

Lecture 14

Classification and K -Nearest Neighbors

- 1 Recap
- 2 Case Study: Breast Tissue Classification
- 3 K -Nearest Neighbors Classification

1 Recap

2 Case Study: Breast Tissue Classification

3 *K*-Nearest Neighbors Classification

Recap

So far, we've been looking at **regression** problems,

Recap

So far, we've been looking at **regression** problems, where the label y we are trying to predict is quantitative (e.g., price of wine).

Recap

So far, we've been looking at **regression** problems, where the label y we are trying to predict is quantitative (e.g., price of wine).

Today, we switch to **classification** problems, where the label y is categorical.

Recap

So far, we've been looking at **regression** problems, where the label y we are trying to predict is quantitative (e.g., price of wine).

Today, we switch to **classification** problems, where the label y is categorical.

We will focus on the differences between regression and classification.

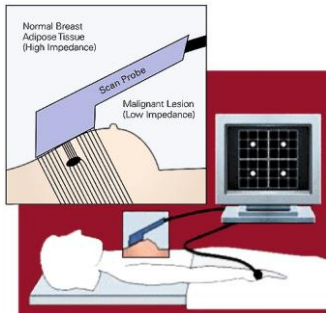
- 1 Recap
- 2 Case Study: Breast Tissue Classification
- 3 K -Nearest Neighbors Classification

Breast Tissue Classification

Electrical signals can be used to detect whether tissue is cancerous.

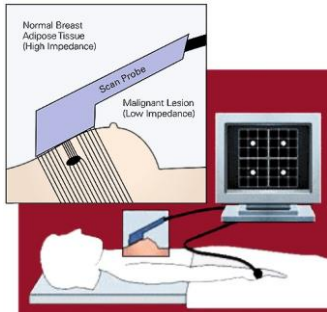
Breast Tissue Classification

Electrical signals can be used to detect whether tissue is cancerous.



Breast Tissue Classification

Electrical signals can be used to detect whether tissue is cancerous.



The goal is to determine whether a sample of breast tissue is:

- | | | | |
|----------------------|------------|------------------|---|
| 1. connective tissue | Bağ dokusu | 4. carcinoma | Karsinom (kanserli tümör türü) |
| 2. adipose tissue | Yağ dokusu | 5. fibro-adenoma | Fibroadenom (iyi huylu meme tümörü) |
| 3. glandular tissue | Bez dokusu | 6. mastopathy | Mastopati (meme dokusu hastalığı / iyi huylu değişiklikler) |

Reading in the Data

```
import pandas as pd
df = pd.read_csv("https://datasci112.stanford.edu/data/BreastTissue.csv")
df
```

Reading in the Data

```
import pandas as pd
df = pd.read_csv("https://datasci112.stanford.edu/data/BreastTissue.csv")
df
```

	Case #	Class	I0	PA500	HFS	DA	Area	A/DA	Max IP	DR	P
0	1	car	524.794072	0.187448	0.032114	228.800228	6843.598481	29.910803	60.204880	220.737212	556.828334
1	2	car	330.000000	0.226893	0.265290	121.154201	3163.239472	26.109202	69.717361	99.084964	400.225776
2	3	car	551.879287	0.232478	0.063530	264.804935	11888.391827	44.894903	77.793297	253.785300	656.769449
3	4	car	380.000000	0.240855	0.286234	137.640111	5402.171180	39.248524	88.758446	105.198568	493.701814
4	5	car	362.831266	0.200713	0.244346	124.912559	3290.462446	26.342127	69.389389	103.866552	424.796503
...
101	102	adi	2000.000000	0.106989	0.105418	520.222649	40087.920984	77.059161	204.090347	478.517223	2088.648870
102	103	adi	2600.000000	0.200538	0.208043	1063.441427	174480.476218	164.071543	418.687286	977.552367	2664.583623
103	104	adi	1600.000000	0.071908	-0.066323	436.943603	12655.342135	28.963331	103.732704	432.129749	1475.371534
104	105	adi	2300.000000	0.045029	0.136834	185.446044	5086.292497	27.427344	178.691742	49.593290	2480.592151
105	106	adi	2600.000000	0.069988	0.048869	745.474369	39845.773698	53.450226	154.122604	729.368395	2545.419744

106 rows x 11 columns

I0 – Sıfır frekansta empedans (ohm)
PA500 – 500 KHz'de faz açısı
HFS – Yüksek frekans eğimi
DA – Spektrum sonları arasındaki mesafe
AREA – Spektrum altındaki alan

A/DA – DA'ya normalizasyonlu alan
MAX IP – Spektrumun maksimum noktası
DR – I0 ile maksimum frekans noktası arasındaki mesafe
P – Spektral eğrinin uzunluğu
Class – Doku sınıfı etiketi (örneğin: car, fad, mas, gla, con, adi)

Reading in the Data

```
import pandas as pd
df = pd.read_csv("https://datasci112.stanford.edu/data/BreastTissue.csv")
df
```

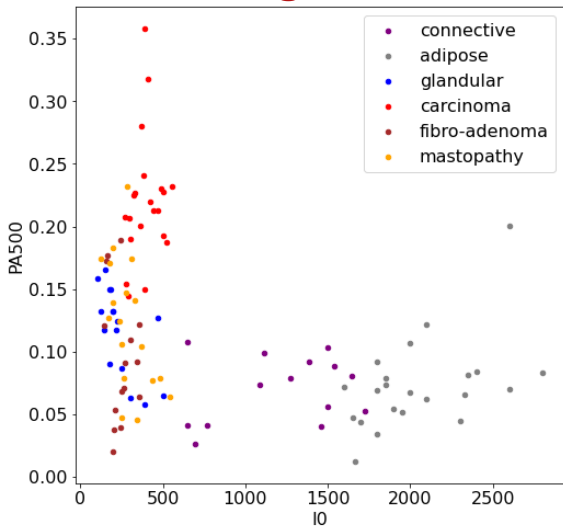
	Case #	Class	IO	PA500	HFS	DA	Area	A/DA	Max IP	DR	P
0	1	car	524.794072	0.187448	0.032114	228.800228	6843.598481	29.910803	60.204880	220.737212	556.828334
1	2	car	330.000000	0.226893	0.265290	121.154201	3163.239472	26.109202	69.717361	99.084964	400.225776
2	3	car	551.879287	0.232478	0.063530	264.804935	11888.391827	44.894903	77.793297	253.785300	656.769449
3	4	car	380.000000	0.240855	0.286234	137.640111	5402.171180	39.248524	88.758446	105.198568	493.701814
4	5	car	362.831266	0.200713	0.244346	124.912559	3290.462446	26.342127	69.389389	103.866552	424.796503
...
101	102	adi	2000.000000	0.106989	0.105418	520.222649	40087.920984	77.059161	204.090347	478.517223	2088.648870
102	103	adi	2600.000000	0.200538	0.208043	1063.441427	174480.476218	164.071543	418.687286	977.552367	2664.583623
103	104	adi	1600.000000	0.071908	-0.066323	436.943603	12655.342135	28.963331	103.732704	432.129749	1475.371534
104	105	adi	2300.000000	0.045029	0.136834	185.446044	5086.292497	27.427344	178.691742	49.593290	2480.592151
105	106	adi	2600.000000	0.069988	0.048869	745.474369	39845.773698	53.450226	154.122604	729.368395	2545.419744

106 rows x 11 columns

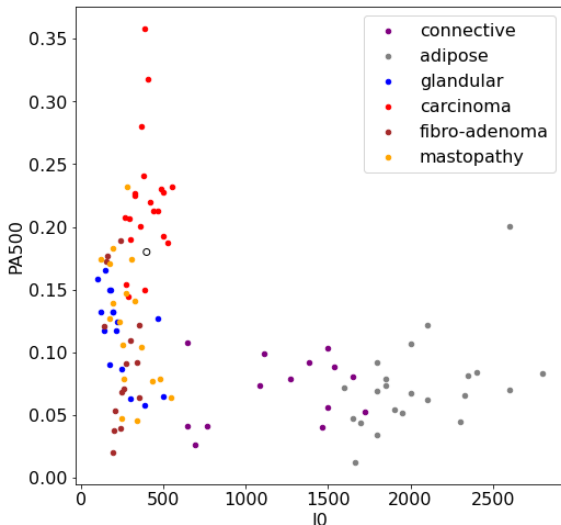
We will focus on two features:

- I_0 : impedivity at 0 kHz, Dokunun düşük frekanstaki elektrik direnci.
- PA_{500} : phase angle at 500 kHz. dokunun akım ve gerilim arasındaki açı farkı.

Visualizing the Data

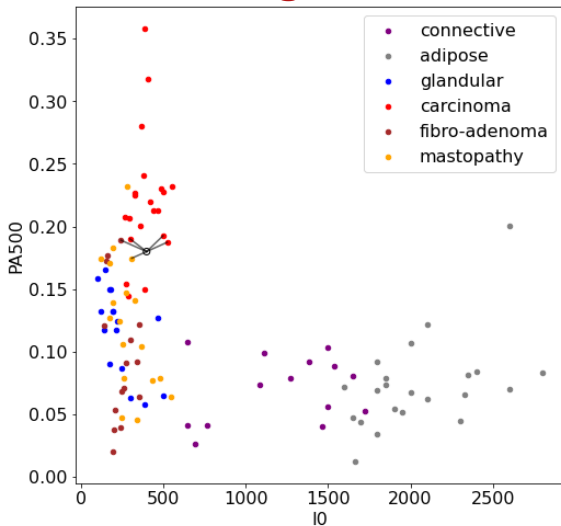


Visualizing the Data

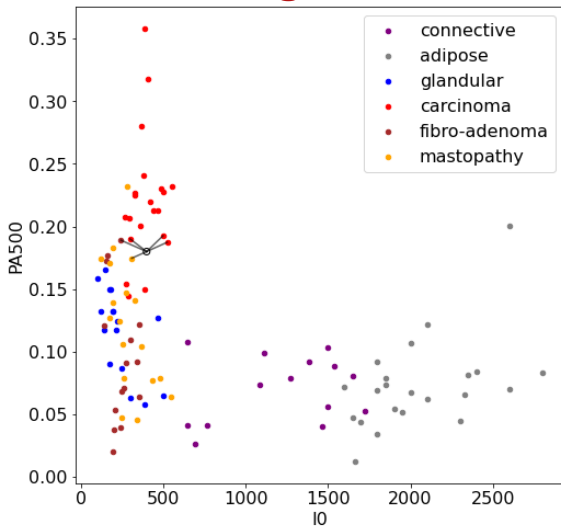


Consider a new sample, with an I_0 of 400 and a PA_{500} of 0.18.
What kind of tissue is it?

Visualizing the Data

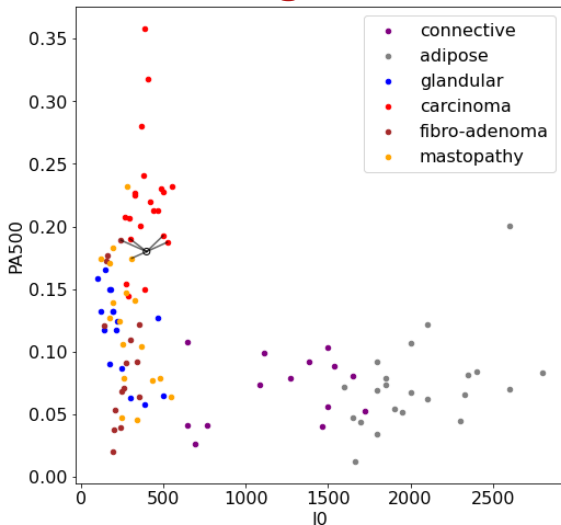


Visualizing the Data



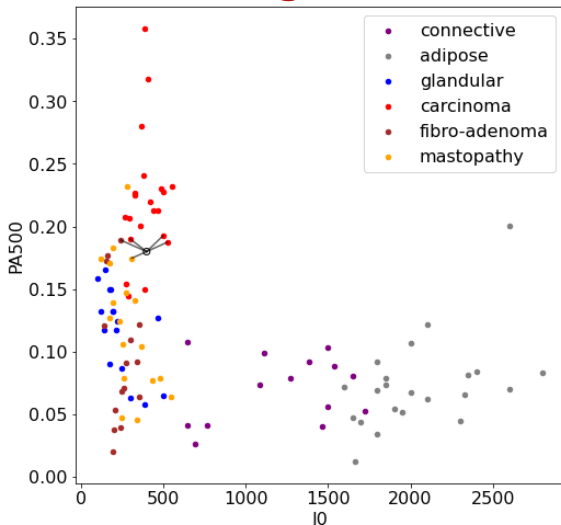
Of its 5 nearest neighbors in the training data:

Visualizing the Data



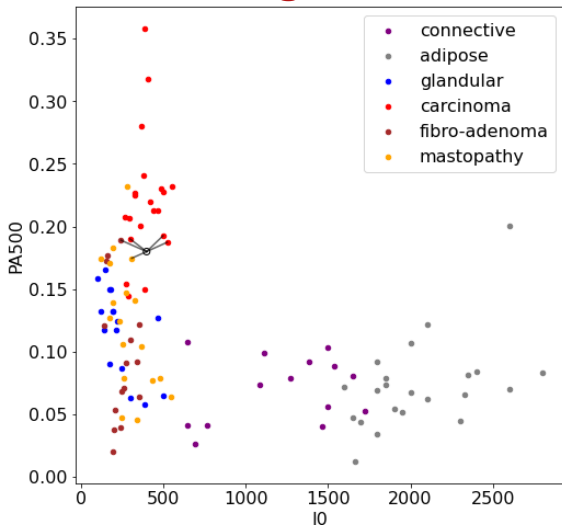
Of its 5 nearest neighbors in the training data:
3 are carcinomas,

Visualizing the Data



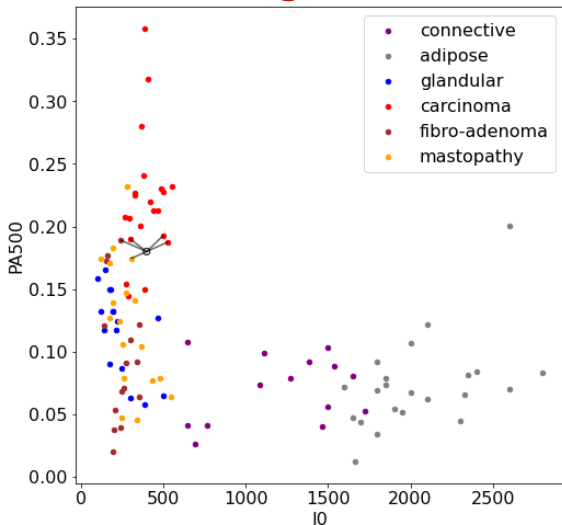
Of its 5 nearest neighbors in the training data:
3 are carcinomas, 1 is fibro-adenoma,

Visualizing the Data



Of its 5 nearest neighbors in the training data:
3 are carcinomas, 1 is fibro-adenoma, 1 is mastopathy,

Visualizing the Data



Of its 5 nearest neighbors in the training data:
3 are carcinomas, 1 is fibro-adenoma, 1 is mastopathy,
so our best guess is that it is a carcinoma.

- 1 Recap
- 2 Case Study: Breast Tissue Classification
- 3 *K*-Nearest Neighbors Classification

```
X_train = df[["I0", "PA500"]]  
y_train = df["Class"]  
  
x_test = pd.Series({"I0": 400, "PA500": .18})  
X_test = x_test.to_frame().T  
X_test
```



```
X_train = df[["I0", "PA500"]]  
y_train = df["Class"]  
  
x_test = pd.Series({"I0": 400, "PA500": .18})  
X_test = x_test.to_frame().T  
X_test
```

	I0	PA500
0	400.0	0.18

```
X_train = df[["I0", "PA500"]]  
y_train = df["Class"]  
  
x_test = pd.Series({"I0": 400, "PA500": .18})  
X_test = x_test.to_frame().T  
X_test
```

	I0	PA500
0	400.0	0.18

Here is code we used for k -nearest neighbor regression.

```
from sklearn.preprocessing import StandardScaler  
from sklearn.neighbors import KNeighborsRegressor  
from sklearn.pipeline import make_pipeline  
  
pipeline = make_pipeline(  
    StandardScaler(),  
    KNeighborsRegressor(n_neighbors=5, metric="euclidean"))  
pipeline.fit(X_train, y_train)  
pipeline.predict(X_test)
```

```
X_train = df[["I0", "PA500"]]  
y_train = df["Class"]  
  
x_test = pd.Series({"I0": 400, "PA500": .18})  
X_test = x_test.to_frame().T  
X_test
```

	I0	PA500
0	400.0	0.18

Here is code we used for k -nearest neighbor regression.

```
from sklearn.preprocessing import StandardScaler  
from sklearn.neighbors import KNeighborsRegressor  
from sklearn.pipeline import make_pipeline  
  
pipeline = make_pipeline(  
    StandardScaler(),  
    KNeighborsRegressor(n_neighbors=5, metric="euclidean"))  
pipeline.fit(X_train, y_train)  
pipeline.predict(X_test)
```

What would need to change for classification?

```
X_train = df[["I0", "PA500"]]
y_train = df["Class"]

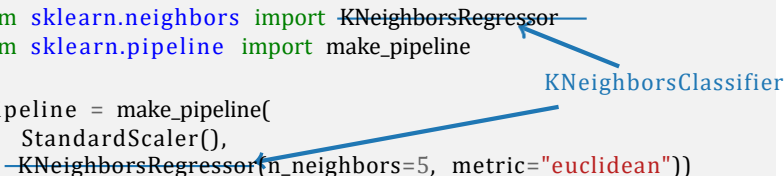
x_test = pd.Series({"I0": 400, "PA500": .18})
X_test = x_test.to_frame().T
X_test
```

	I0	PA500
0	400.0	0.18

Here is code we used for k -nearest neighbor regression.

```
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsRegressor
from sklearn.pipeline import make_pipeline

pipeline = make_pipeline(
    StandardScaler(),
    KNeighborsRegressor(n_neighbors=5, metric="euclidean"))
pipeline.fit(X_train, y_train)
pipeline.predict(X_test)
```



What would need to change for classification?

```
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.pipeline import make_pipeline

pipeline = make_pipeline(
    StandardScaler(),
    KNeighborsClassifier(n_neighbors=5, metric="euclidean"))

pipeline.fit(X_train, y_train)
pipeline.predict(X_test)
```

```
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.pipeline import make_pipeline

pipeline = make_pipeline(
    StandardScaler(),
    KNeighborsClassifier(n_neighbors=5, metric="euclidean"))

pipeline.fit(X_train, y_train)
pipeline.predict(X_test)

array(['car'], dtype=object)
```

```
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.pipeline import make_pipeline

pipeline = make_pipeline(
    StandardScaler(),
    KNeighborsClassifier(n_neighbors=5, metric="euclidean"))

pipeline.fit(X_train, y_train)
pipeline.predict(X_test)

array(['car'], dtype=object)
```

Instead of returning a single predicted class, we can ask it to return the predicted probabilities for each class.

```
pipeline.predict_proba(X_test)
```

```
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.pipeline import make_pipeline

pipeline = make_pipeline(
    StandardScaler(),
    KNeighborsClassifier(n_neighbors=5, metric="euclidean"))

pipeline.fit(X_train, y_train)
pipeline.predict(X_test)

array(['car'], dtype=object)
```

Instead of returning a single predicted class, we can ask it to return the predicted probabilities for each class.

```
pipeline.predict_proba(X_test)

array([[0. , 0.6, 0. , 0.2, 0. , 0.2]])
```



```
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.pipeline import make_pipeline

pipeline = make_pipeline(
    StandardScaler(),
    KNeighborsClassifier(n_neighbors=5, metric="euclidean"))

pipeline.fit(X_train, y_train)
pipeline.predict(X_test)

array(['car'], dtype=object)
```

Instead of returning a single predicted class, we can ask it to return the predicted probabilities for each class.

```
pipeline.predict_proba(X_test)

array([[0. , 0.6, 0. , 0.2, 0. , 0.2]])

pipeline.classes_
```

```
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.pipeline import make_pipeline

pipeline = make_pipeline(
    StandardScaler(),
    KNeighborsClassifier(n_neighbors=5, metric="euclidean"))

pipeline.fit(X_train, y_train)
pipeline.predict(X_test)
```

```
array(['car'], dtype=object)
```

Instead of returning a single predicted class, we can ask it to return the predicted probabilities for each class.

```
pipeline.predict_proba(X_test)
array([[0. , 0.6, 0. , 0.2, 0. , 0.2]])
```

```
pipeline.classes_
array(['adi', 'car', 'con', 'fad', 'gla', 'mas'], dtype=object)
```

```
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.pipeline import make_pipeline

pipeline = make_pipeline(
    StandardScaler(),
    KNeighborsClassifier(n_neighbors=5, metric="euclidean"))

pipeline.fit(X_train, y_train)
pipeline.predict(X_test)
```

```
array(['car'], dtype=object)
```

Instead of returning a single predicted class, we can ask it to return the predicted probabilities for each class.

```
pipeline.predict_proba(X_test)
array([[0. , 0.6, 0. , 0.2, 0. , 0.2]])

pipeline.classes_
```

```
array(['adi', 'car', 'con', 'fad', 'gla', 'mas'], dtype=object)
```

How did Scikit-Learn calculate these predicted probabilities?

Cross-Validation for Classification

Cross-Validation for Classification

Here is code we used to cross-validate a regression model.

```
from sklearn.model_selection import cross_val_score

cross_val_score(
    pipeline, X_train, y_train,
    scoring="neg_mean_squared_error",
    cv=10)
```

Cross-Validation for Classification

Here is code we used to cross-validate a regression model.

```
from sklearn.model_selection import cross_val_score

cross_val_score(
    pipeline, X_train, y_train,
    scoring="neg_mean_squared_error",
    cv=10)
```

What would need to change for classification?

Cross-Validation for Classification

Here is code we used to cross-validate a regression model.

```
from sklearn.model_selection import cross_val_score

cross_val_score(
    pipeline, X_train, y_train,
    scoring="neg_mean_squared_error",
    cv=10)
```

What would need to change for classification?

Cross-Validation for Classification

Here is code we used to cross-validate a regression model.

```
from sklearn.model_selection import cross_val_score

cross_val_score(
    pipeline, X_train, y_train,
    scoring="neg_mean_squared_error",
    cv=10)
```

What would need to change for classification?

We need a different scoring method for classification. A simple one is **accuracy**:

$$\text{accuracy} = \frac{\text{\# correct predictions}}{\text{\# predictions}}.$$

Cross-Validation for Classification

```
from sklearn.model_selection import cross_val_score

scores = cross_val_score(
    pipeline, X_train, y_train,
    scoring="accuracy",
    cv=10)
scores
```

Cross-Validation for Classification

```
from sklearn.model_selection import cross_val_score

scores = cross_val_score(
    pipeline, X_train, y_train,
    scoring="accuracy",
    cv=10)
scores

array([0.63636364, 0.81818182, 0.45454545, 0.54545455, 0.63636364,
       0.54545455, 0.5         , 0.6         , 0.4         , 0.7         ])
```

Cross-Validation for Classification

```
from sklearn.model_selection import cross_val_score

scores = cross_val_score(
    pipeline, X_train, y_train,
    scoring="accuracy",
    cv=10)
scores
array([0.63636364, 0.81818182, 0.45454545, 0.54545455, 0.63636364,
       0.54545455, 0.5         , 0.6         , 0.4         , 0.7         ])
```

As before, we can get an overall estimate of test accuracy by averaging the cross-validation accuracies:

Cross-Validation for Classification

```
from sklearn.model_selection import cross_val_score

scores = cross_val_score(
    pipeline, X_train, y_train,
    scoring="accuracy",
    cv=10)
scores
array([0.63636364, 0.81818182, 0.45454545, 0.54545455, 0.63636364,
       0.54545455, 0.5         , 0.6         , 0.4         , 0.7         ])
```

As before, we can get an overall estimate of test accuracy by averaging the cross-validation accuracies:

```
scores.mean()
```

Cross-Validation for Classification

```
from sklearn.model_selection import cross_val_score
```

```
scores = cross_val_score(  
    pipeline, X_train, y_train,  
    scoring="accuracy",  
    cv=10)
```

```
scores
```

```
array([0.63636364, 0.81818182, 0.45454545, 0.54545455, 0.63636364,  
       0.54545455, 0.5         , 0.6         , 0.4         , 0.7         ])
```

As before, we can get an overall estimate of test accuracy by averaging the cross-validation accuracies:

```
scores.mean()
```

```
0.5836363636363637
```

Cross-Validation for Classification

```
from sklearn.model_selection import cross_val_score

scores = cross_val_score(
    pipeline, X_train, y_train,
    scoring="accuracy",
    cv=10)
scores
array([0.63636364, 0.81818182, 0.45454545, 0.54545455, 0.63636364,
       0.54545455, 0.5         , 0.6         , 0.4         , 0.7         ])
```

As before, we can get an overall estimate of test accuracy by averaging the cross-validation accuracies:

```
scores.mean()

0.5836363636363637
```

Accuracy is not always the best measure of a classification model.

Cross-Validation for Classification

```
from sklearn.model_selection import cross_val_score

scores = cross_val_score(
    pipeline, X_train, y_train,
    scoring="accuracy",
    cv=10)
scores
array([0.63636364, 0.81818182, 0.45454545, 0.54545455, 0.63636364,
       0.54545455, 0.5         , 0.6         , 0.4         , 0.7         ])
```

As before, we can get an overall estimate of test accuracy by averaging the cross-validation accuracies:

```
scores.mean()

0.5836363636363637
```

Accuracy is not always the best measure of a classification model. We'll talk about some better measures next time.