

⑤アルゴリズム_実装演習

In [1]:

```
%matplotlib inline
import numpy as np
import matplotlib.pyplot as plt
from scipy import stats
```

=====

1.k近傍法(knn)

- ・訓練データ生成

In [2]:

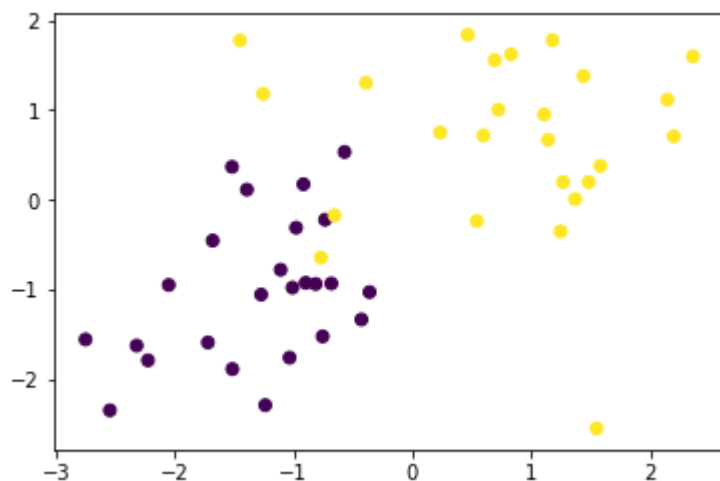
```
def gen_data():
    x0 = np.random.normal(size=50).reshape(-1, 2) - 1
    x1 = np.random.normal(size=50).reshape(-1, 2) + 1.
    x_train = np.concatenate([x0, x1])
    y_train = np.concatenate([np.zeros(25), np.ones(25)]).astype(np.int)
    return x_train, y_train
```

In [3]:

```
X_train, ys_train = gen_data()
plt.scatter(X_train[:, 0], X_train[:, 1], c=ys_train)
```

Out[3]:

<matplotlib.collections.PathCollection at 0x26648d0dd48>



- ・予測 →予測するデータ点との、距離が最も近い k 個の、訓練データのラベルの最頻値を割り当てる(教師あり学習)

In [4]:

```
def distance(x1, x2):
    return np.sum((x1 - x2)**2, axis=1)

def knn_predict(n_neighbors, x_train, y_train, X_test):
    y_pred = np.empty(len(X_test), dtype=y_train.dtype)
    for i, x in enumerate(X_test):
        distances = distance(x, X_train)
        nearest_index = distances.argsort()[:n_neighbors]
        mode, _ = stats.mode(y_train[nearest_index])
        y_pred[i] = mode
    return y_pred

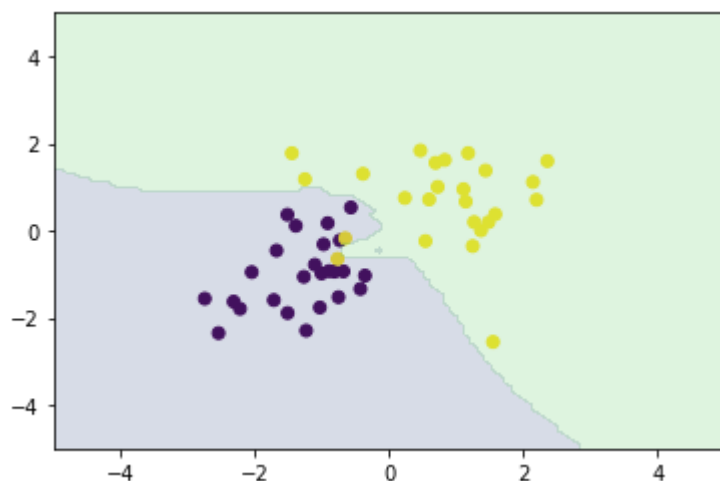
def plot_result(x_train, y_train, y_pred):
    xx0, xx1 = np.meshgrid(np.linspace(-5, 5, 100), np.linspace(-5, 5, 100))
    xx = np.array([xx0, xx1]).reshape(2, -1).T
    plt.scatter(x_train[:, 0], x_train[:, 1], c=y_train)
    plt.contourf(xx0, xx1, y_pred.reshape(100, 100).astype(dtype=np.float), alpha=0.2, levels=np.linspace(0, 1, 3))
```

In [5]:

```
n_neighbors = 3

xx0, xx1 = np.meshgrid(np.linspace(-5, 5, 100), np.linspace(-5, 5, 100))
X_test = np.array([xx0, xx1]).reshape(2, -1).T

y_pred = knn_predict(n_neighbors, X_train, y_train, X_test)
plot_result(X_train, y_train, y_pred)
```



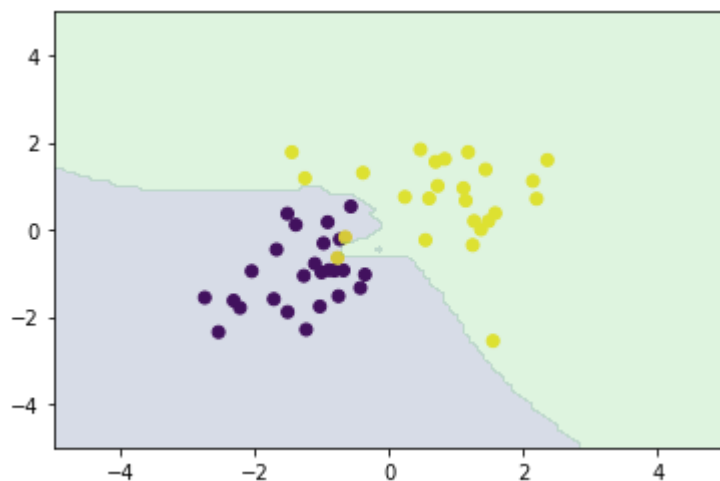
• scikit-learnによる実行

In [6]:

```
xx0, xx1 = np.meshgrid(np.linspace(-5, 5, 100), np.linspace(-5, 5, 100))
xx = np.array([xx0, xx1]).reshape(2, -1).T
```

In [7]:

```
from sklearn.neighbors import KNeighborsClassifier
knc = KNeighborsClassifier(n_neighbors=n_neighbors).fit(X_train, ys_train)
plt_resut(X_train, ys_train, knc.predict(xx))
```



=====

2.k平均クラスタリング(k-means)

- ・教師なし学習

In [9]:

```
%matplotlib inline
import numpy as np
import matplotlib.pyplot as plt
```

- ・データ生成

In [10]:

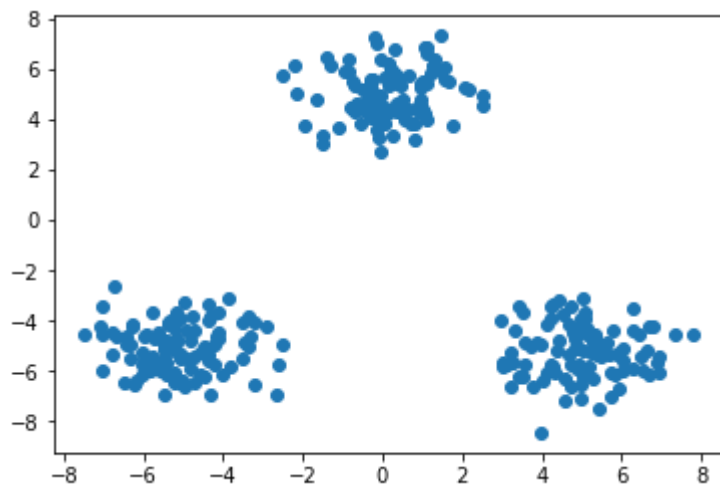
```
def gen_data():
    x1 = np.random.normal(size=(100, 2)) + np.array([-5, -5])
    x2 = np.random.normal(size=(100, 2)) + np.array([5, -5])
    x3 = np.random.normal(size=(100, 2)) + np.array([0, 5])
    return np.vstack((x1, x2, x3))
```

In [11]:

```
#データ作成
X_train = gen_data()
#データ描画
plt.scatter(X_train[:, 0], X_train[:, 1])
```

Out[11]:

<matplotlib.collections.PathCollection at 0x266493fa6c8>



・学習

k-meansアルゴリズムは以下のとおりである

- 1) 各クラスタ中心の初期値を設定する
- 2) 各データ点に対して、各クラスタ中心との距離を計算し、最も距離が近いクラスタを割り当てる
- 3) 各クラスタの平均ベクトル（中心）を計算する
- 4) 収束するまで2, 3の処理を繰り返す

In [13]:

```
def distance(x1, x2):
    return np.sum((x1 - x2)**2, axis=1)

n_clusters = 3
iter_max = 100

# 各クラスタ中心をランダムに初期化
centers = X_train[np.random.choice(len(X_train), n_clusters, replace=False)]

for _ in range(iter_max):
    prev_centers = np.copy(centers)
    D = np.zeros((len(X_train), n_clusters))
    # 各データ点に対して、各クラスタ中心との距離を計算
    for i, x in enumerate(X_train):
        D[i] = distance(x, centers)
    # 各データ点に、最も距離が近いクラスタを割り当
    cluster_index = np.argmin(D, axis=1)
    # 各クラスタの中心を計算
    for k in range(n_clusters):
        index_k = cluster_index == k
        centers[k] = np.mean(X_train[index_k], axis=0)
    # 収束判定
    if np.allclose(prev_centers, centers):
        break
```

・クラスタリング結果

In [14]:

```
def plt_result(X_train, centers, xx):
    # データを可視化
    plt.scatter(X_train[:, 0], X_train[:, 1], c=y_pred, cmap='spring')
    # 中心を可視化
    plt.scatter(centers[:, 0], centers[:, 1], s=200, marker='X', lw=2, c='black', edgecolor="white")
    # 領域の可視化
    pred = np.empty(len(xx), dtype=int)
    for i, x in enumerate(xx):
        d = distance(x, centers)
        pred[i] = np.argmin(d)
    plt.contourf(xx0, xx1, pred.reshape(100, 100), alpha=0.2, cmap='spring')
```

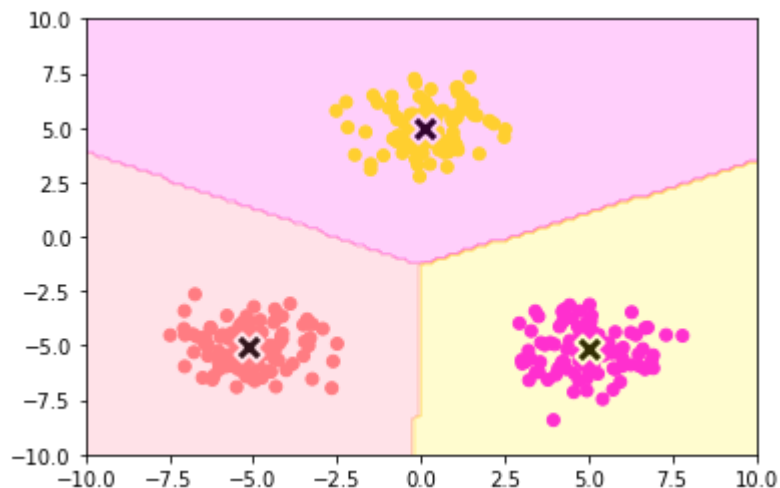
In [15]:

```
y_pred = np.empty(len(X_train), dtype=int)
for i, x in enumerate(X_train):
    d = distance(x, centers)
    y_pred[i] = np.argmin(d)
```

```
array([[ 0.13926641,  5.0093696 ],
       [-5.11387211, -5.06088169],
       [ 4.98960769, -5.18133892]])
```

In [19]:

```
plt_result(X_train, kmeans.cluster_centers_, xx)
```



=====

・k近傍法とk平均クラスタリングを実施。(各自、前半numpy・後半scikit-learnによる実行)・k近傍法は、予測するデータ点との、距離が最も近い k 個の、訓練データのラベルの最頻値を割り当てる方法。教師あり学習である。・k平均クラスタリングは、クラスタの平均(重心)を用い、与えられたクラスタ数 k 個に分類する方法。教師なし学習である。

In []: