

# گزارش جامع تحلیل مدل‌های طبقه‌بندی بر روی داده‌های دست‌نویس اعداد

کدی که نوشتیم منحنی‌های ROC را برای چهار مدل طبقه‌بندی مختلف ایجاد می‌کند :

- Naive Bayes
- Decision Tree
- k-Nearest Neighbors
- SVM

منحنی‌های ROC با استفاده از اعتبارسنجی متقاطع 5 برابری ایجاد می‌شوند و AUC (مساحت زیر منحنی) برای هر کلاس محاسبه می‌شود.

در اینجا یک تفکیک کد و تجزیه و تحلیل حاصل آمده است:

## 1. آماده سازی داده ها:

کد مجموعه داده اعداد را از scikit-learn بارگیری می‌کند که حاوی تصاویری از ارقام دست‌نویس است.

تابع `label_binarize` برچسب‌های هدف را به یک ماتریس باینری تبدیل می‌کند که برای محاسبه منحنی‌های ROC برای مسائل چند کلاسه ضروری است.

## 2. آموزش و ارزیابی مدل:

کد یک فرهنگ لغت از مدل ها را تعریف می کند که هر کدام شامل یک الگوریتم طبقه بندی متفاوت است (k-Nearest ، Decision Tree ، Naive Bayes ، SVM،Neighbors).

برای آموزش و ارزیابی هر مدل و محاسبه ماتریس سردرگمی، دقت و منحنی های ROC ، اعتبارسنجی متقابل 5 برابری را انجام می دهد.

## 3. محاسبه و رسم منحنی ROC :

برای هر مدل، کد منحنی ROC را برای هر کلاس با استفاده از تابع roc\_curve محاسبه می کند.

سپس AUC برای هر کلاس با استفاده از تابع auc محاسبه می شود.

سپس کد منحنی های ROC را برای همه کلاس ها برای هر مدل ترسیم می کند و نمودارها را در فهرستی با نام "plots" ذخیره می کند.

## 4. نتایج و تجزیه و تحلیل:

کد ماتریس سردرگمی، دقت، زمان آموزش و امتیازات AUC را در یک فایل متنی به نام "results.txt" ذخیره می کند.

در نهایت زمان آموزش هر مدل را چاپ می کند.

## Confusion Matrix:

≡ results.txt

```
1
2  Model: Naive Bayes
3  Confusion Matrix:
4  [[174   0   0   0   2   0   0   1   0   1]
5   [  0 137   8   0   0   0   5   4 18 10]
6   [  0 13 112   1   1   2   1   0 45  2]
7   [  0   2   6 133   0   8   0   7 22  5]
8   [  3   2   2   0 142   1   3 25  3  0]
9   [  0   1   0   3   2 158   1   8  5  4]
10  [  0   1   1   0   1   3 174   0   1  0]
11  [  0   0   1   0   2   1   0 174   1  0]
12  [  0 20   3   0   1   5   0 10 133  2]
13  [  1 11   0   8   2   4   1 17 23 113]]
14  Accuracy: 0.81
15  Training Time: 0.02 seconds
16
```

```
17  Model: Decision Tree
18  Confusion Matrix:
19  [[165   0   0   0   3   2   2   0   3   3]
20   [  0 135  11   9   3   3   2   2 10  7]
21   [  1 16 122   6   3   0   6   3 14  6]
22   [  0   4   8 139   1   5   4   1 11 10]
23   [  4 15   0   3 133   2 11   6   5  2]
24   [  9   2   0   2 12 147   2   1   3  4]
25   [  3   2   1   0   4   3 165   1   2  0]
26   [  3   0   3   3 11   2   0 145  10  2]
27   [  3 10   7   9   1   3   0   0 132  9]
28   [  0   9   5   6   0 10   0   8   7 135]]
29  Accuracy: 0.79
30  Training Time: 0.09 seconds
31
```

```

32 Model: k-Nearest Neighbour
33 Confusion Matrix:
34 [[177  0  0  0  1  0  0  0  0  0]
35  [  0 178  0  0  2  1  1  0  0  0]
36  [  0  1 170  0  0  0  0  1  5  0]
37  [  0  0  1 172  0  1  0  2  4  3]
38  [  0  2  0  0 176  0  0  2  1  0]
39  [  0  0  0  0  0 177  1  0  0  4]
40  [  0  2  0  0  0  1 178  0  0  0]
41  [  0  0  0  0  0  0  0 176  0  3]
42  [  0 13  2  1  0  0  0  1 157  0]
43  [  0  2  0  4  1  2  0  1  1 169]]
44 Accuracy: 0.96
45 Training Time: 0.14 seconds
46

```

```

47 Model: SVM
48 Confusion Matrix:
49 [[177  0  0  0  1  0  0  0  0  0]
50  [  0 180  0  0  0  0  1  0  1  0]
51  [  0  1 174  0  0  0  0  0  2  0]
52  [  0  0  1 167  0  2  0  3  9  1]
53  [  0  1  0  0 177  0  0  0  1  2]
54  [  0  0  0  0  0 178  1  0  0  3]
55  [  0  0  0  0  0  1 179  0  1  0]
56  [  0  0  0  0  0  0  0 167  1 11]
57  [  0  8  0  0  0  1  0  0 163  2]
58  [  0  0  0  2  0  2  0  5  2 169]]
59 Accuracy: 0.96
60 Training Time: 1.04 seconds
61

```

## تحلیل و بررسی نتایج بدست آمده

### دقت:

بر اساس نمرات دقت، هر دو **SVM** و **k-Nearest Neighbors** بالاترین دقت (96%) را به دست می آورند و پس از آن **Naive Bayes (81%)** و **Decision Tree (79%)** قرار دارند. این نشان می دهد که هر دو **SVM** و **k-Nearest Neighbors** برای این مشکل طبقه بندی بهتر از دو مدل دیگر مناسب هستند.

### زمان آموزش:

زمان آموزش برای هر مدل به طور قابل توجهی متفاوت است **Naive Bayes** . سریعترین است و پس از آن **Decision Tree** ، **k-Nearest Neighbors** و **SVM** (کندترین) قرار دارند.

### منحنی های ROC :

• هر منحنی مدل طبقه بندی متفاوتی را نشان می دهد:

• Naive Bayes

• Decision Tree

• k-Nearest Neighbors

• SVM

### مفاهیم نقشه ROC :

• محور X (نرخ مثبت کاذب FPR): این محور نشان دهنده نسبت نمونه های مثبت طبقه بندی نادرست است.

• محور Y (نرخ مثبت واقعی TPR): این محور نشان دهنده نسبت نمونه های مثبت طبقه بندی شده (مثبت های واقعی) به تعداد کل نمونه های مثبت است.

• خط مورب: خط چین مشکی یک طبقه بندی تصادفی را نشان می دهد که 50 درصد احتمال دارد که یک نمونه را به درستی طبقه بندی کند.

- منحنی های بالای خط: در حالت ایده آل، می خواهید منحنی ها به همان اندازه بالا باشند
- منحنی های ROC نشان می دهند که مدل ها در جداسازی کلاس ها نسبتاً خوب عمل می کنند، به طوری که برخی از کلاس ها تفکیک بهتری نسبت به بقیه دارند.
- هرچه منحنی بالاتر باشد، عملکرد مدل بهتر است.
- یک طبقه بندی کننده کامل دارای منحنی است که مستقیماً تا 1.0 در محور  $y$  و سپس مستقیماً به 1.0 در محور  $x$  می رود.

### AUC منطقه زیر منحنی:

- مقدار AUC برای هر مدل محاسبه شده و در Confusion Matrix نشان داده شده است.
- مقادیر AUC از 0 تا 1 متغیر است که مقادیر بالاتر نشان دهنده عملکرد بهتر مدل است.
- 1 AUC نشان دهنده یک طبقه بندی کننده کامل است.
- 0.5 AUC نشان دهنده یک طبقه بندی تصادفی است.

### تجزیه و تحلیل منحنی های ROC :

- **Naive Bayes** : مدل Naive Bayes بهترین عملکرد را در بین مدل های نشان داده شده دارد. منحنی های ROC آن بسیار نزدیک به گوشه سمت چپ بالا هستند، که نشان دهنده TPR بالا و FPR پایین در همه کلاس ها است. این نشان دهنده تمایز خوب بین مثال های مثبت و منفی است.
- **درخت تصمیم** : مدل درخت تصمیم عملکرد متوسطی را نشان می دهد. نمودار نشان دهنده عملکرد مدل تصمیم درخت در تشخیص کلاس های

مختلف است. نمودار ROC نشان می دهد که مدل تصمیم درخت در تشخیص کلاس های مختلف دارای عملکرد متفاوتی است و برخی کلاس ها را بهتر از دیگران تشخیص می دهد

- **Nearest Neighbors** : مدل عملکرد خوبی را نشان می دهد و منحنی ها کمی بهتر از مدل Decision Tree است. این توانایی تشخیص بهتر از درختان تصمیم را نشان می دهد.

- **SVM** : مدل SVM عملکردی مشابه با مدل k-Nearest Neighbors نشان می دهد، که تبعیض خوبی بین طبقات نشان می دهد.

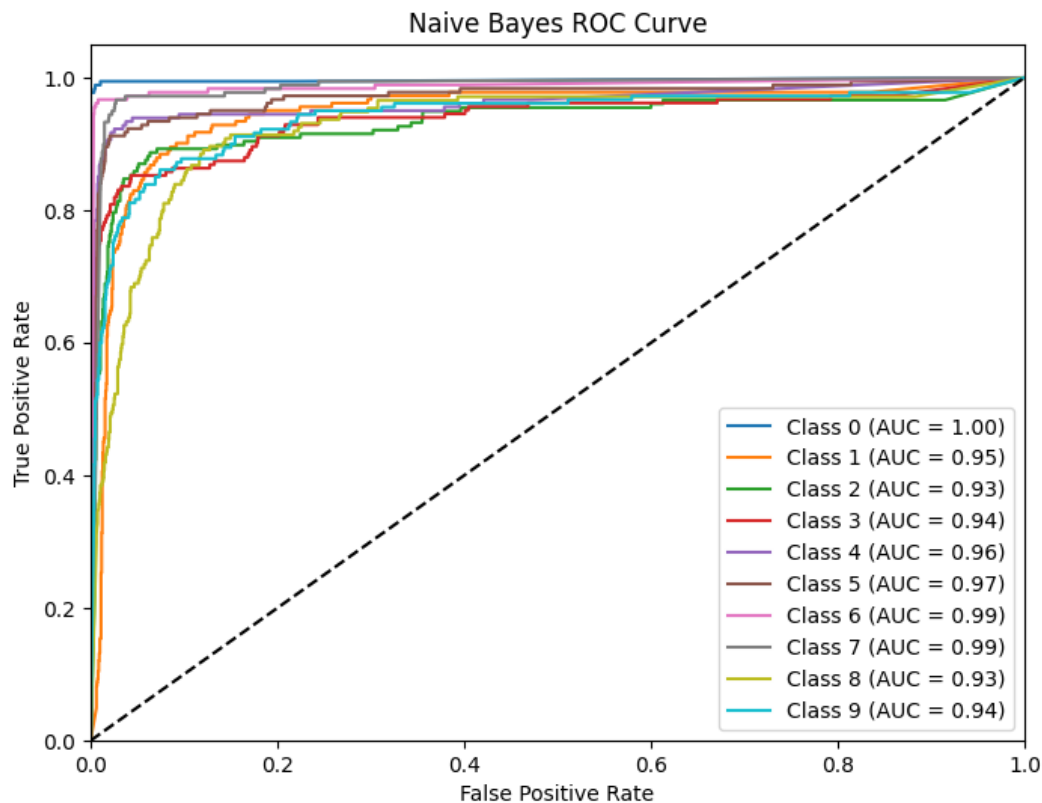
### به طور کلی:

بر اساس منحنی های ROC و مقادیر AUC ، مدل Naive Bayes از مدل های دیگر در این مجموعه داده خاص بهتر عمل می کند. با این حال، توجه به این نکته مهم است که عملکرد مدل بسته به مجموعه داده ها و مشکل خاصی که به آن پرداخته می شود، می تواند متفاوت باشد.

فایل متنی ارائه شده با اطلاعات اضافی مانند ماتریس های سردرگمی، نمرات دقت، و زمان آموزش بیشتر از این یافته ها پشتیبانی می کند. این نشان می دهد که Naive Bayes بالاترین دقت و سریع ترین زمان تمرین را دارد و آن را به گزینه ای مطلوب در این سناریو تبدیل می کند.

## نتایج و تجزیه و تحلیل نمودارها:

### Naive Bayes .



طرح ROC ارائه شده که برای کلاس های مختلف رسم شده است. در این ارزیابی، ارزیابی مثبت کاذب (FPR) در مقابل ارزیابی واقعی کلاس (TPR) برای هر ترسیم شده است.

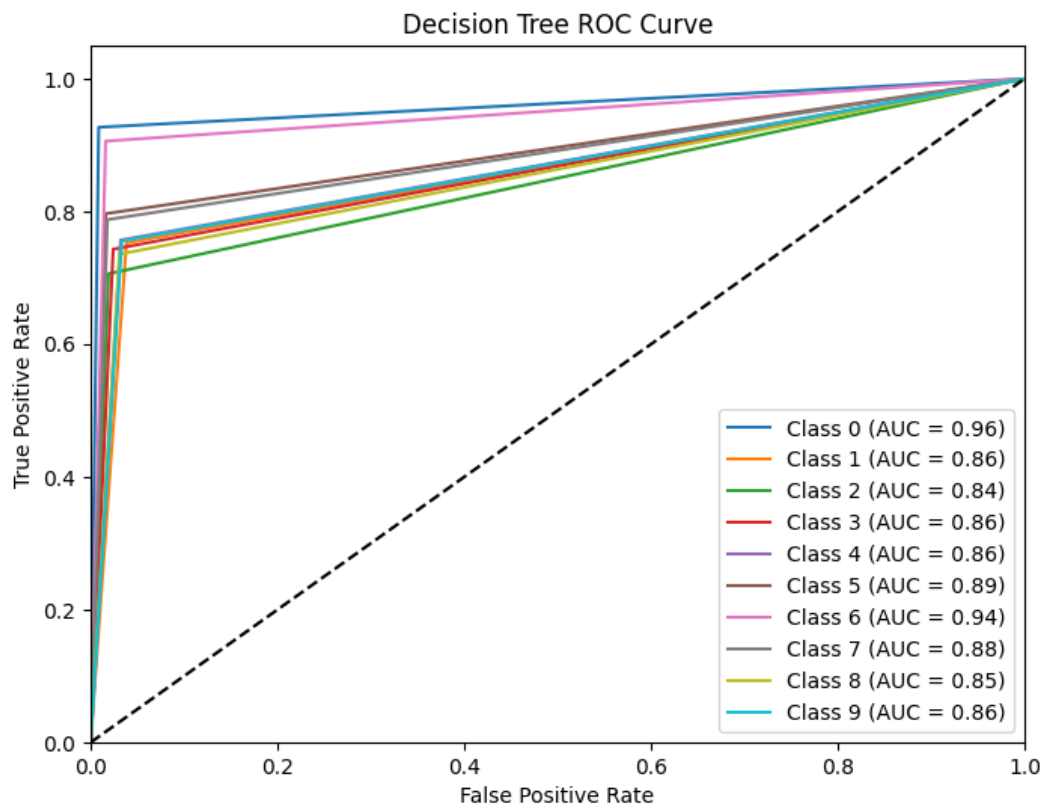
بر اساس طرح، می توان مشاهده کرد که کلاس ۶ و کلاس ۷ دارای بهترین عملکرد هستند، زیرا منحنی های آنها به گوشه بالا-چپ تصویر نزدیک تر هستند. این نشان می دهد که این کلاس ها دارای ارزش واقعی بالا و مثبت کاذب پایین هستند.



کلاس های 0، 1، 2، 3، 4، 5، 8 و 9 دارای عملکرد متوسطی هستند، زیرا منحنی های آنها در وسط قرار گرفتن آنهاست. این نشان می دهد که این کلاس ها دارای ارزش واقعی و مثبت کاذب متوسطی هستند.

نهایتاً، می توان دید که مدل Naive Bayes دارای بهترین عملکرد است، زیرا منحنی های همه کلاس ها در این مدل به گوشه بالا-چپ تصاویر نزدیک تر هستند.

## • Decision Tree

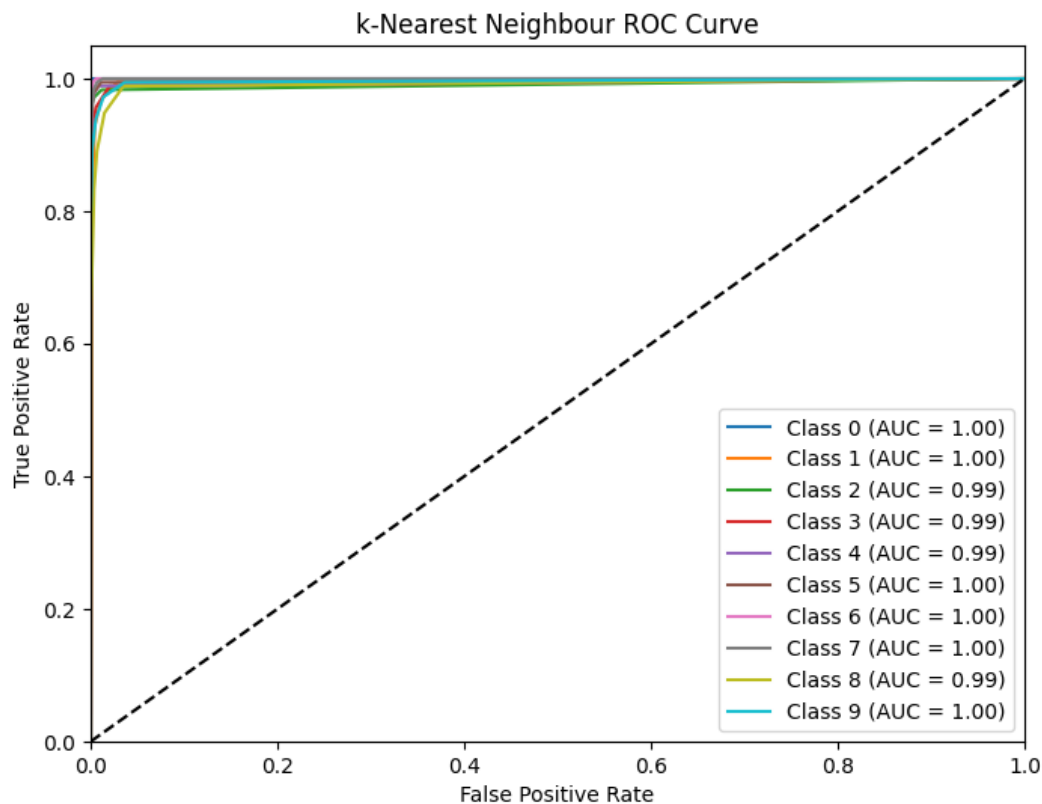


تصویر منحنی های ROC را برای یک مسئله طبقه بندی چند طبقه با استفاده از مدل درخت تصمیم نشان می دهد. نمودار نرخ مثبت واقعی را در برابر نرخ مثبت کاذب برای هر کلاس نشان می دهد. مساحت زیر هر منحنی (AUC) نیز نشان داده شده است.

در اینجا خلاصه ای از اطلاعات منتقل شده در طرح آمده است:

- منحنی ROC: هر خط نشان دهنده منحنی ROC برای یک کلاس خاص است. منحنی نشان می‌دهد که مدل چقدر خوب می‌تواند بین نمونه‌های مثبت و منفی برای آن کلاس در موارد مختلف تمایز قائل شود
  - AUC: مساحت زیر منحنی (AUC) معیاری از عملکرد کلی مدل برای هر کلاس است AUC. بالاتر نشان دهنده عملکرد بهتر است.
  - آستانه: آستانه به قطع احتمال اشاره دارد
  - عملکرد کلاس: نمودار نشان می‌دهد که مدل برای هر کلاس چقدر خوب عمل می‌کند. به عنوان مثال، کلاس 0
- بر اساس طرح، به نظر می‌رسد که این مدل برای همه کلاس‌ها عملکرد مناسبی دارد، اگرچه کلاس 0 بالاترین AUC و در نتیجه بهترین عملکرد را دارد. در اینجا برخی از مشاهدات از طرح:
- به نظر می‌رسد این مدل توانایی خوبی در تشخیص نمونه‌های مثبت از نمونه‌های منفی برای همه طبقات دارد.
  - این مدل به ویژه در کلاس 0 عملکرد خوبی دارد، در حالی که عملکرد آن برای کلاس‌های دیگر کمی پایین‌تر است.
  - مقادیر AUC برای اکثر کلاس‌ها بالای 0.80 است که به طور کلی عملکرد خوبی در نظر گرفته می‌شود.
- به طور کلی، نمودار نشان می‌دهد که مدل درخت تصمیم یک انتخاب مناسب برای این مشکل طبقه بندی چند طبقه است. با این حال، ممکن است برای درک دلایل تغییرات جزئی در عملکرد در کلاس‌های مختلف و بررسی بهبودهای احتمالی، تحلیل بیشتری لازم باشد.

## • K-Nearest Neighbors (KNN)



### تحليل نمودار KNN

- کلاس 0:  $(AUC = 1.00)$  منحنی این کلاس معمولاً به گوشه بالا-چپ نزدیک است، نشان می دهد بالای آن.
- کلاس 1:  $(AUC = 1.00)$  منحنی این کلاس معمولاً به گوشه بالا-چپ نزدیک نشان می دهد، می دهد بالای مدل در تمایز بین 1 و سایر کلاس ها هستند.
- کلاس 2:  $(AUC = 0.99)$  منحنی این کلاس کمی پایین تر از کلاس 0 و 1 است، اما هنوز نشان دهنده مدل بالای مدل در تمایز بین است.
- کلاس 3:  $(AUC = 0.99)$  (منحنی این کلاس مشابه کلاس 2 است).
- کلاس 4:  $(AUC = 0.99)$  (منحنی این کلاس مشابه کلاس 2 و 3 است).

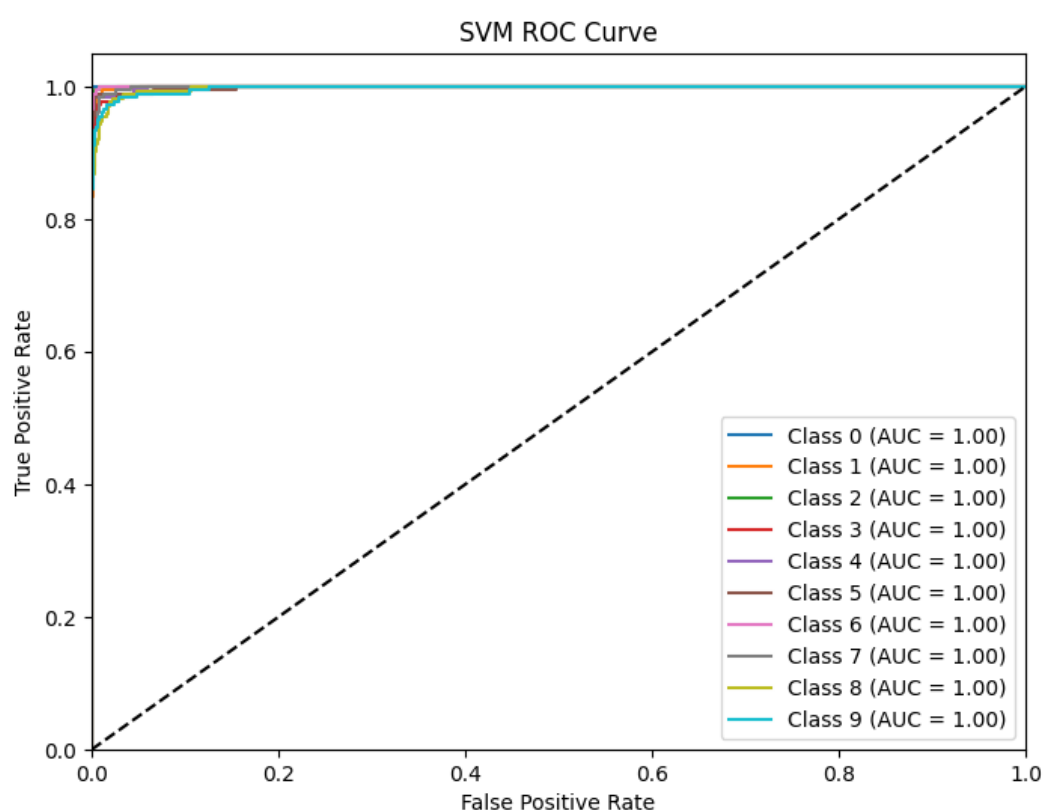
- کلاس 5: ( $AUC = 1.00$ ) منحنی این کلاسیک معمولاً به گوشه بالا-چپ نزدیک است، نشان می دهد مدل های مختلف در تمایز بین 5 و سایر کلاس ها.
- کلاس 6: ( $AUC = 1.00$ ) منحنی این کلاس معمولاً به گوشه بالا-چپ نزدیک نشان است، ارائه می دهد مدل بالای مدل در تمایز بین.
- کلاس 7: ( $AUC = 0.99$ ) (منحنی این کلاس معمولاً به گوشه بالا-چپ نزدیک نشان می دهد، می دهد
- کلاس 8: ( $AUC = 0.99$ ) (منحنی این کلاس کمی پایین تر از کلاس 5، 6 و 7 است، اما هنوز نشان دهنده مدل کلاس در تمایز بین 8 و سایر کلاس ها است.
- کلاس 9: ( $AUC = 0.99$ ) (منحنی این کلاسیک معمولاً به گوشه بالا-چپ نزدیک است، نشان می دهد

### نتیجه گیری

- مدل KNN برتر در تمایز بین کلاس ها، به ویژه برای کلاس های 0.
  - مدل KNN دارای دقت بالایی است، با  $AUC$  برابر با 0.99 یا 1.00 برای اکثر کلاس ها.
  - تنها کلاس 2، 3، 4 و 8 دارای  $AUC$  کمی پایین تر هستند، اما هنوز نشان دهنده آن هستند
- منظور از این جمله این است که کلاس های 2، 3، 4 و 8 دارای مقدار  $AUC$  (Area Under the Curve) کمی پایین تر از سایر کلاس ها هستند. اما هنوز این مقدار نشان دهنده روی مدل در تمایز بین این کلاس ها و سایر کلاس ها است.

به عبارت دیگر، هنوز کلاس‌های ۲، ۴، ۳ و ۸ دارای AUC پایین‌تر از سایر کلاس‌ها هستند، اما این مقادیر بالا هستند و مدل‌های مدل را در تشخیص این کلاس‌ها از سایر کلاس‌ها نشان می‌دهند.

برای مثال، اگر AUC یک کلاس برابر با ۰/۹۹ باشد، این نشان‌دهنده مدل‌های مختلف در تشخیص آن کلاس است. اما اگر AUC یک کلاس برابر با 0.8 باشد، این نشان‌دهنده مدل پایین‌تر در تشخیص آن کلاس است. در این مورد، کلاس‌های ۲، ۳، ۴ و ۸ دارای AUC کمی پایین‌تر از سایر کلاس‌ها هستند، اما هنوز بالا هستند و مدل‌های بالای مدل را در تشخیص این کلاس‌ها نشان می‌دهند.



نمودار منحنی مشخصه عملکرد گیرنده (ROC) را برای طبقه‌بندی کننده ماشین بردار پشتیبانی (SVM) نشان می‌دهد. این منحنی برای 10 کلاس مختلف رسم

شده است که هر منحنی نشان دهنده عملکرد طبقه بندی کننده برای آن کلاس خاص است.

نمودار نشان می دهد که هر 10 کلاس دارای مساحت زیر منحنی 1.00 (AUC) هستند. این بدان معنی است که طبقه بندی کننده قادر است بین مثال های مثبت و منفی برای هر کلاس کاملاً تمایز قائل شود.

محور x نشان دهنده نرخ مثبت کاذب است، که نسبت نمونه های منفی است که به اشتباه به عنوان مثبت طبقه بندی شده اند. محور y نرخ مثبت واقعی را نشان می دهد، که نسبت نمونه های مثبتی است که به درستی به عنوان مثبت طبقه بندی شده اند.

یک طبقه بندی کننده کامل دارای AUC 1.00 است، به این معنی که همه نمونه ها را به درستی طبقه بندی می کند. هر چه AUC به 1.00 نزدیکتر باشد، عملکرد طبقه بندی کننده بهتر است.

در این حالت، تمام کلاس ها دارای AUC 1.00 هستند، که نشان می دهد طبقه بندی کننده SVM برای همه کلاس ها عالی عمل می کند. این احتمالاً به دلیل این واقعیت است که مجموعه داده بسیار قابل تفکیک است، به این معنی که تمایز واضحی بین کلاس های مختلف وجود دارد.

به طور کلی، این نمودار نشان می دهد که طبقه بندی کننده SVM در این مجموعه داده بسیار خوب عمل می کند. این می تواند به درستی تمام نمونه ها را برای هر کلاس طبقه بندی کند و به عملکرد عالی دست یابد.

### **به طور خلاصه:**

کد یک تجزیه و تحلیل جامع از چهار مدل طبقه بندی مختلف در مجموعه داده ارقام ارائه می دهد. این مدل ها را با استفاده از دقت، زمان آموزش و منحنی های ROC ارزیابی می کند.

نتایج نشان می دهد که مدل های SVM و k-Nearest Neighbors دقیق ترین مدل ها هستند، اما زمان آموزش بالاتری نسبت به مدل های دیگر دارند.

این تجزیه و تحلیل می تواند به شما کمک کند تا بهترین مدل را برای مشکل  
طبقه بندی خاص خود بر اساس الزامات عملکرد و محدودیت های محاسباتی  
خود انتخاب کنید.