
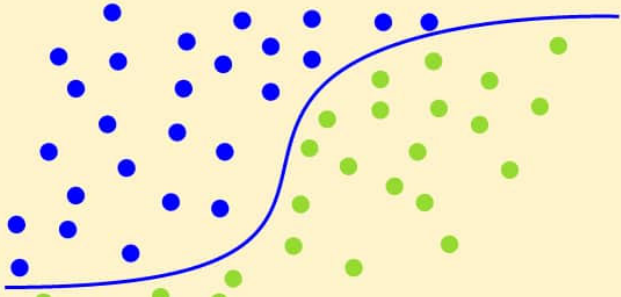


Artificial Intelligence and Soft computing



K.N. TOOSI  
University of  
Technology

# طبقه بندی چند کلاسی




## Multi Classification

Advanced Application of Artificial Intelligence  
and Digital Transformation



K.N. Toosi  
University of  
Technology



کاربرد پیشرفته هوش مصنوعی و تحول دیجیتال

5G

AI

Hasan Ghasemzadeh  
<http://wp.kntu.ac.ir/ghasemzadeh>

### نمونه رگرسیون لجستیک با بیش از یک ویژگی

Age	Heart Rate	Heart Attack
45	50	N
50	50	N
55	50	Y
60	50	N
65	70	N
70	70	Y
75	90	Y
80	90	Y
85	90	N
90	90	Y
95	90	Y

مثال: احتمال حمله قلبی با دو ویژگی ضربان قلب و سن فرد احتمال حمله قلبی شخصی ۲۰ ساله با ضربان قلب ۸۰ چقدر است؟

بسیاری از عوامل دیگر مانند فشار خون، سطح کلسترول، سابقه خانوادگی، و عادات زندگی (مثل سیگار کشیدن و ورزش نکردن) نیز در تعیین خطر حمله قلبی مؤثرند، که در این جدول لحاظ نشده است.

$$P(Y = 1 | X = x) = p(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots)}}$$

$$P(\text{Heart Attack} = 1 | X = \text{Age}, \text{Heart rate}) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 \cdot \text{Age} + \theta_2 \cdot \text{Heart rate})}}$$

Adv. App. of AI and DT

3

### نمونه رگرسیون لجستیک با بیش از یک ویژگی

Age	Heart Rate	Heart Attack
45	50	N
50	50	N
55	50	Y
60	50	N
65	70	N
70	70	Y
75	90	Y
80	90	Y
85	90	N
90	90	Y
95	90	Y

مثال: احتمال حمله قلبی با دو ویژگی ضربان قلب و سن فرد

بسیاری از عوامل دیگر مانند فشار خون، سطح کلسترول، سابقه خانوادگی، و عادات زندگی (مثل سیگار کشیدن و ورزش نکردن) نیز در تعیین خطر حمله قلبی مؤثرند، که در این جدول لحاظ نشده است.

$$P(Y = 1 | X = x) = p(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots)}}$$

$$P(\text{Heart Attack} = 1 | X = \text{Age}, \text{Heart rate}) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 \cdot \text{Age} + \theta_2 \cdot \text{Heart rate})}}$$

Adv. App. of AI and DT

4

### نمونه رگرسیون لجستیک با بیش از یک ویژگی

فایل: 8 regression logistic two feature.py

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split

# Define the dataset
data = {
    'Age': [45, 50, 55, 60, 65, 70, 75, 80, 85, 90, 95],
    'Heart Rate': [50, 50, 50, 50, 70, 70, 90, 90, 90, 90, 90],
    'Heart Attack': [0, 0, 1, 0, 0, 1, 1, 1, 0, 1, 1]
}
df = pd.DataFrame(data)

# Prepare the data
X = df[['Age', 'Heart Rate']]
y = df['Heart Attack']
```

Adv. App. of AI and DT

5

### نمونه رگرسیون لجستیک با بیش از یک ویژگی

```
# Train the logistic regression model
model = LogisticRegression()
model.fit(X, y)

# model's parameters
print("Coefficients:", model.coef_)
print("Intercept:", model.intercept_)
theta_0 = model.intercept_[0]
theta_1 = model.coef_[0][0]
theta_2 = model.coef_[0][1]
print("Intercept (θ0):", theta_0)
print("Coefficient for Age (θ1):", theta_1)
print("Coefficient for Heart Rate (θ2):", theta_2)
```

```
Coefficients: [[-0.17241598  0.14740334]]
Intercept: [0.89730001]
Intercept (θ0): 0.8973000136808484
Coefficient for Age (θ1): -0.17241597828106092
Coefficient for Heart Rate (θ2): 0.14740333833885552
```

Adv. App. of AI and DT

6

### نمونه رگرسیون لجستیک با بیش از یک ویژگی

```
# Example prediction for a person with Age = 20 and Heart Rate = 80
example_data = np.array([[20, 80]])
probability_heart_attack = model.predict_proba(example_data)[0][1]
print("Probability of Heart Attack for Age=20 and Heart Rate=80:",
      probability_heart_attack)
```

Probability of Heart Attack for Age=20 and Heart Rate=80: 0.9999030528396969

```
Intercept (θ0): 0.8973000136808484
Coefficient for Age (θ1): -0.17241597828106092
Coefficient for Heart Rate (θ2): 0.1474033383388552
```

$$P(\text{Heart Attack} = 1 | X = \text{Age}, \text{Heart rate}) = \frac{1}{1 + e^{-(0.8973 - 0.1724 \cdot \text{Age} + 0.1474 \cdot \text{Heart rate})}}$$

Adv. App. of AI and DT

7

### ارزیابی الگوریتم

پس از مدلسازی و تحلیل داده ها نیاز به ارزیابی الگوریتم است

**اگر الگوریتم ارزیابی نشود نمی توان آنرا بهبود بخشید.**



Adv. App. of AI and DT

8

## ارزیابی الگوریتم

مثال: الگوریتمی بیمارهای دارای سرطان را شناسایی می کند و نتایج زیر را داده است.

Individual Number	1	2	3	4	5	6	7	8	9	10	11	12
Actual Classification	1	1	1	1	1	1	1	1	0	0	0	0
Predicted Classification	0	0	1	1	1	1	1	1	1	0	0	0

تعداد جوابهای غلط: ۳

تعداد جوابهای صحیح: ۹

صحت و دقت (Accuracy): نسبت جوابهای صحیح مدل به کل نمونه ها ۷۵٪

Adv. App. of AI and DT

9

## ارزیابی الگوریتم

پس از مدلسازی و تحلیل داده های دودویی نتایج را به صورت زیر می توان دسته بندی نمود:

### مثبت صحیح (TP: True Positive):

تعداد حالت هایی که مقدار پیش بینی شده با مقدار واقعی برابر بوده و متغیر پاسخ برابر با یک است. (شناسایی درست - معیار مثبت).

### مثبت کاذب (FP: False Positive):

تعداد حالت هایی که مقدار پیش بینی شده برای متغیر پاسخ یک و مقدار واقعی برابر صفر است. (شناسایی غلط - معیار منفی).

### منفی صحیح (TN: True Negative):

تعداد حالت هایی که مقدار پیش بینی شده با مقدار واقعی برابر بوده و متغیر پاسخ برابر با صفر است. (مردودی درست - معیار مثبت).

### منفی کاذب (FN: False Negative):

تعداد حالت هایی که مقدار پیش بینی شده برای متغیر پاسخ صفر و مقدار واقعی برابر یک است. (مردودی غلط - معیار منفی).

Adv. App. of AI and DT

10

### ارزیابی الگوریتم

مثال: الگوریتمی بیمارهای دارای سرطان را شناسایی می کند و نتایج زیر را داده است.

Individual Number	1	2	3	4	5	6	7	8	9	10	11	12
Actual Classification	1	1	1	1	1	1	1	1	0	0	0	0
Predicted Classification	0	0	1	1	1	1	1	1	1	0	0	0
Result	FN	FN	TP	TP	TP	TP	TP	TP	FP	TN	TN	TN

مثبت صحیح (TP: True Positive)      مثبت کاذب (FP: False Positive)  
 منفی صحیح (TN: True Negative)      منفی کاذب (FN: False Negative)

چگونه از این شاخص ها استفاده کنیم؟

### ارزیابی الگوریتم

#### ماتریس اختلاط یا ماتریس درهم ریختگی (Confusion Matrix or Error Matrix)

به منظور سنجش مدل از یک ماتریس دو در دو به نام ماتریس درهم ریختگی استفاده می شود

مثبت صحیح (TP)      مثبت کاذب (FP)  
 منفی صحیح (TN)      منفی کاذب (FN)

		Predicted condition	
		Positive (PP) $PP = TP + FP$	Negative (PN) $PN = FN + TN$
Actual condition	Positive (P) $P = TP + FN$	True positive (TP)	False negative (FN)
	Negative (N) $N = FP + TN$	False positive (FP)	True negative (TN)

ماتریس درهم ریختگی

### ارزیابی الگوریتم

ماتریس درهم ریختگی برای مثال بیماران سرطانی

Individual Number	1	2	3	4	5	6	7	8	9	10	11	12
Actual Classification	1	1	1	1	1	1	1	1	0	0	0	0
Predicted Classification	0	0	1	1	1	1	1	1	1	0	0	0
Result	FN	FN	TP	TP	TP	TP	TP	TP	FP	TN	TN	TN

		Predicted condition	
		Positive (PP)	Negative (PN)
Actual condition	Total population = P + N		
	Positive (P)	True positive (TP)	False negative (FN)
		Predicted condition	
		Cancer	Non-cancer
Actual condition	Total 8 + 4 = 12	7	5
	Cancer 8	6	2
		Actual condition	
		Cancer	Non-cancer
		4	3

Adv. App. of AI and DT

13

### ارزیابی الگوریتم

معیارها در ماتریس اختلاط یا ماتریس درهم ریختگی

به منظور اندازه گیری کارایی مدل از معیارهای زیر برای حالت مثبت (مقادیر یک) استفاده می شود

**صحت یا حساسیت (Recall, Sensitivity)**

$$TPR (Recall) = \frac{TP}{P} = \frac{TP}{TP+FN}$$

**Hit rate, True Positive Rate (TPR)**

صحت برابرست با تعداد مثبت صحیح به تعداد صحیح در واقعیت

این معیار نسبت پیش بینی مثبت درست مدل به کل مقادیر مثبت درست واقعی را نشان می دهد

**دقت (Precision)**

$$PPV (Precision) = \frac{TP}{PP} = \frac{TP}{TP+FP}$$

**Positive Predictive Value (PPV)**

دقت برابرست با تعداد مثبت صحیح به تعداد صحیح در مدل (پیش بینی)

این معیار نسبت پیش بینی مثبت درست مدل به کل مقادیر مثبت درست مدل را نشان می دهد

Adv. App. of AI and DT

14

## ارزیابی الگوریتم

### معیارها در ماتریس اختلاط یا ماتریس درهم ریختگی

مشابه حالت قبل از معیارهای زیر برای حالت منفی (مقدار صفر) استفاده می شود

$$TNR = \frac{TN}{N} = \frac{TN}{TN+FP}$$

#### Specificity

True Negative Rate (TNR), Selectivity

برابرست با تعداد منفی صحیح به تعداد منفی در واقعیت  
این معیار در برابر حساسیت است

$$NPV = \frac{TN}{TN+FN}$$

#### Negative Predictive Value (NPV)

Adv. App. of AI and DT

15

## ارزیابی الگوریتم

### معیارها در ماتریس اختلاط یا ماتریس درهم ریختگی

به منظور اندازه گیری کارایی مدل از معیارهای زیر نیز استفاده می شود

$$\text{Support} = P$$

#### تکیه گاه (Support)

تکیه گاه تعداد مشاهده مقدار ۱ (هر کلاسی از داده ها) در مشاهدات است.

$$F_1 = \frac{1}{\frac{1}{2\text{Precision}} + \frac{1}{2\text{Recall}}} = 2 \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

#### اندازه $F_1$

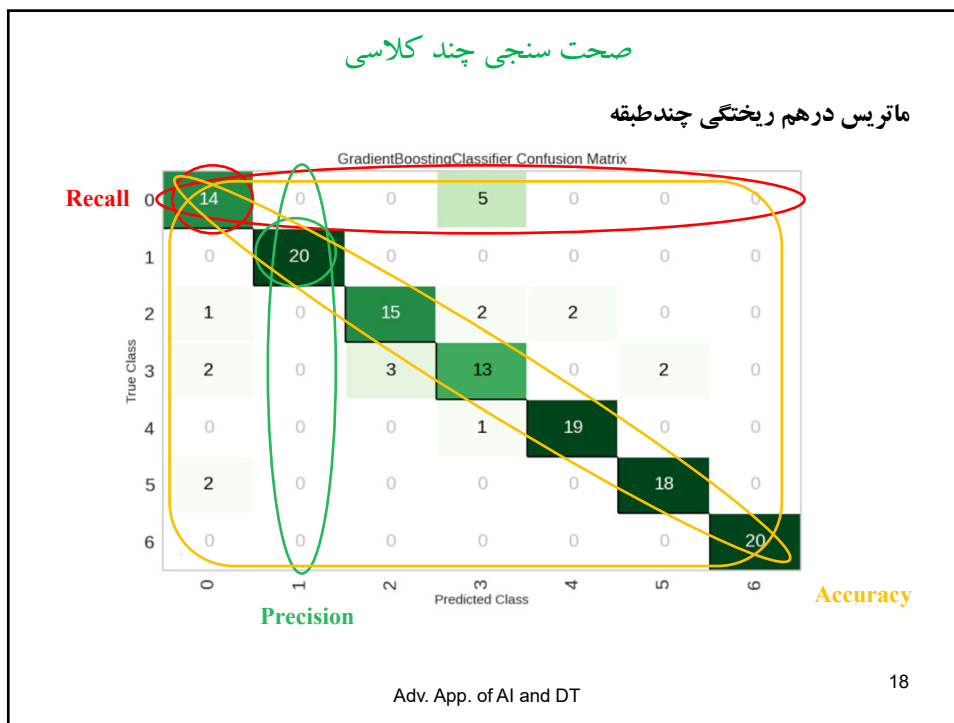
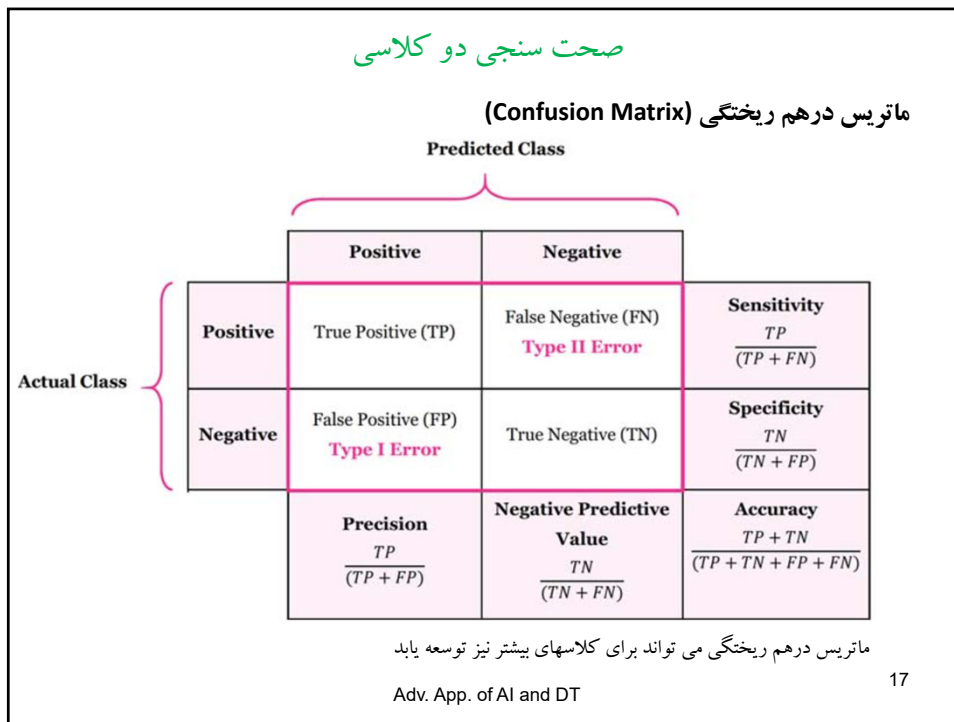
میانگین همساز دو معیار دقت و صحت است.  
مقدار صفر بدترین حالت و مقدار یک بهترین  
حالت است و تعادل دقت و صحت را دربردارد

$$\text{ACC} = \frac{TP+TN}{P+N} = \frac{TP+TN}{TP+TN+FP+FN}$$

#### صحت و دقت (Accuracy)

Adv. App. of AI and DT

16



## صحت سنجی چند کلاسی

## ماتریس درهم ریختگی چندطبقه

**Recall** is the fraction of samples in class  $i$  that the model correctly classified as class  $i$

$$\text{Recall} = \frac{C_{ii}}{\sum_j C_{ij}}$$

$C_{ii}$ : The count of samples correctly classified as class  $i$  (True Positives for class  $i$ ).

$\sum_j C_{ij}$ : The total number of samples that actually belong to class  $i$ , regardless of how they were classified (True Positives + False Negatives for class  $i$ ).

**Precision** is the fraction of samples that the model assigned to class  $i$  that actually belong to class  $i$ .

$$\text{Precision} = \frac{C_{ii}}{\sum_j C_{ji}}$$

$C_{ii}$ : The count of samples correctly classified as class  $i$  (True Positives for class  $i$ ).

$\sum_j C_{ji}$ : The total number of samples predicted as class  $i$ , regardless of their actual class (True Positives + False Positives for class  $i$ ).

**Accuracy** is the overall fraction of correctly classified samples across all classes

$$\text{Accuracy} = \frac{\sum_i C_{ii}}{\sum_i \sum_j C_{ij}}$$

$\sum_i C_{ii}$ : The sum of True Positives across all classes (correctly classified samples for each class).

$\sum_i \sum_j C_{ij}$ : The total number of samples across all classes (including both correct and incorrect classifications).

19

## صحت سنجی چند کلاسی

## میانگین گیری بدون وزن بین چند کلاس

- **Macroaveraging**: treats each class equally by computing the performance measure (like F1 score) for each class independently and then averaging them.
- **Implication**: Macroaveraging gives equal weight to all classes, regardless of the number of samples in each class. This is useful when you want each class to have an **equal impact on the overall metric**, even if some classes have fewer samples.

## میانگین گیری وزن دار بین چند کلاس

- **Microaveraging**: Microaveraging aggregates all true positives (TP), false positives (FP), and false negatives (FN) across all classes and then computes a single performance measure.
- **Implication**: Microaveraging takes into account the imbalance of class distributions by weighing each sample equally. It is especially useful in cases where **class imbalance is significant**, as it doesn't prioritize one class over another.

## صحت سنجی رگرسیون لجستیک

### ماتریس درهم ریختگی (Confusion Matrix)

به منظور اندازه گیری کارایی مدل ابتدا داده های مساله را به دو قسمت آموزش و تست تقسیم می کنیم. در دستور زیر داده های تست ۳۰ درصد داده های مساله در نظر گرفته شده اند

```
# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
```

محاسبه معیارها

```
from sklearn import metrics
print(metrics.classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	0.91	0.97	0.94	10981
1	0.52	0.25	0.34	1376

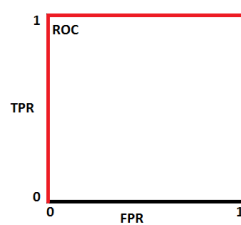
Adv. App. of AI and DT

21

## صحت سنجی رگرسیون لجستیک

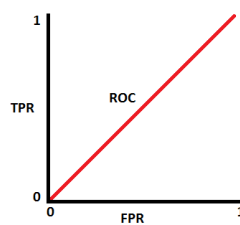
### منحنی ROC (Receiver Operating Characteristic)

تحلیل بازدهی مدل برای جدا کردن دو دسته از یکدیگر است. و با رسم TPR در برابر FPR نشان داده می شود



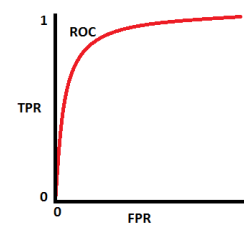
Ideal situation.

AUC = 1.0, It is perfectly able to distinguish between positive class and negative class.



Worst situation (Random predictions)

AUC = 0.5, model has no discrimination capacity to distinguish between positive class and negative class.



Usual situation

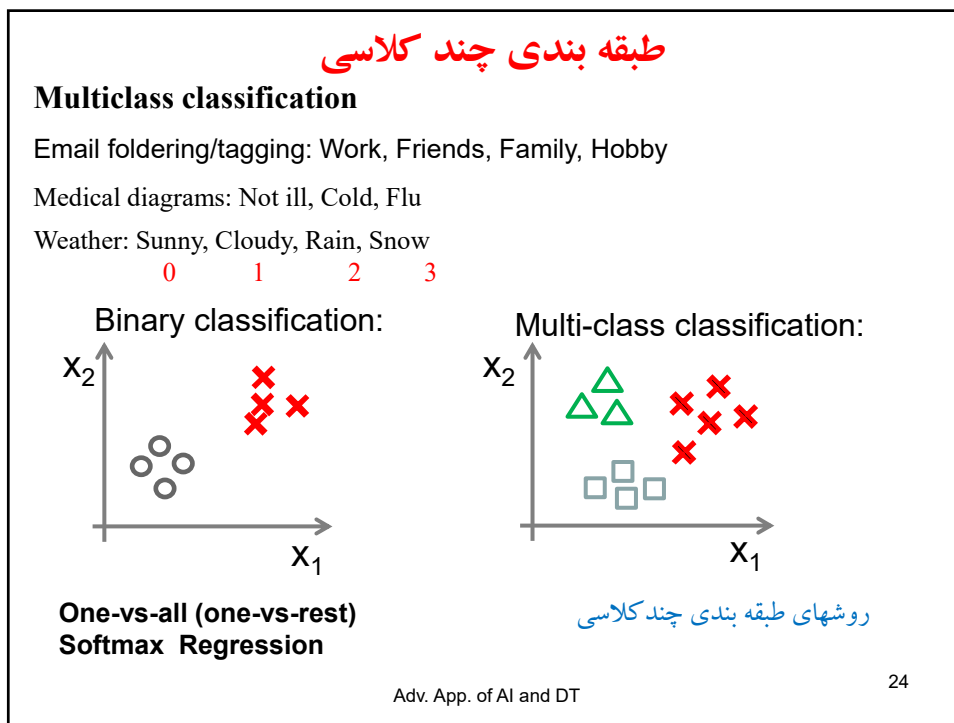
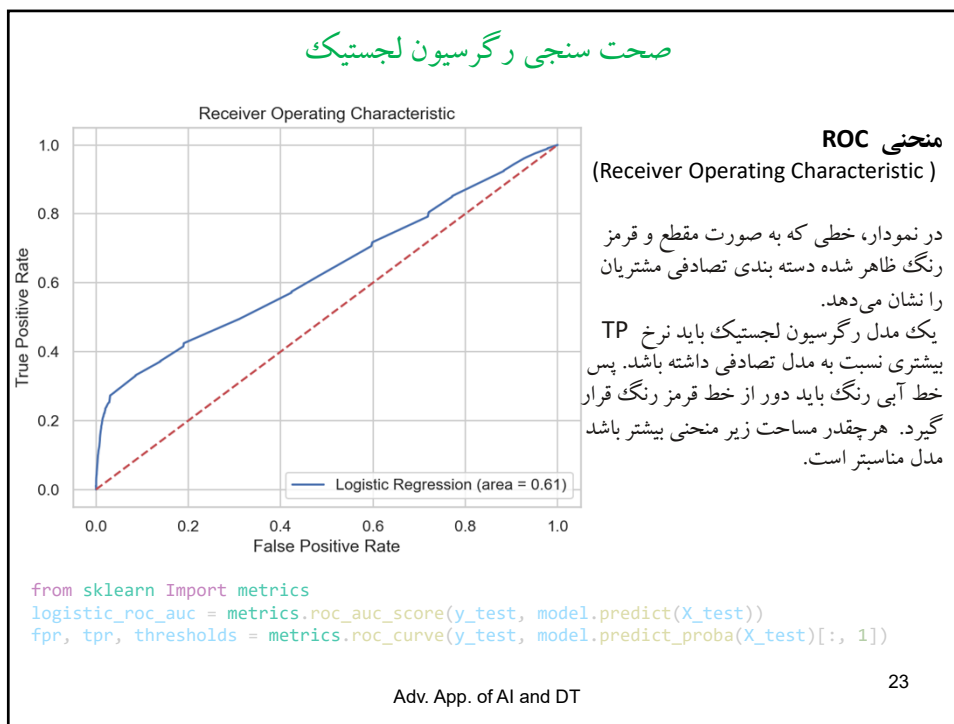
$0.5 < \text{AUC} < 1$

### مساحت زیر منحنی AUC

The Area Under the Curve (AUC) is the measure of the ability of a classifier to distinguish between classes and is used as a summary of the ROC curve.

Adv. App. of AI and DT

22



### طبقه بندی چند کلاسی

**One-vs-all (one-vs-rest):**

Class 1: △  
 Class 2: □  
 Class 3: ×

$$h_{\theta}^{(i)}(x) = P(y = i|x; \theta) \quad (i = 1, 2, 3)$$

25

Adv. App. of AI and DT

### طبقه بندی چند کلاسی

**Softmax Regression:**

یک تابع ریاضی است برای مسائل طبقه بندی چند کلاسی که یک بردار از اعداد حقیقی (مقادیر خام مدل) را به یک توزیع احتمال تبدیل می کند

$$\text{Softmax}(z_j) = \frac{e^{z_j}}{\sum_{j=1}^K e^{z_j}} \quad K: \text{تعداد کلاس ها با ابعاد } z$$

z = [2.0, 1.0, 0.1]      نمونه ورودی

$e^{2.0} = 7.389$	$\text{Softmax}(z_1) = \frac{7.389}{11.212} \approx 0.659$
$e^{1.0} = 2.718$	$\text{Softmax}(z_2) = \frac{2.718}{11.212} \approx 0.242$
$e^{0.1} = 1.105$	$\text{Softmax}(z_3) = \frac{1.105}{11.212} \approx 0.099$
$\sum_{j=1}^K e^{z_j} = 11.212$	

Softmax(z) = [0.659, 0.242, 0.099]      کلاس ۱ با احتمال ۶۵,۹٪ به عنوان پیش بینی نهایی انتخاب می شود

26

Adv. App. of AI and DT

## طبقه بندی چند کلاسی

Train a logistic regression classifier  $h_{\theta}^{(i)}(x)$  for each class  $i$  to predict the probability that  $y = i$  .

On a new input  $x$ , to make a prediction, pick the class  $i$  that maximizes

$$\max_i h_{\theta}^{(i)}(x)$$

```
# Initialize logistic regression model for multiclass
logreg = LogisticRegression(multi_class='multinomial')
logreg.fit(X_train, y_train)
```

Adv. App. of AI and DT

27

## گل‌های زنبق مختلف

Iris Setosa



Iris Versicolour

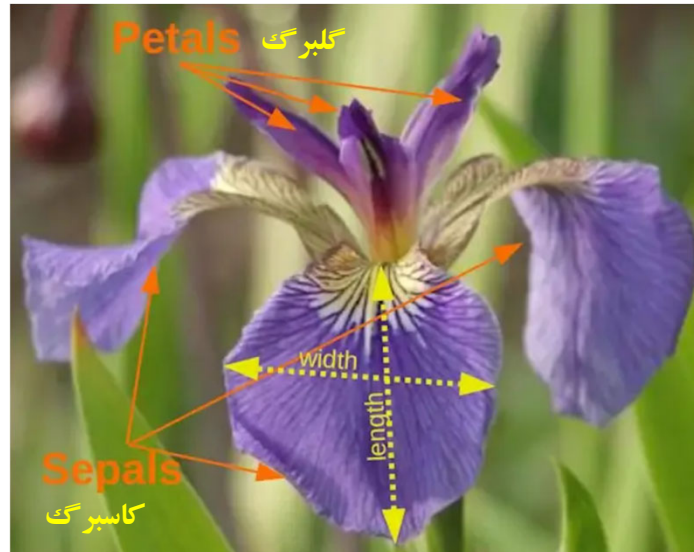


Iris Virginica



Adv. App. of AI and DT

## ابعاد کاسبرگ و گلبرگ زنبق مختلف



29

## ابعاد کاسبرگ و گلبرگ زنبق مختلف

فایل دیتای ایریس دارای طول و عرض کاسبرگ و گلبرگ است برچسب نام گل است که دارای سه نوع زنبق است

	sepal_length	sepal_width	petal_length	petal_width	species
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa

طول و عرض کاسبرگ و گلبرگ است چهار ویژگی گلها است ترسیم ویژگیها در چهار بعد میسر نیست و از روشهای دیگری باید بهره گرفت

Adv. App. of AI and DT

30

## طبقه بندی چند کلاسی

مثال: طبقه بندی چند کلاسی فایل دیتای ایریس

فایل: 8 multiclass regression iris petal.py

```
import seaborn as sns
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.datasets import load_iris
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, classification_report

# Load dataset
data = load_iris()
X = data.data[:, 2:4] # Select only petal length and petal width
y = data.target
# Map target numbers to iris species names
species_mapping = {0: 'Setosa', 1: 'Versicolor', 2: 'Virginica'}
species_names = [species_mapping[label] for label in y]
# Convert to DataFrame for easier plotting and visualization
df = pd.DataFrame(X, columns=['petal length (cm)', 'petal width (cm)'])
df['species'] = species_names # Use species names instead of numbers
# Split the data into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                    test_size=0.3, random_state=0)
```

Adv. App. of AI and DT

31

## طبقه بندی چند کلاسی

مثال: طبقه بندی چند کلاسی فایل دیتای ایریس

```
# Initialize and train logistic regression model for multiclass classification
logreg = LogisticRegression(multi_class='multinomial', solver='lbfgs',
                             max_iter=200)
logreg.fit(X_train, y_train)

# Predict and evaluate
y_pred = logreg.predict(X_test)
print("Accuracy:", accuracy_score(y_test, y_pred))
print(classification_report(y_test, y_pred,
                             target_names=species_mapping.values()))

# Convert test labels and predictions to species names for plotting
y_test_names = [species_mapping[label] for label in y_test]
y_pred_names = [species_mapping[label] for label in y_pred]
# Plotting classification results for petal length and petal width
plt.figure(figsize=(8, 6))
sns.scatterplot(x=X_test[:, 0], y=X_test[:, 1], hue=y_test_names,
                style=y_pred_names, palette='viridis', s=100, edgecolor="k")
plt.xlabel('Petal Length (cm)')
plt.ylabel('Petal Width (cm)')
plt.title('Logistic Regression Classification: Petal Length vs Petal Width')
plt.legend(title="True Species")
plt.show()
```

Adv. App. of AI and DT

32

## طبقه بندی چند کلاسی

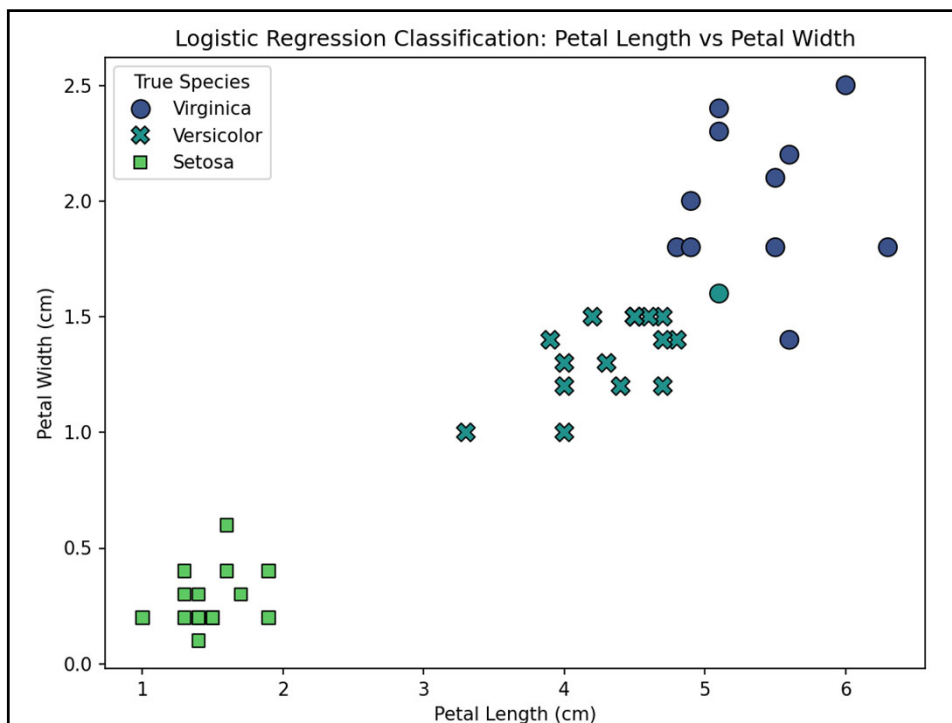
مثال: طبقه بندی چند کلاسی فایل دیتای ایریس

```
Accuracy: 0.9777777777777777
```

	precision	recall	f1-score	support
Setosa	1.00	1.00	1.00	16
Versicolor	1.00	0.94	0.97	18
Virginica	0.92	1.00	0.96	11
accuracy			0.98	45
macro avg	0.97	0.98	0.98	45
weighted avg	0.98	0.98	0.98	45

Adv. App. of AI and DT

33



## طبقه بندی چند کلاسی

مثال: طبقه بندی چند کلاسی فایل دیتای ایریس با چهار ویژگی طول و عرض گلبرگ و کاسبرگ

فایل: 8 multiclass regression iris pairplot.py

```
import seaborn as sns
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.datasets import load_iris
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, classification_report

# Load dataset
data = load_iris()
X = data.data
y = data.target
# Map target numbers to iris species names
species_mapping = {0: 'Setosa', 1: 'Versicolor', 2: 'Virginica'}
species_names = [species_mapping[label] for label in y]

# Convert to DataFrame for easier plotting with species names
df = pd.DataFrame(X, columns=data.feature_names)
df['species'] = species_names # Use species names instead of numbers
```

35

Adv. App. of AI and DT

## طبقه بندی چند کلاسی

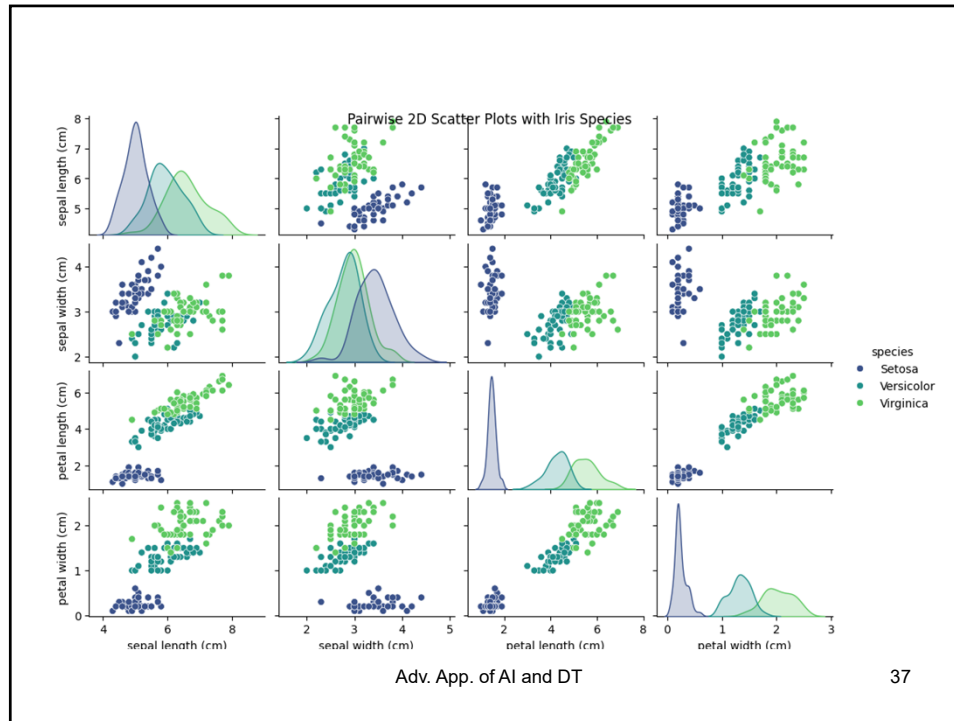
مثال: طبقه بندی چند کلاسی فایل دیتای ایریس با چهار ویژگی طول و عرض گلبرگ و کاسبرگ

```
# Split the data into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.3, random_state=0)
# Initialize and train logistic regression model for multiclass
classification
logreg = LogisticRegression(multi_class='multinomial',
solver='lbfgs', max_iter=200)
logreg.fit(X_train, y_train)
# Predict and evaluate
y_pred = logreg.predict(X_test)
print("Accuracy:", accuracy_score(y_test, y_pred))
print(classification_report(y_test, y_pred,
target_names=species_mapping.values()))

# Pairplot for 4D visualization using seaborn with species names
plt.figure(figsize=(8, 5))
sns.pairplot(df, hue='species', palette='viridis')
plt.suptitle("Pairwise 2D Scatter Plots with Iris Species", y=0.99)
plt.show()
```

36

Adv. App. of AI and DT



## تمرین برنامه نویسی

تمرین هشتم:

یک برنامه به زبان پایتون بنویسید که یک فایل داده را خوانده و رگرسیون لجستیک برای سه ویژگی و دو کلاس را حساب نماید.

۱- داده ها را طبقه بندی کنید

۲- طبقه بندی انجام شده را سه بعدی رسم کنید

۳- دقت طبقه بندی را بدست آورید

۲- برای یک حالت با ویژگیهای مشخص، کلاس را مشخص کنید

Adv. App. of AI and DT

38

