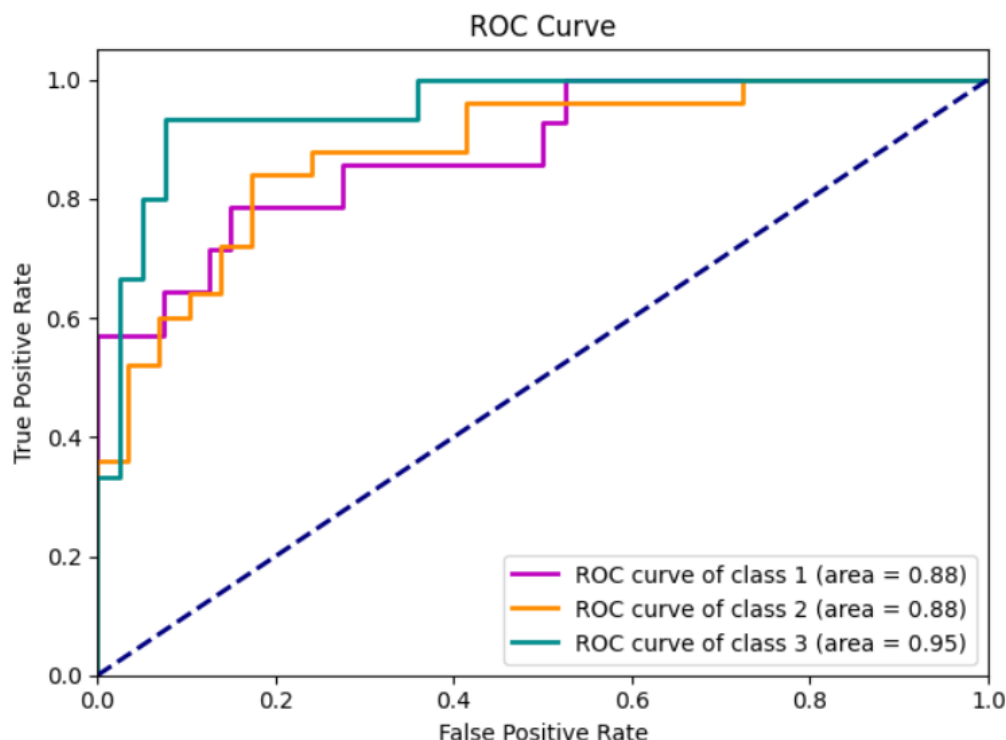
دقت به دست آمده با روش 6-fold cross validation : ۹۳,۴۵٪

### ب)

در این قسمت ۷۰ درصد از داده‌ها را برای آموزش و ۳۰٪ را برای تست در نظر گرفتیم.

دقت به دست آمده : ۹۴,۳۳٪

نمودار ROC :



### تحلیل نمودار ROC:

هر چقدر نمودار ROC از خط  $y=x$  فاصله بیشتری داشته باشد بهتر است، یعنی هر چقدر مساحت زیر منحنی ROC بیشتر باشد دقت بالاتر است، با توجه به نمودار بالا می‌توان مشاهده کرد که دسته بندی برای کلاس ۳ با دقت بیشتری انجام شده چون مساحت زیر منحنی ROC برای آن بزرگتر است.

### سوال سوم:

#### (الف)

در این قسمت ۷۰٪ از داده‌ها برای آموزش و ۳۰٪ برای تست در نظر گرفته شده است. آموزش به این صورت است که ۱۰ تا دسته بند (به ازای هر رقم یک دسته بند) آموزش می‌دهیم. برای آموزش هر دسته بند برچسب داده‌ی ترین مخصوص خود آن دسته بند تولید می‌کنیم بدین صورت که اگر برچسب داده ترین برابر آن رقم بود آن را ۱ می‌کنیم و اگر غیر از آن رقم بود آن را ۰ می‌کنیم. بدین صورت برای هر رقم یک دسته بند باینری آموزش می‌دهیم. وقتی یک داده تست می‌آید هر ۱۰ تا دسته بند را بر روی آن اعمال می‌کنیم و احتمال تعلق آن برای هر دسته که بیشتر بود نتیجه می‌گیریم رقم مربوط به آن دسته است.



دقت مدل بر روی داده‌های تست: ۹۵,۵۵

خطای آموزش و تست:

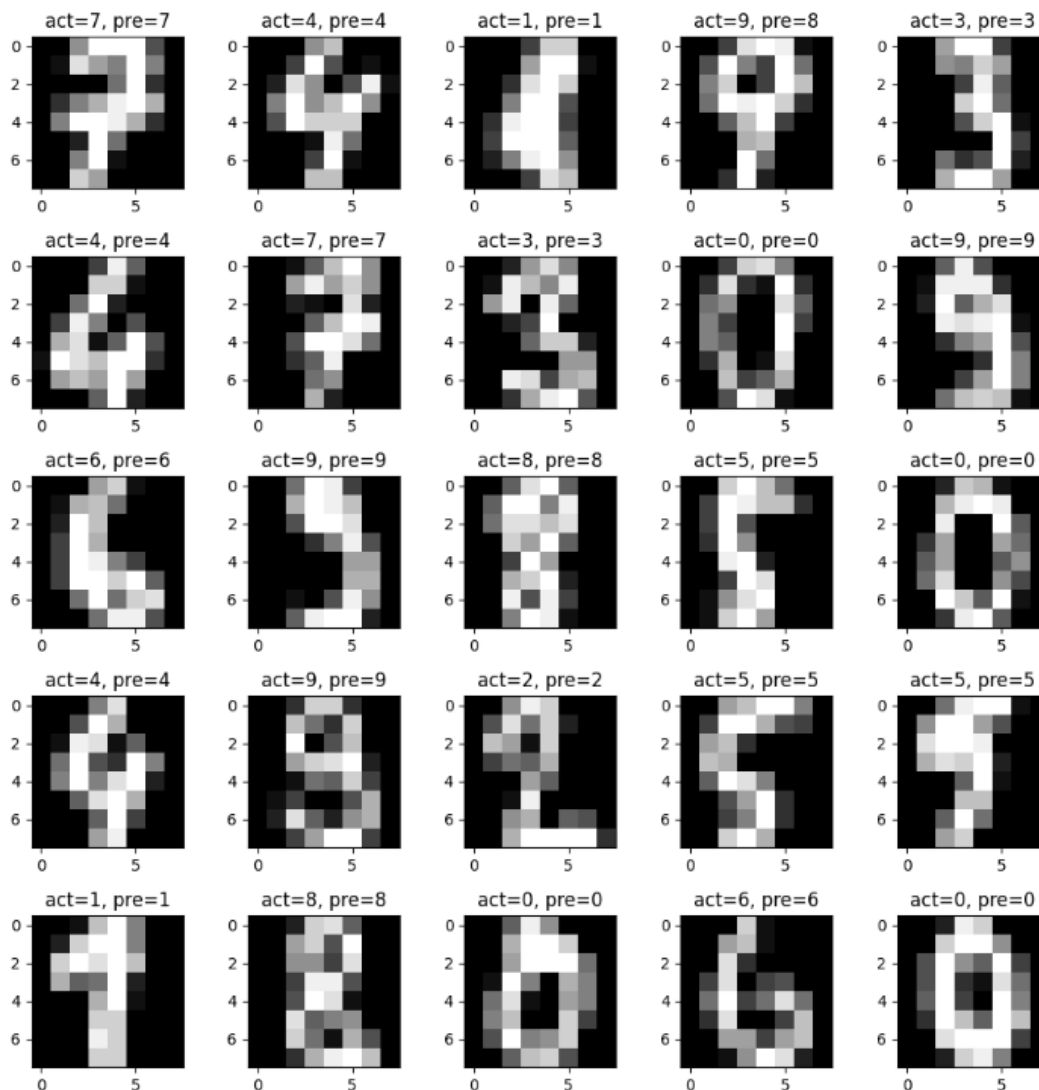
Model	Train Error	Test Error
Logistic Regression	0.4 %	4.44%

ماتریس در هم ریختگی:

		Predicted Class									
True Class		0	1	2	3	4	5	6	7	8	9
	0	53	0	0	0	0	0	0	0	0	0
	1	0	42	1	0	0	0	0	0	7	0
	2	0	0	47	0	0	0	0	0	0	0
	3	0	0	0	52	0	1	0	0	1	0
	4	0	1	0	0	59	0	0	0	0	0
	5	0	1	1	0	0	62	0	1	1	0
	6	0	0	0	0	0	1	52	0	0	0
	7	0	0	0	0	0	0	0	54	0	1
	8	0	1	0	0	0	1	0	0	41	0
	9	0	0	0	0	0	0	0	0	5	54

(ب)

مقدار واقعی و مقدار پیش بینی شده برای ۲۵ رقم که به طور تصادفی از داده های تست انتخاب شده اند رسم شده است:



(ج)

### مقایسه عملکرد KNN با Logistic Regression :

مرحله training در KNN بسیار سریع‌تر از logistic Regression است اما مرحله test در KNN بسیار کندتر است زیرا KNN یک lazy learner است و بیشتر محاسبات مربوط به دسته‌بندی را در زمانی که داده تست به آن داده می‌شود انجام می‌گیرد.

دقت به دست آمده در روش KNN از روش Logistic Regression بیشتر است.

(د)

### توضیح یادگیری نامتوازن:

به این معناست که فرکانس کلاس‌های هدف در داده‌ها به طرز چشمگیری نامتوازن است. به عنوان مثال وقوع یکی از کلاس‌ها در داده‌ها در مقایسه با سایر کلاس‌ها بسیار بیشتر است و الگوریتم دسته‌بندی تمایل پیدا می‌کند که داده‌های تست را در کلاسی که تکرار بالایی دارد جای دهد. فرض کنید یک مسئله دسته‌بندی دو کلاس با کلاس‌های A و B داریم. فرض کنید داده‌های آموزشی به گونه‌ای است که ۹۵٪ آن متعلق به کلاس A و فقط ۵٪ آن متعلق به کلاس B است. وقتی مدل را روی این داده‌های آموزشی، آموزش بدهیم اگر حتی همه داده‌ها را در کلاس A دسته‌بندی کند دقت مدل ۹۵٪ است و دقت قابل قبولی است اما داده‌های کلاس B نادیده گرفته شده‌اند! به بیان دیگر تعداد سмпل‌های کلاس B به اندازه‌ای کافی نیستند که مدل بتواند بفهمد چه زمانی یک داده را در کلاس B جای دهد. حالا فرض کنید در یک مسئله دسته‌بندی ۵ کلاس دارید که در متوازن‌ترین حالت ممکن هر کلاس دارای ۲۰٪ از داده‌ها است و اگر از روش one-vs-all استفاده کنیم در هر بار استفاده از مدل، یک توزیع ۲۰-۸۰ درصدی از داده‌ها داریم و یک کلاس بندی نامتوازن است. (البته به بدی مثال اول نیست اما باز هم نامتوازن است) اگر تعداد کلاس‌ها بیشتر باشد این مشکل نامتوازن بودن هم افزایش می‌یابد. به این مسئله، یادگیری متوازن گفته می‌شود که یکی از مشکلات روش one-vs-all است.

## راه حل:

- ۱- یک راه برای حل این **نسبت دادن وزن‌هایی** به کلاس اقلیت و اکثریت است. بدین صورت که اهمیت کلاس اقلیت را زیاده‌تر و اهمیت کلاس اکثریت را کمتر جلوه دهد و فرآیند آموزش مدل با توجه به این وزن‌دهی صورت گیرد.
- ۲- **Random Under-Sampling**: به طور رندوم بعضی از سмпل‌های کلاسی که اکثریت دارد را حذف کنیم تا جایی که توازن میان کلاس‌ها برقرار شود.
- ۳- **Random Over-Sampling**: افزایش تعداد داده‌های کلاس اقلیت با جایگزین کردن و تکرار آن‌ها به صورت رندوم.
- ۴- **Cluster-Based Over Sampling**: الگوریتم خوشه‌بندی k-means به طور مجزا بر روی داده‌های کلاس اقلیت و اکثریت اجرا می‌شود تا خوشه‌های موجود در دیتاست مشخص شوند. پس از آن هر خوشه به گونه‌ای **oversample** می‌شود که تمام خوشه‌هایی که کلاس یکسان دارند دارای تعداد یکسانی داده باشند و سائز همه کلاس‌ها هم باهم برابر باشد.