

KNN

→ dated 3-6-25

from sklearn.neighbors import

KNeighborsClassifier

knn = KNeighborsClassifier(n_neighbors=3)

knn.fit(x_train, y_train)

3) Experiment with diff. k values

Use cross-val score to evaluate

~~diff~~ scores (stability). for different 'k' values.

Visualize decision boundaries

X_2d = x.iloc[:, 1:3]

→ take the first 2 features for 2D plotting ~~gives~~

x_min, x_max = x_2d.iloc[:, 0].min() - 1,

x_2d.iloc[:, 0].max() + 1

y_min, y_max = x_2d.iloc[:, 1].min() - 1,

x_2d.iloc[:, 1].max() + 1

} getting min and max value from each column. k use ± 1 to ensure

`xx, yy = np.meshgrid(np.arange(x.min, x.max,`

`0.02),`

`np.arange(y.min, y.max,`
`0.02))`

Eg,

`x = [1, 2, 3]`

`y = [4, 5]`

`xx, yy = np.meshgrid(x, y)`

generates value from
 x_{\min} to x_{\max} with

step as 0.02

(i.e) if $x_{\min} = 0$ $x_{\max} = 1$

$\Rightarrow [0 \ 0.02 \ 0.04 \ 0.06 \ \dots \ 1]$

$\Rightarrow xx = \begin{bmatrix} [1, 2, 3] \\ [1, 2, 3] \end{bmatrix}$

$yy = \begin{bmatrix} [4, 4, 4] \\ [5, 5, 5] \end{bmatrix}$

2D Array

`mesh_points = pd.DataFrame(np.c_[xx.ravel(), yy.ravel()])`

Output

`[1, 2, 3, 1, 2, 3] \rightarrow xx`

`[4, 4, 4, 5, 5, 5] \rightarrow yy`

columns = $x_{2d} \cdot \text{columns}$

\rightarrow makes xx and yy
as 1D array

And `np.c` does this `(1, 4) (2, 4) (3, 4) ...`

And `pd.DataFrame` creates dataset with
 $x_{2d} \cdot \text{columns}$ names as column names.

Final Output is like see in Notebook

decision boundary, plot nicely, surrounds all data points

$I = \text{labels} = \text{knn.predict (mesh.predict)}$

D is nothing but

Sepall width

Petal length

1.00

0.00

1.02

0.00

Like this and it do our normal

knn thing (ie) finding Euclidean distance for each row of this mesh-points with ~~trained dataset of knn~~ the original dataset

Note:

In KNN model, it doesn't learn any pattern and then predict like normal models.

D simply stores the full dataset (training dataset) during modelling or training.

And later if we want to predict it with unseen data (test set), it simply computes

Euclidean distance, value of any unseen set

Eg: $X = (x_1, x_2, \dots, x_k)$

$Y = (y_1, y_2, \dots, y_k)$

→ each row value of train sets

$$\sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_k - y_k)^2}$$

After this we obtain lot of values (for each row), now we need to consider k nearest values (distance) and then take majority voting for the outcome (output)

Eg: Predict: (2.5, 3)

Ex

Point \Rightarrow	x	y	Label
1	1	2	A
2	2	3	A
3	3	3	B
6	6	5	B

in our example

Compute distances to all points

(1, 2) and (2.5, 3) $\Rightarrow \sqrt{(2.5-1)^2 + (3-2)^2} \approx 1.8$

for all other

(2, 3)	$\Rightarrow 0.5$
(3, 3)	$\Rightarrow \approx 0.5$
(6, 5)	$\Rightarrow \approx 4.0$

if $k=3$

(1, 2), (2, 3), (3, 3) is nearest, And labels for that are A, A and B

majority is A

So the answer is A.

Z-labels

I. $\text{class_to_int} = \{\text{label: idx for idx, label in enumerate(knn.class)}\}$
this creates a dictionary that has each class and with value as integer like 0, 1, 2... like that.

$Z = \text{np.array}([\text{class_to_int}[\text{label}] \text{ for label in } Z\text{-labels}])$

Z converts our Z -label's string output to integers as class-to-int.

$Z = Z.\text{reshape}(xx.\text{shape})$

we have to convert Z into 2D array because only if Z is 2D arr. then only we can plot Decision Boundary

For Eg:

we plot Decision Boundary by computing $x[i][j]$ and $y[i][j] \rightarrow$ so we need $Z[i][j]$
2D array

plt.contourf (xx, yy, Z, alpha=0.4) → diagram's transparency

① Defines x and y axis values &

x-axis \Rightarrow xx.min to xx.max

y-axis \Rightarrow yy.min to yy.max

We have to do this, because whenever,

$\Rightarrow Z[i][j] = 0$, it fills that patch of the plot with color for class 0.

$\Rightarrow Z[i][j] = 1$, it fills with color for class 1.

⋮
And so on

② Contourf looks at 2D array Z and draws regions where values are the same.

③ It automatically finds boundaries where Z changes (ie) class 0 ends and class 1 begins

This is why how we get our Decision boundary.