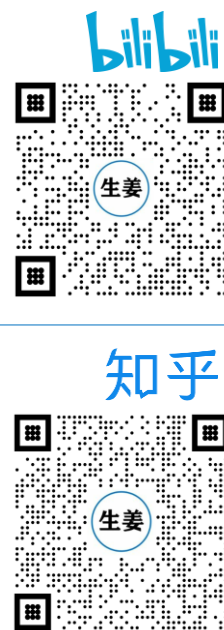
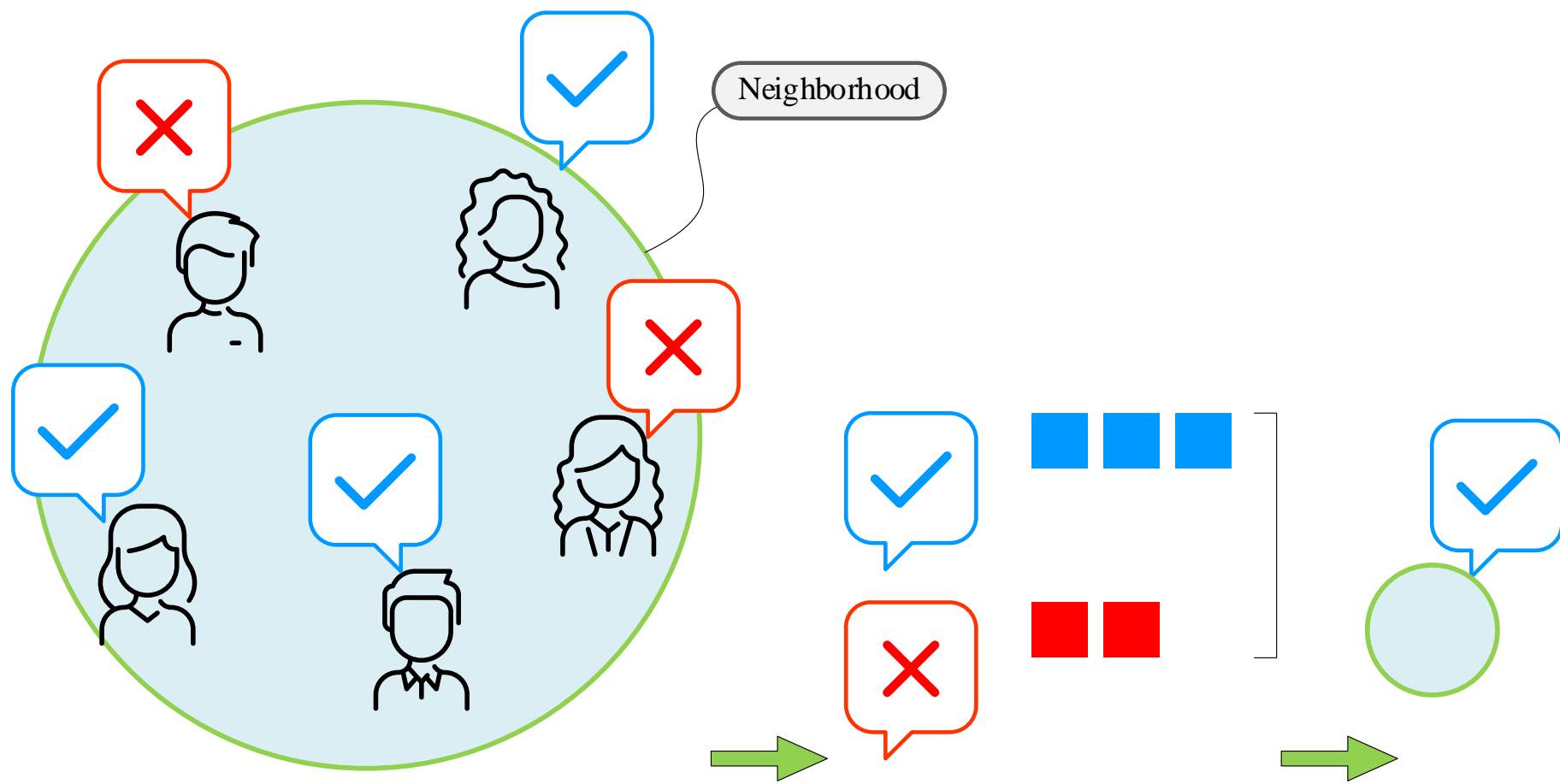
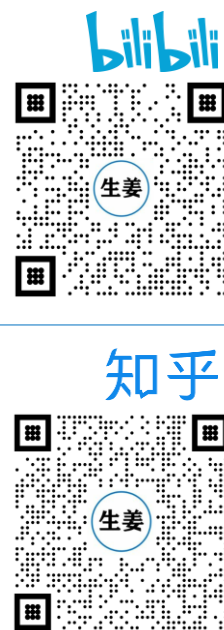
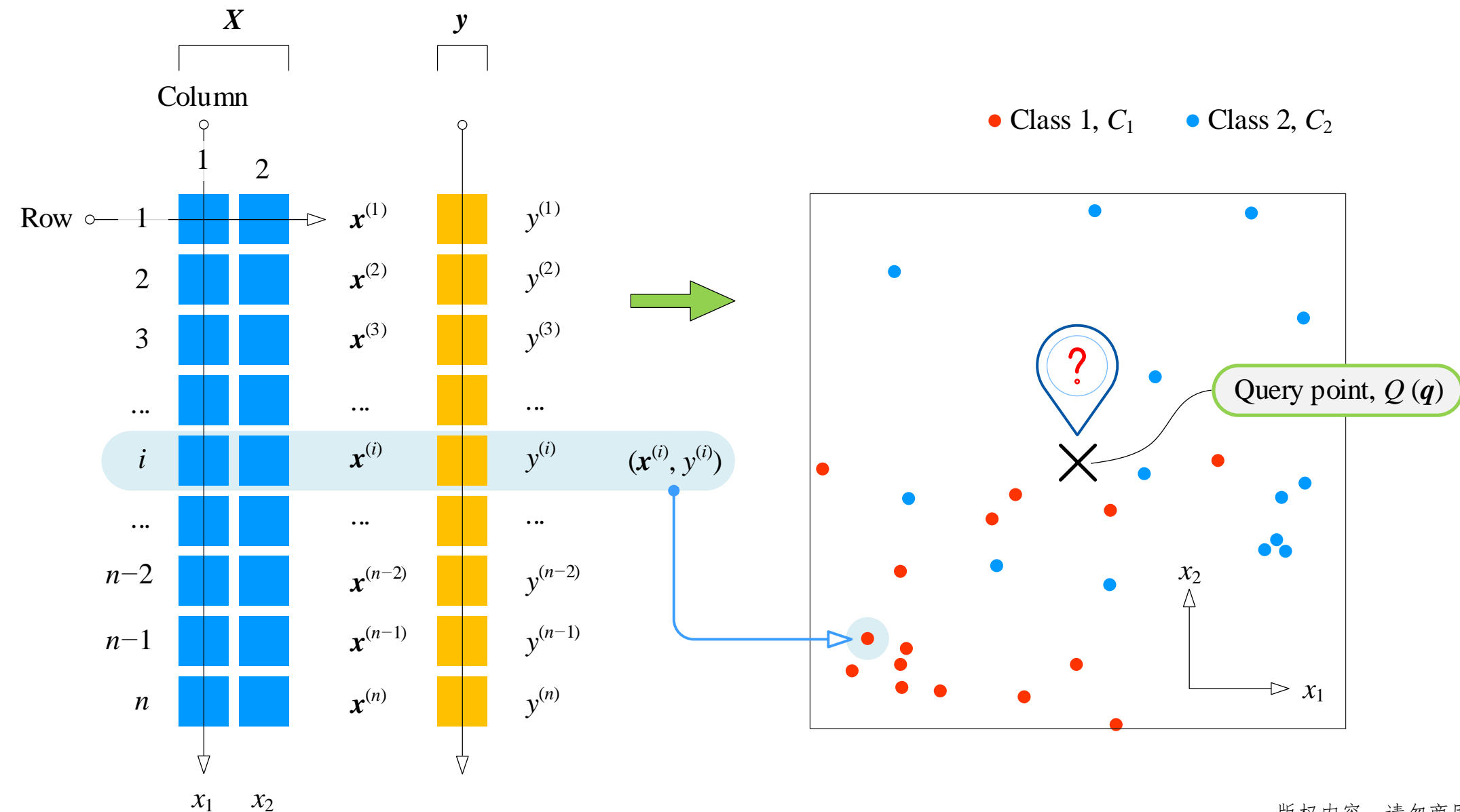


k临近分类kNN——小范围投票，少数服从多数

1





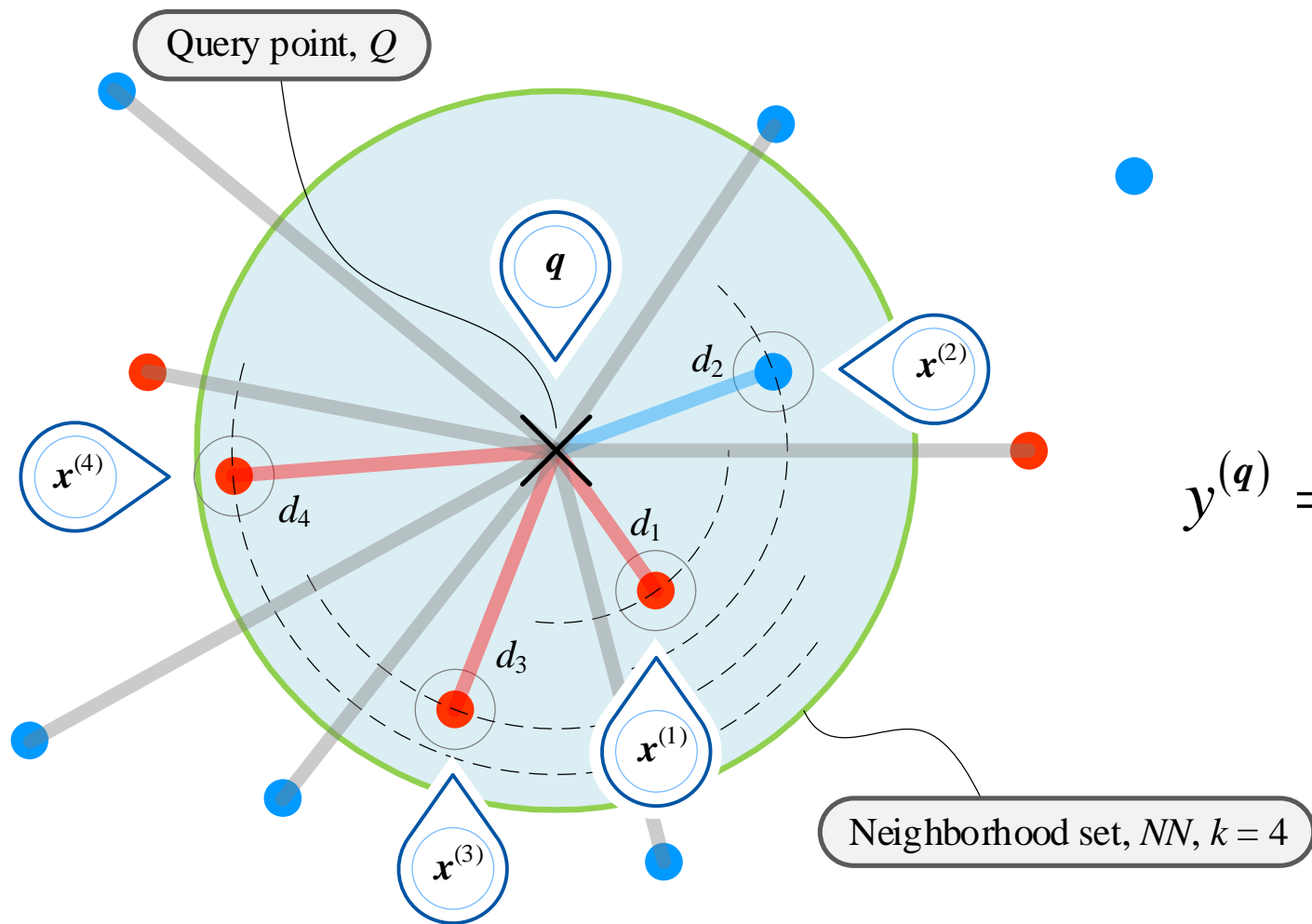
$$y^{(q)} = \arg \max_{C_l} \sum_{i \in kNN(q)} I(y^{(i)} = C_l)$$

$$K = 2$$

$$y^{(q)} = \max_{C_1, C_2} \left\{ \sum_{i \in kNN(q)} I(y^{(i)} = C_1), \sum_{i \in kNN(q)} I(y^{(i)} = C_2) \right\}$$



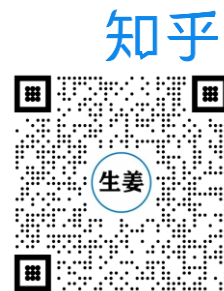
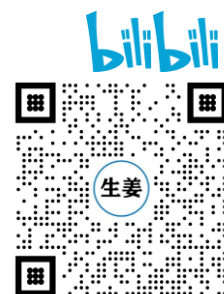
● Class 1, C_1 ● Class 2, C_2



$$\sum_{i \in kNN(q)} I(y^{(i)} = C_1) = 3_{(C_1)}$$

$$\sum_{i \in kNN(q)} I(y^{(i)} = C_2) = 1_{(C_2)}$$

$$y^{(q)} = \max_{C_1, C_2} \{3_{(C_1)}, 1_{(C_2)}\} = C_1$$

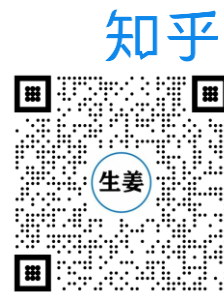
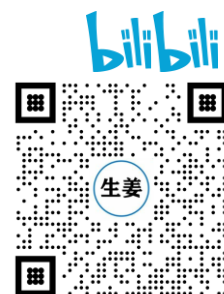


- `sklearn.neighbors.KNeighborsClassifier`
- 函数默认的近邻数量`n_neighbors`为5，默认距离度量`metric`为欧氏距离 (Euclidean distance)。
- 这个函数常用的methods为`fit(X, y)` 和 `predict(q)`;
- `fit(X, y)`用来加载样本数据， `predict(q)` 用来预测查询点 q 的分类

Scikit-learn (曾叫做scikits.learn还叫做sklearn) 是用于Python编程语言的自由软件机器学习库。它的特征是具有各种分类、回归和聚类算法



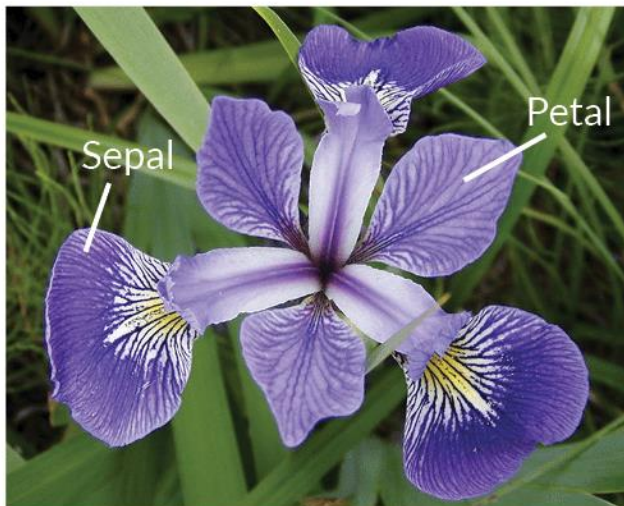
<https://en.wikipedia.org/wiki/Scikit-learn>



数据集包含了150个样本，都属于鸢尾属下的三个亚属，分别是山鸢尾、变色鸢尾和维吉尼亚鸢尾。四个特征被用作样本的定量分析，它们分别是花萼和花瓣的长度和宽度。



Iris Setosa



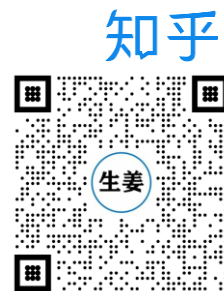
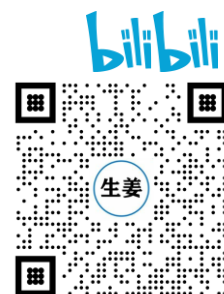
Iris Versicolor

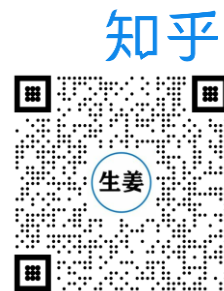
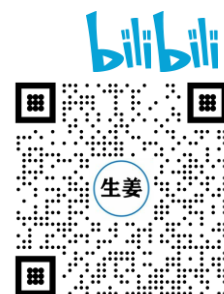
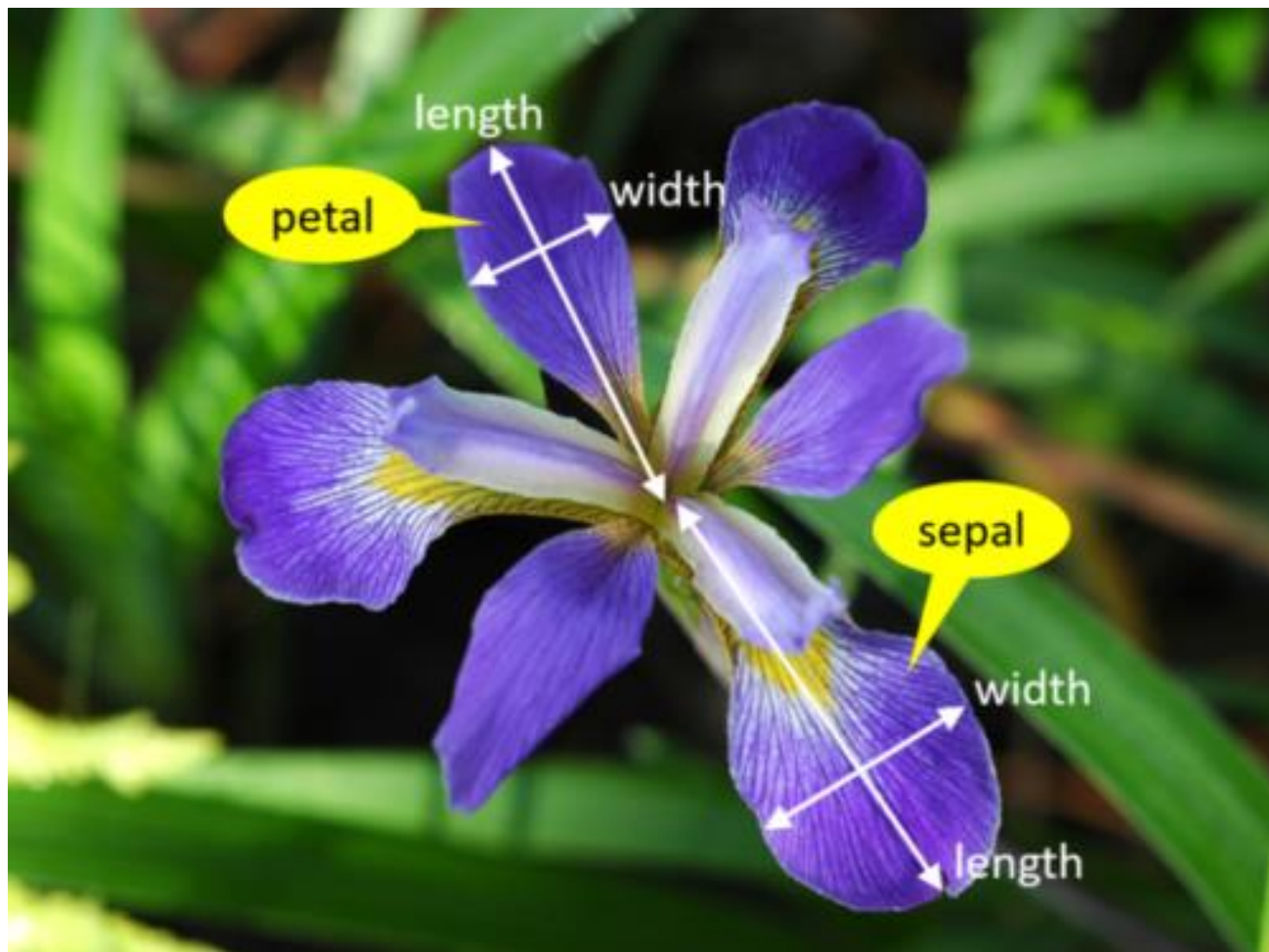


Iris Virginica

● Setosa, C_1 ● Versicolor, C_2 ● Virginica, C_3

https://en.wikipedia.org/wiki/Iris_flower_data_set





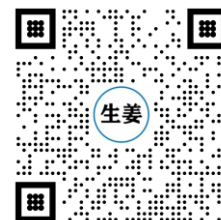
度量单位为厘米:

费雪鸢尾花卉数据集

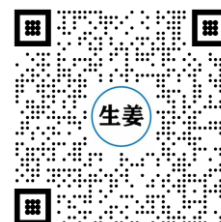
花萼长度 ◆	花萼宽度 ◆	花瓣长度 ◆	花瓣宽度 ◆	属种 ◆
5.1	3.5	1.4	0.2	<i>setosa</i>
4.9	3.0	1.4	0.2	<i>setosa</i>
4.7	3.2	1.3	0.2	<i>setosa</i>
4.6	3.1	1.5	0.2	<i>setosa</i>
5.0	3.6	1.4	0.2	<i>setosa</i>
5.4	3.9	1.7	0.4	<i>setosa</i>
4.6	3.4	1.4	0.3	<i>setosa</i>

https://en.wikipedia.org/wiki/Iris_flower_data_set

bilibili



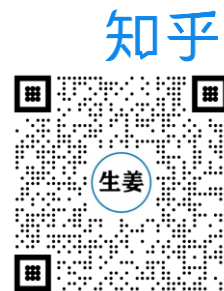
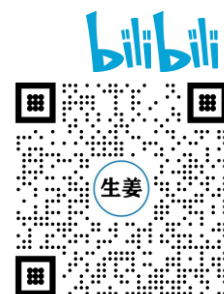
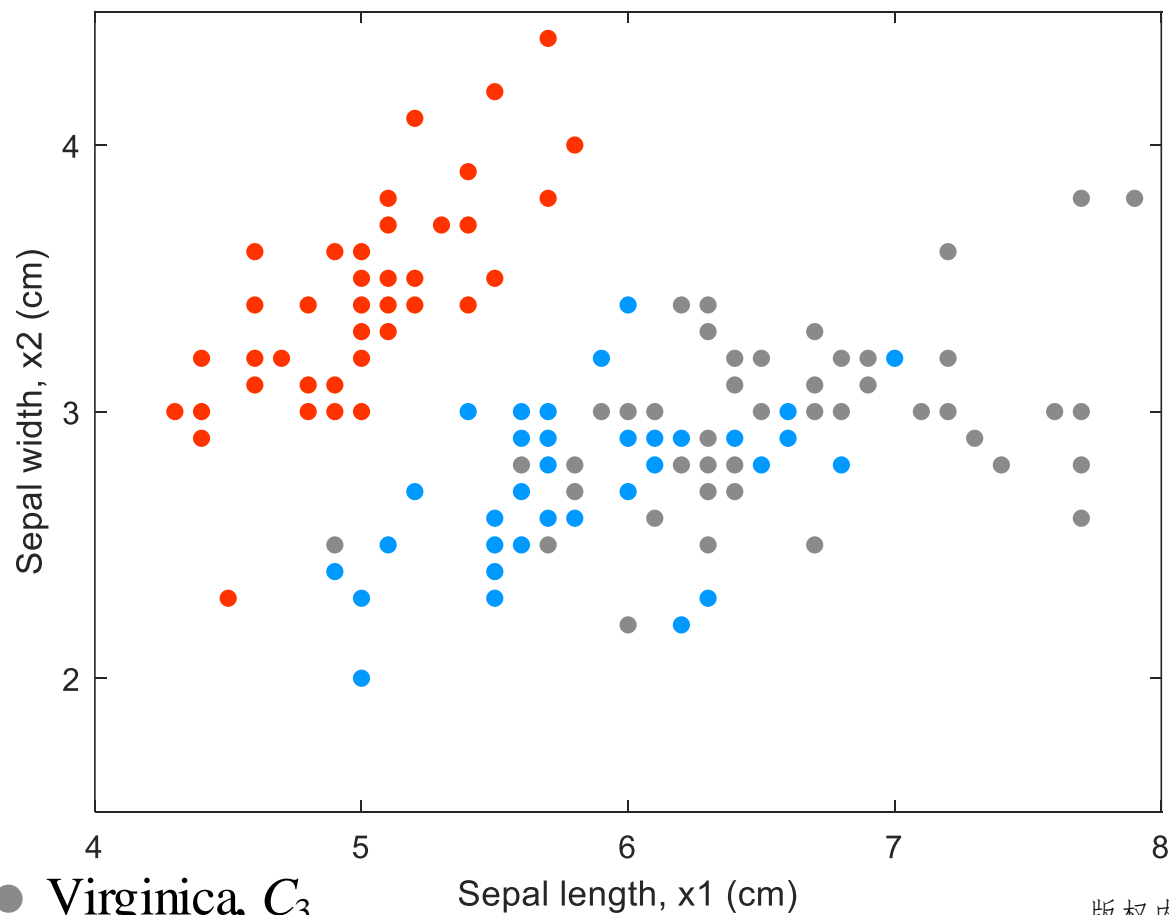
知乎




```
# import the iris data
iris = datasets.load_iris()

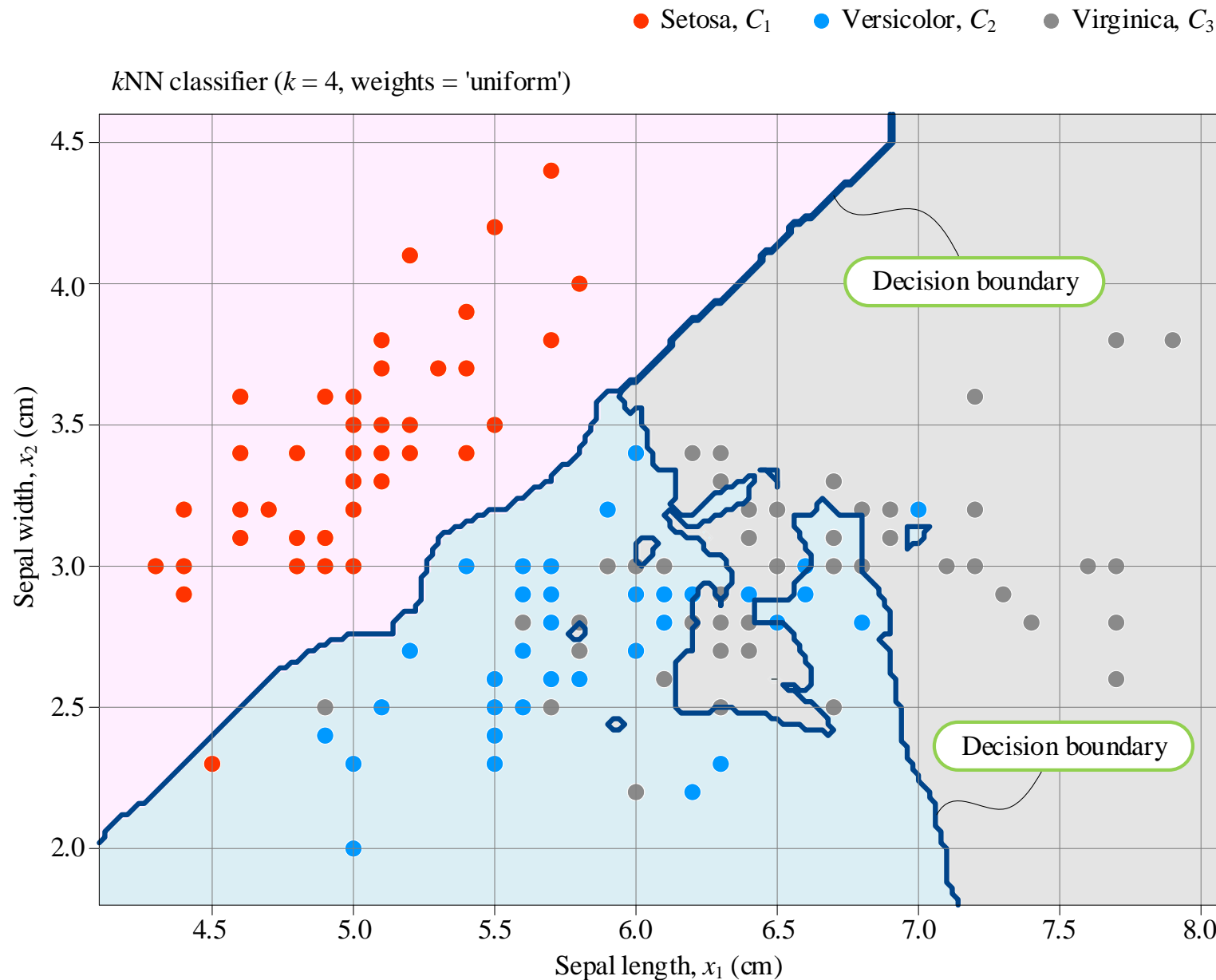
# Only use the first two features: sepal length, sepal width
X = iris.data[:, :2]

# Vector of labels
y = iris.target
```

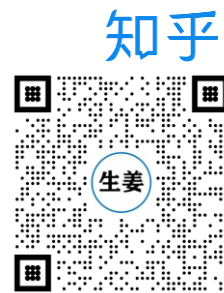
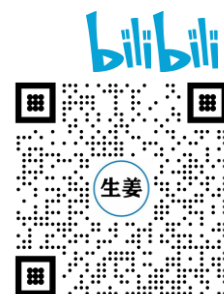


k临近分类结果

10

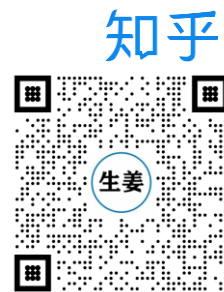
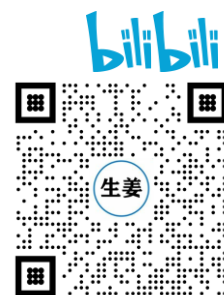
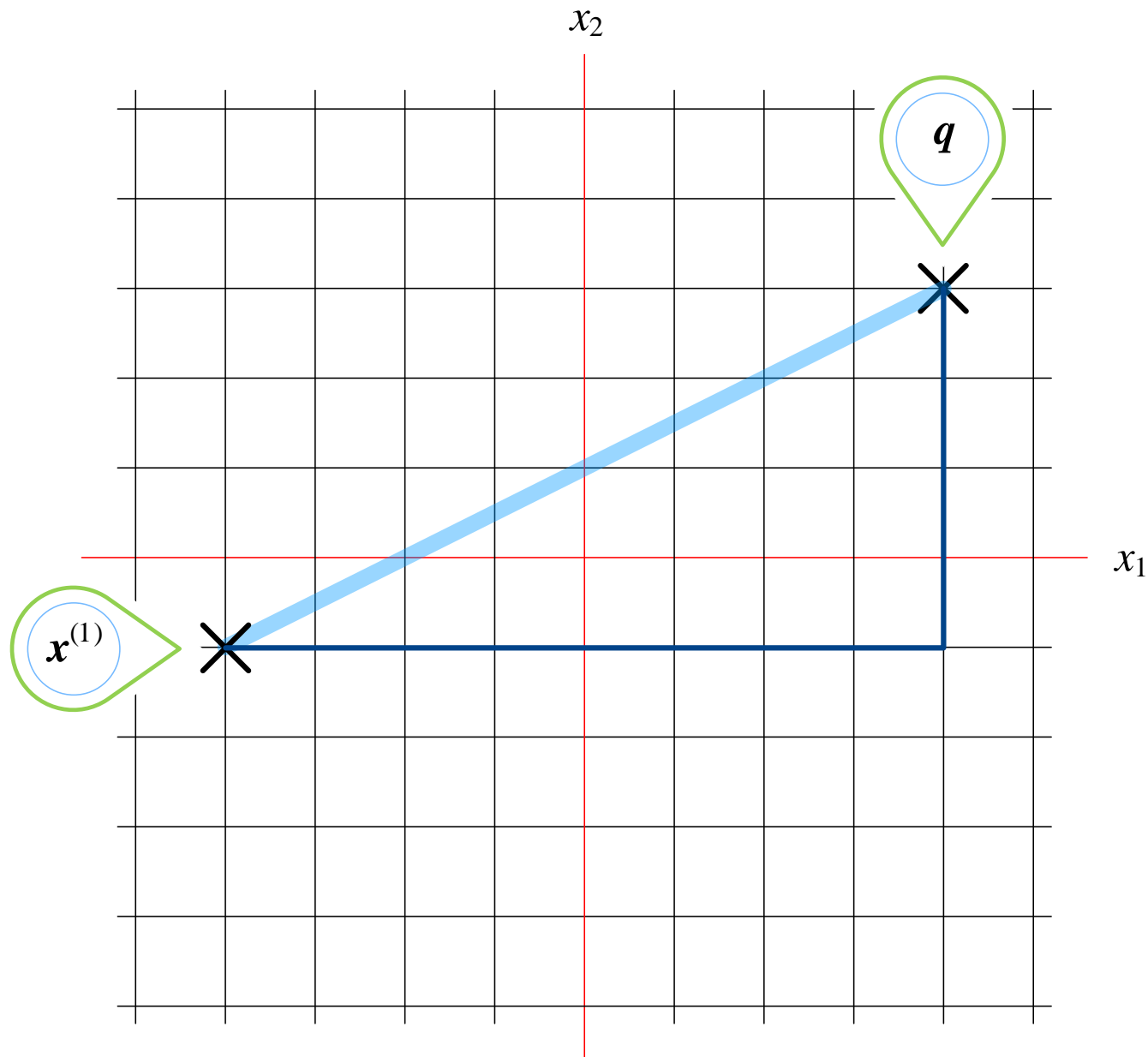


采用2个特征 (花萼长度 和 花萼宽度) 分类三种鸢尾花



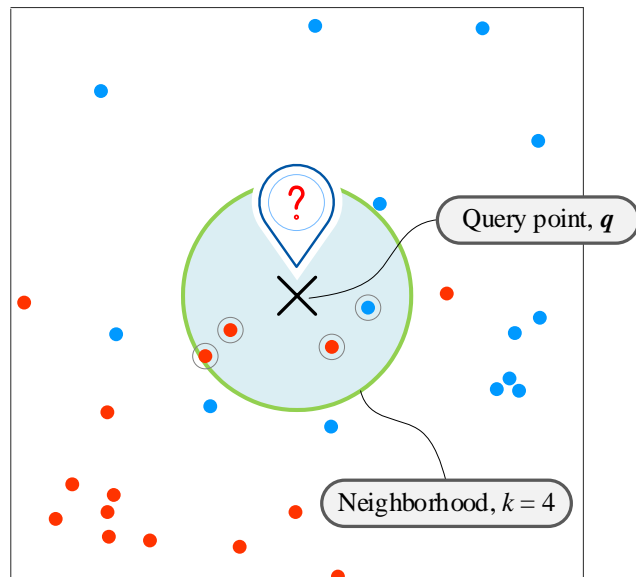
欧氏距离

11

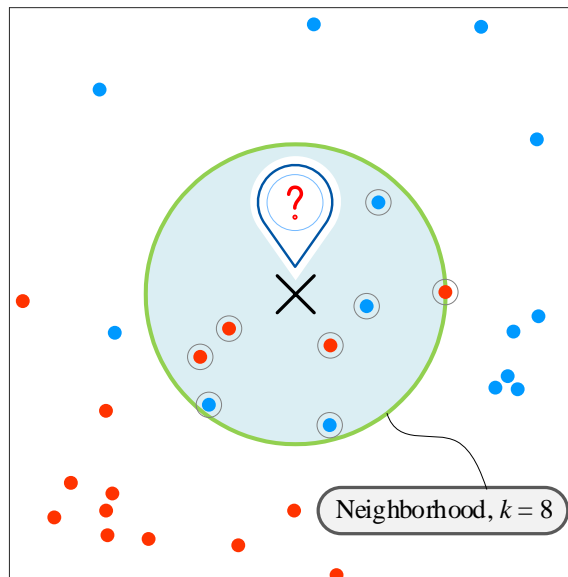


k 近邻数量

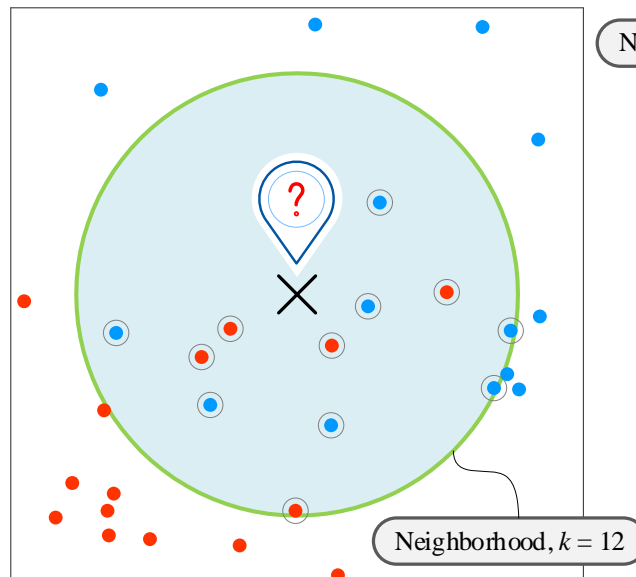
(a) $k = 4$; $\bullet C_1$ (3/4, 75%); $\bullet C_2$ (1/4, 25%)



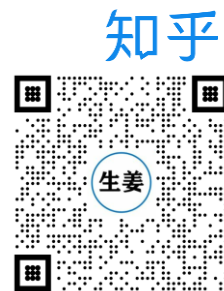
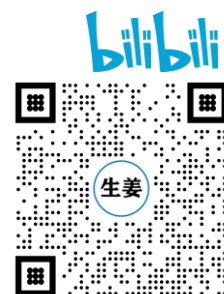
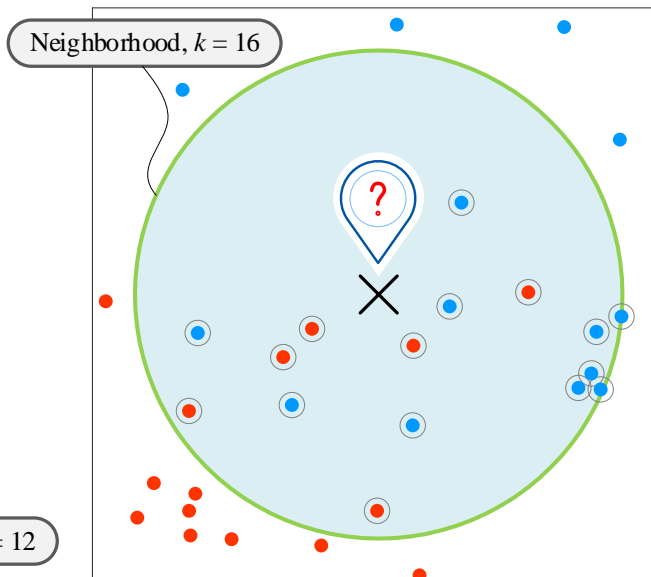
(b) $k = 8$; $\bullet C_1$ (4/8, 50%); $\bullet C_2$ (4/8, 50%)



(c) $k = 12$; $\bullet C_1$ (5/12, 41.67%); $\bullet C_2$ (7/12, 58.33%)



(d) $k = 16$; $\bullet C_1$ (6/16, 37.5%); $\bullet C_2$ (10/16, 62.5%)

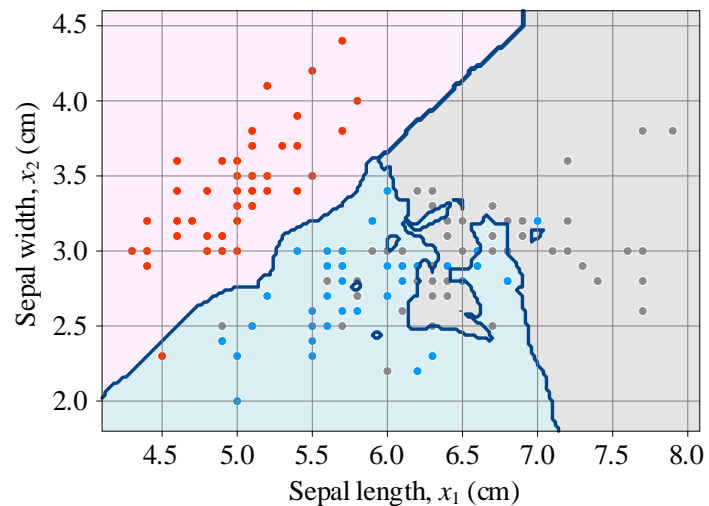


k 选取不同值时对鸢尾花分类影响

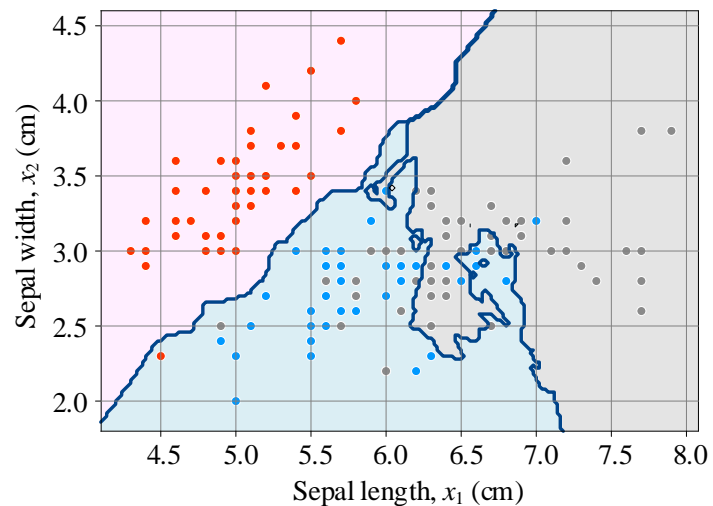
13

● Setosa, C_1 ● Versicolor, C_2 ● Virginica, C_3

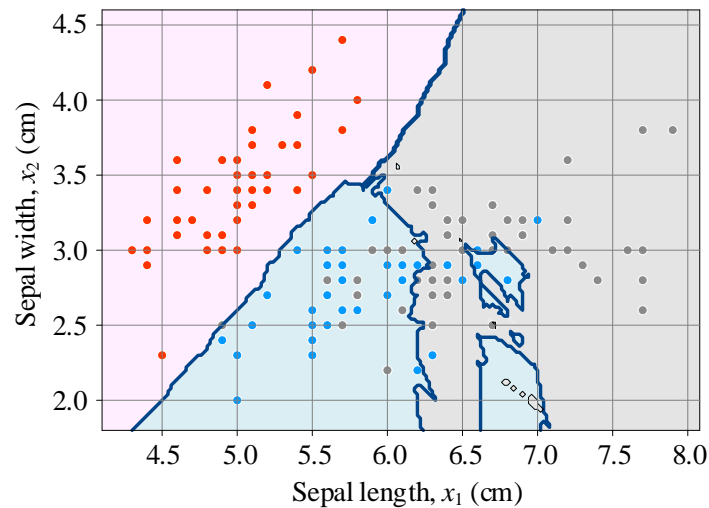
(a) k NN classifier ($k = 4$, weights = 'uniform')



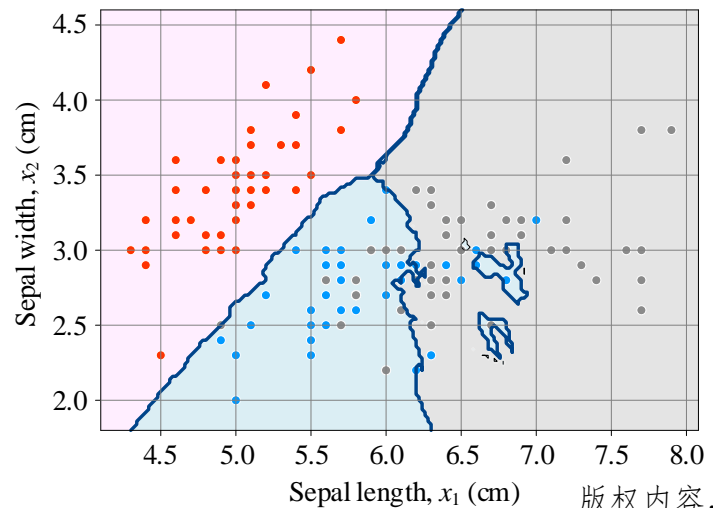
(b) k NN classifier ($k = 8$, weights = 'uniform')



(c) k NN classifier ($k = 12$, weights = 'uniform')



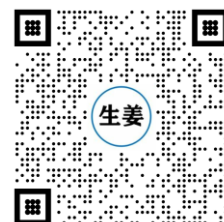
(d) k NN classifier ($k = 16$, weights = 'uniform')



bilibili

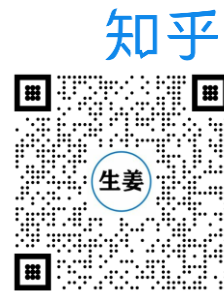
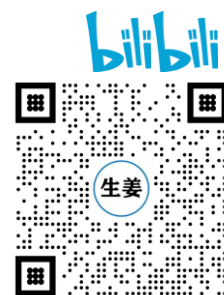
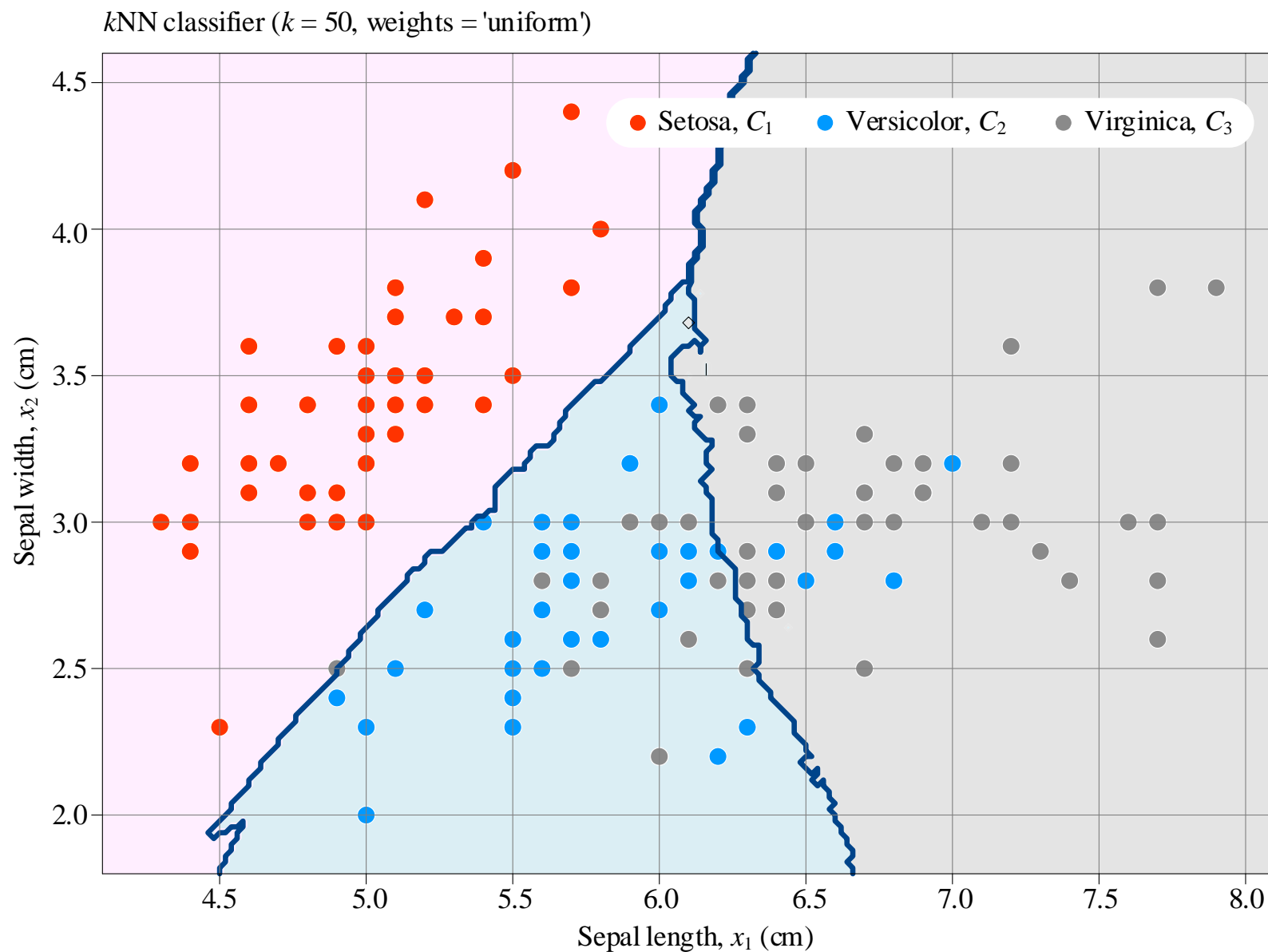


知乎

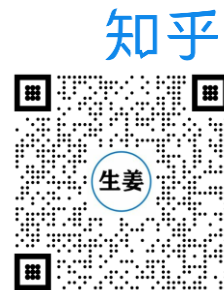
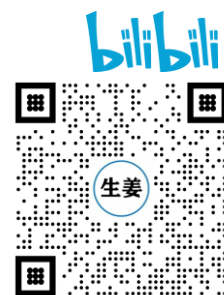


$k = 50$ 时，鸢尾花分类决策边界，kNN，等权重投票

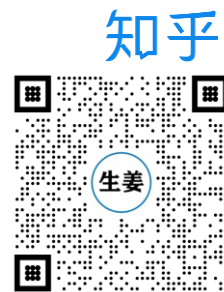
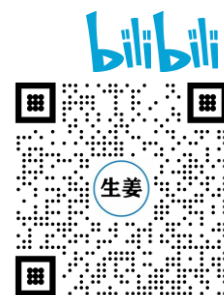
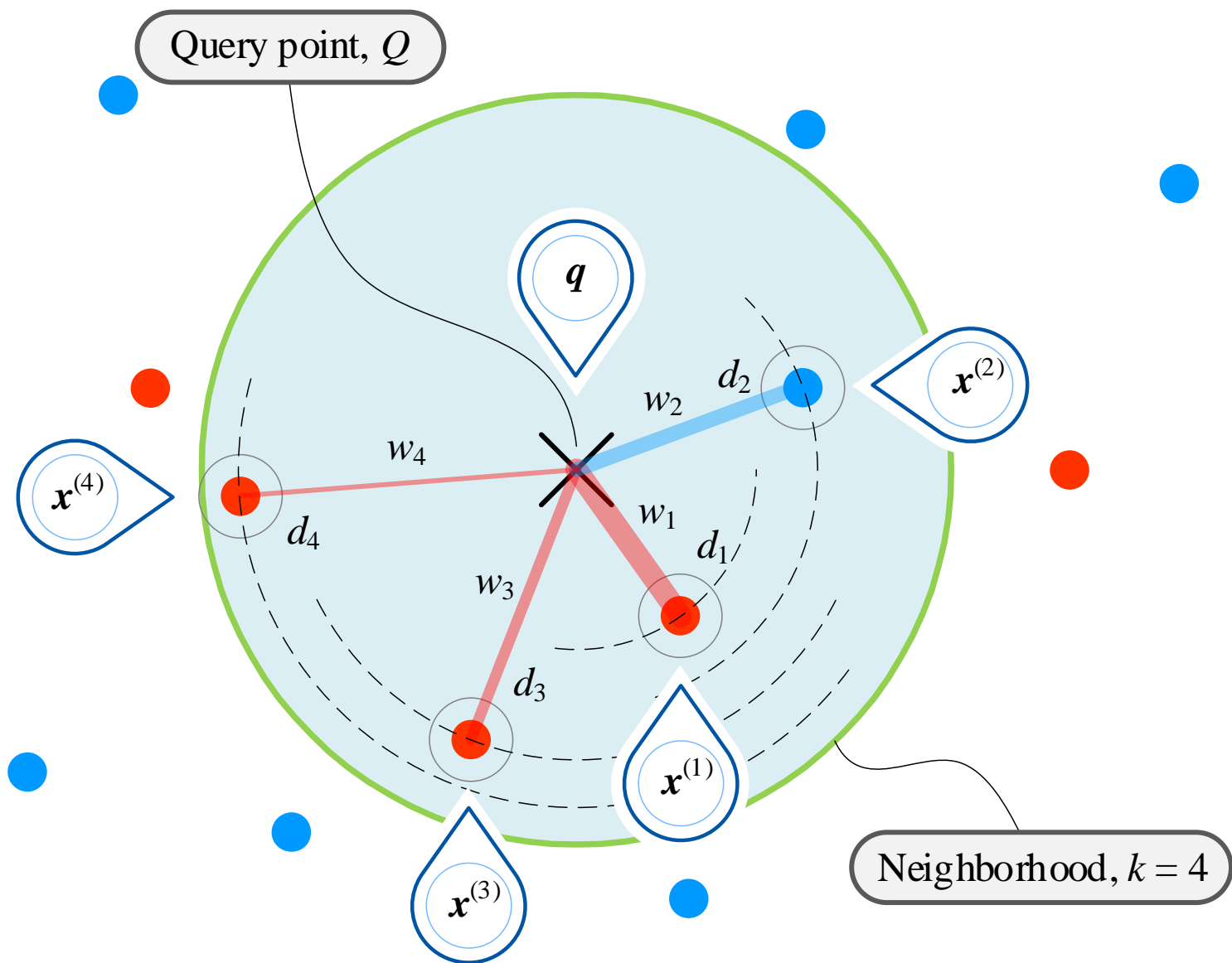
14



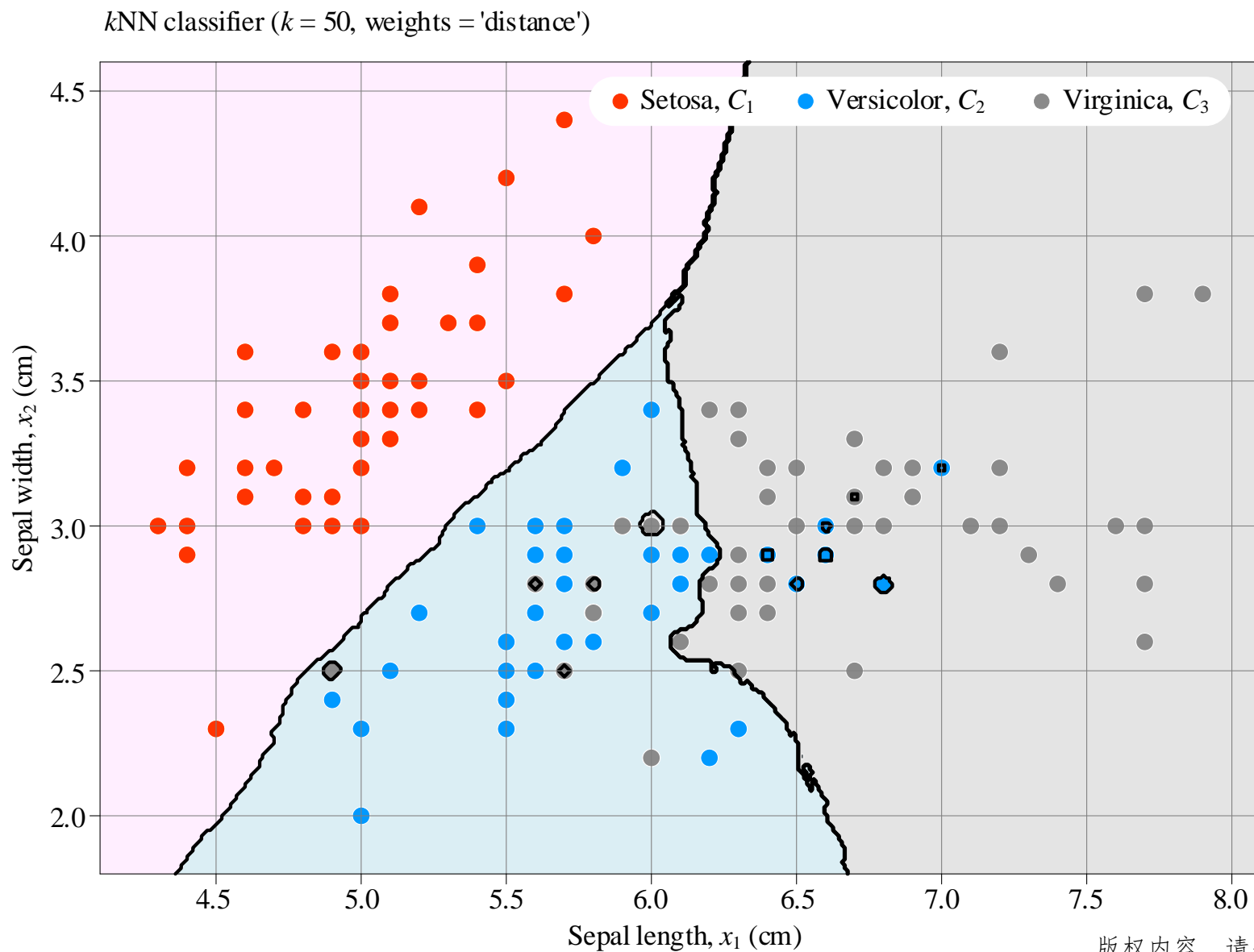
$$y^{(q)} = \arg \max_{C_l} \sum_{i \in kNN(q)} w_i \cdot I(y^{(i)} = C_l)$$



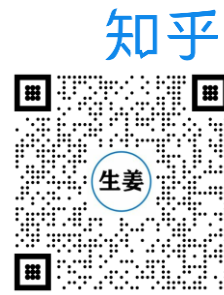
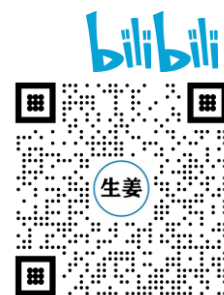
● Class 1, C_1 ● Class 2, C_2



$k = 50$ 时，鸢尾花分类决策边界，kNN



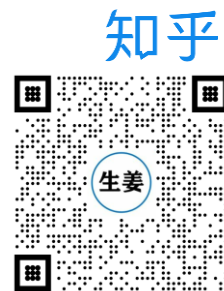
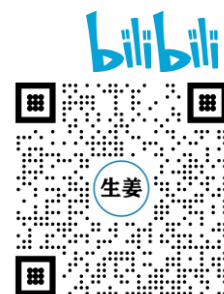
投票权重与查询点距离成反



$$\mu_l = \frac{1}{\text{count}(C_l)} \sum_{i \in C_l} \mathbf{x}_i$$

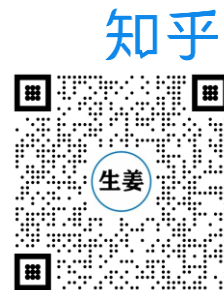
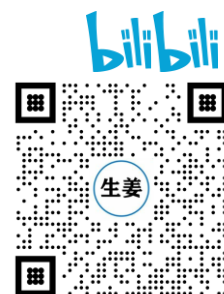
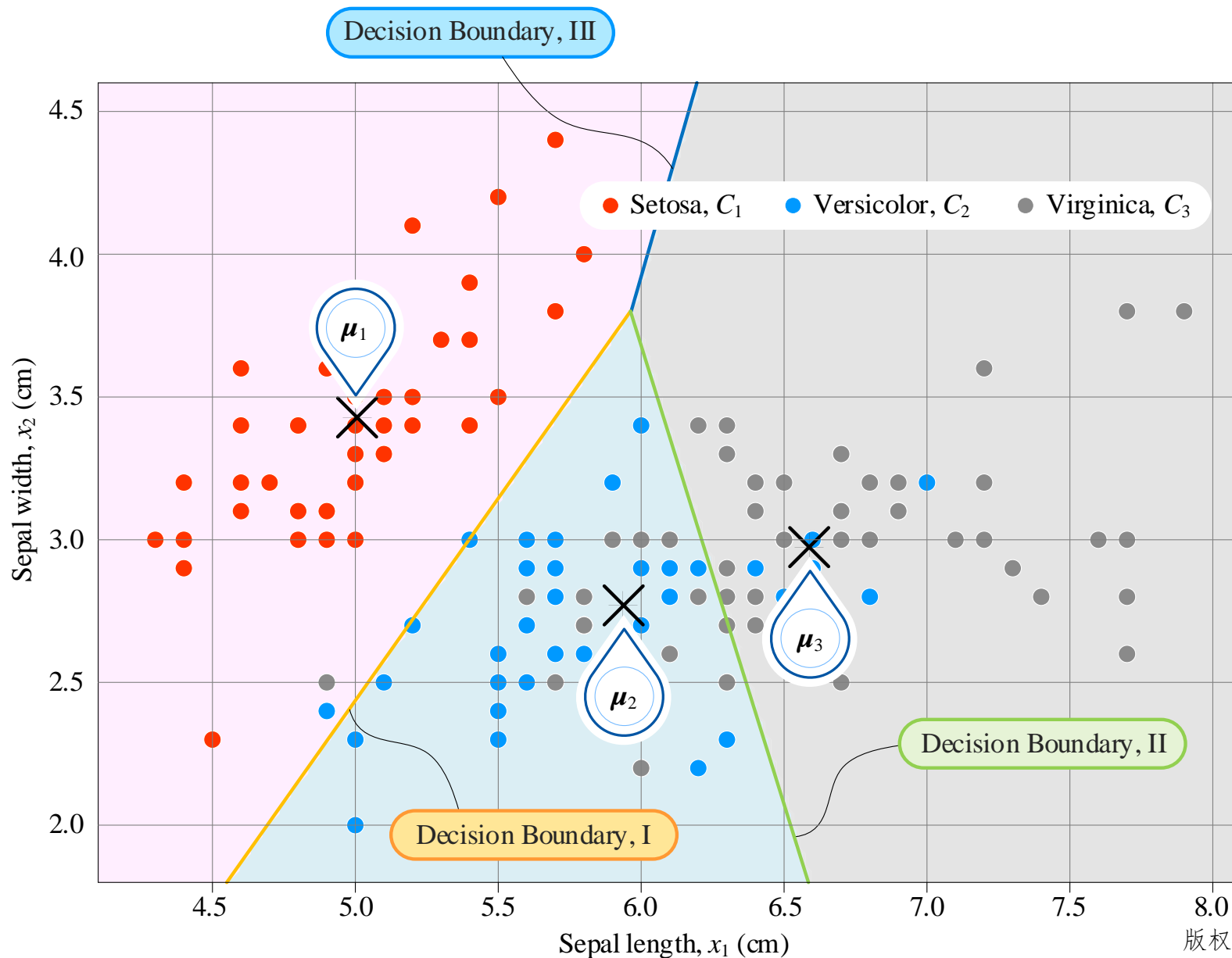
$$\mathbf{x}^{(2)} = [2 \ 3], \quad \mathbf{x}^{(5)} = [3 \ 1], \quad \mathbf{x}^{(6)} = [-2 \ 2], \quad \mathbf{x}^{(9)} = [1 \ 6]$$

$$\begin{aligned} \mu_1 &= \frac{1}{\text{count}(C_1)} \sum_{i \in C_1} \mathbf{x}_i = \frac{1}{\text{count}(C_1)} (\mathbf{x}^{(2)} + \mathbf{x}^{(5)} + \mathbf{x}^{(6)} + \mathbf{x}^{(9)}) \\ &= \frac{1}{4} ([2 \ 3] + [3 \ 1] + [-2 \ 2] + [1 \ 6]) \\ &= [1 \ 3] \end{aligned}$$



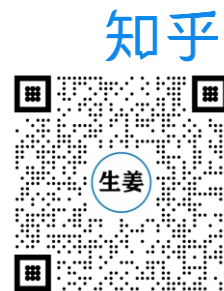
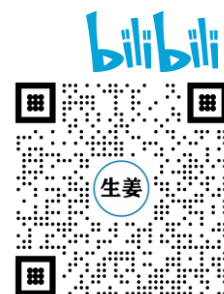
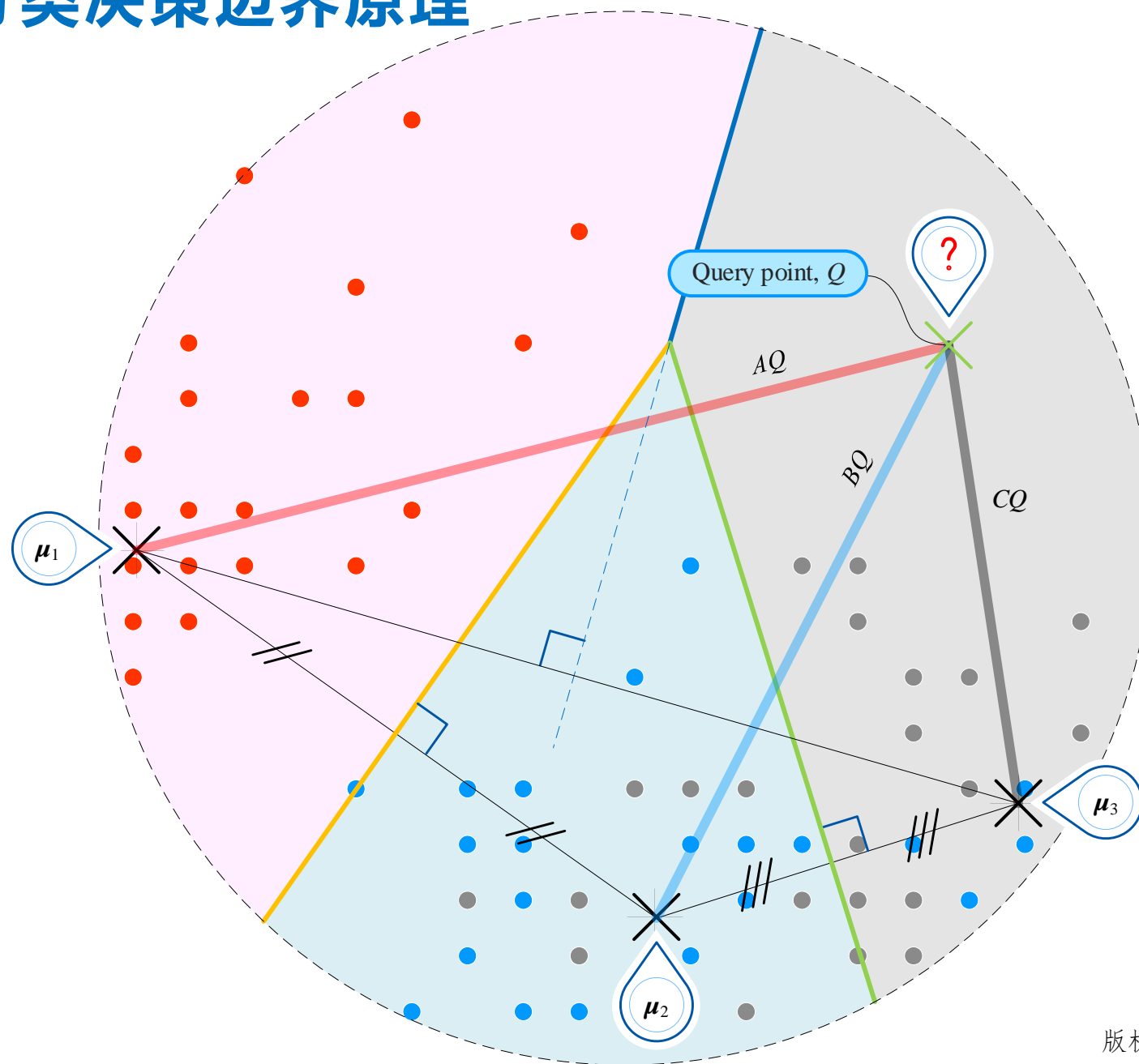
鸢尾花分类决策边界，最近质心分类

19



最近质心分类决策边界原理

20

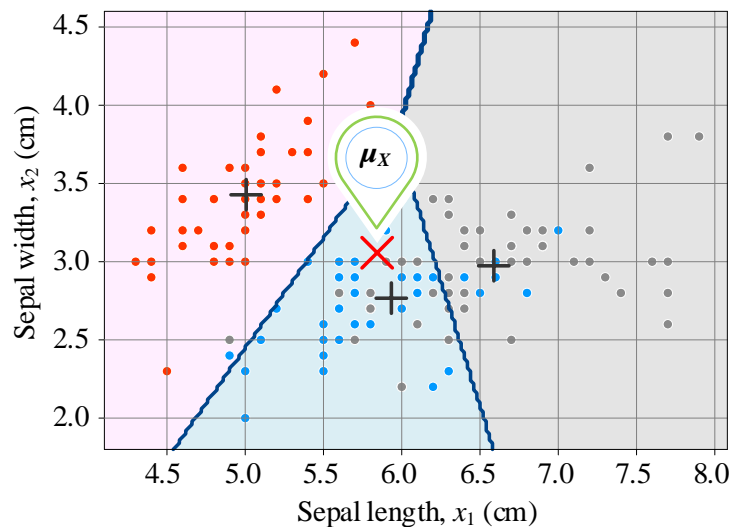


收缩阈值增大对决策边界影响

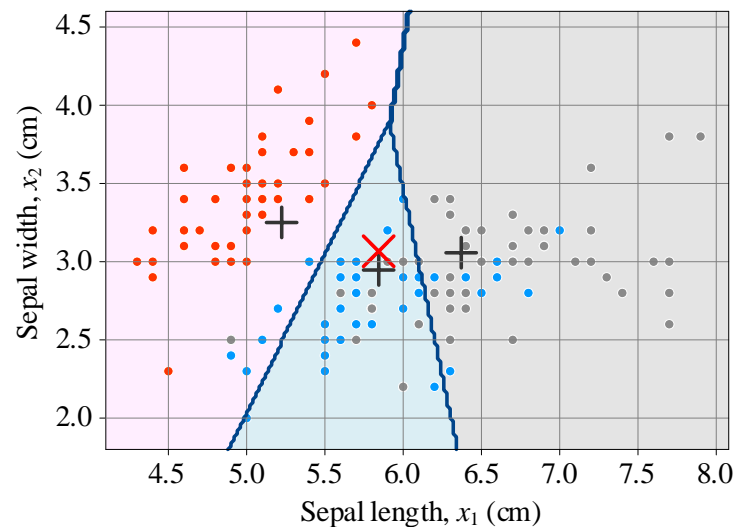
21

● Setosa, C_1 ● Versicolor, C_2 ● Virginica, C_3

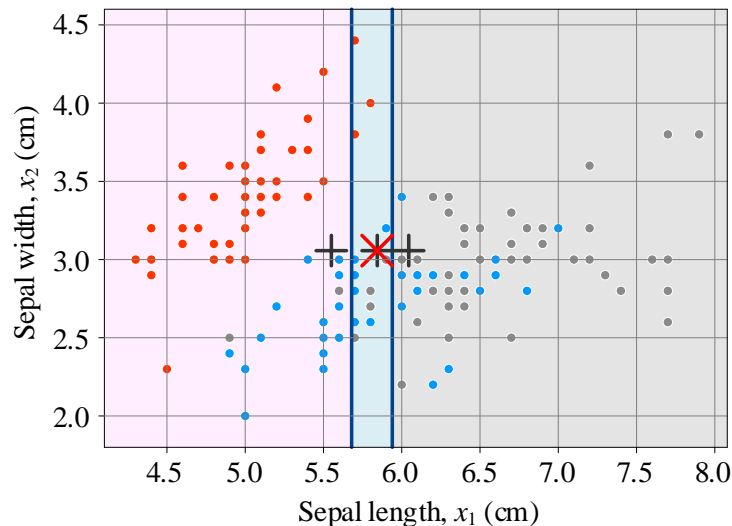
(a) NCC, shrink threshold = None



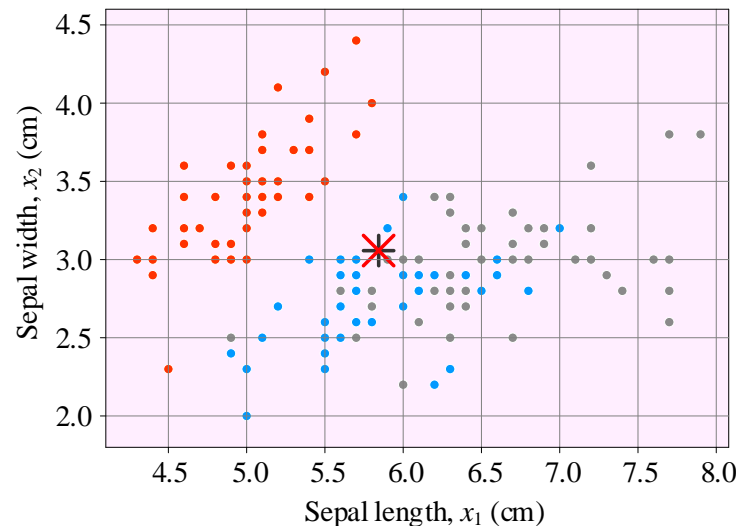
(b) NCC, shrink threshold = 2



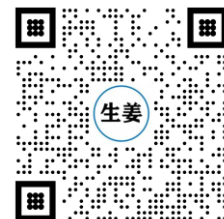
(c) NCC, shrink threshold = 5



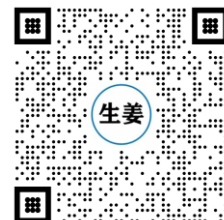
(d) NCC, shrink threshold = 8



bilibili

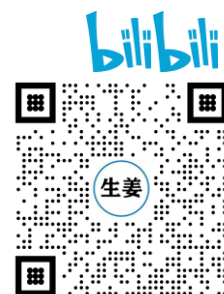
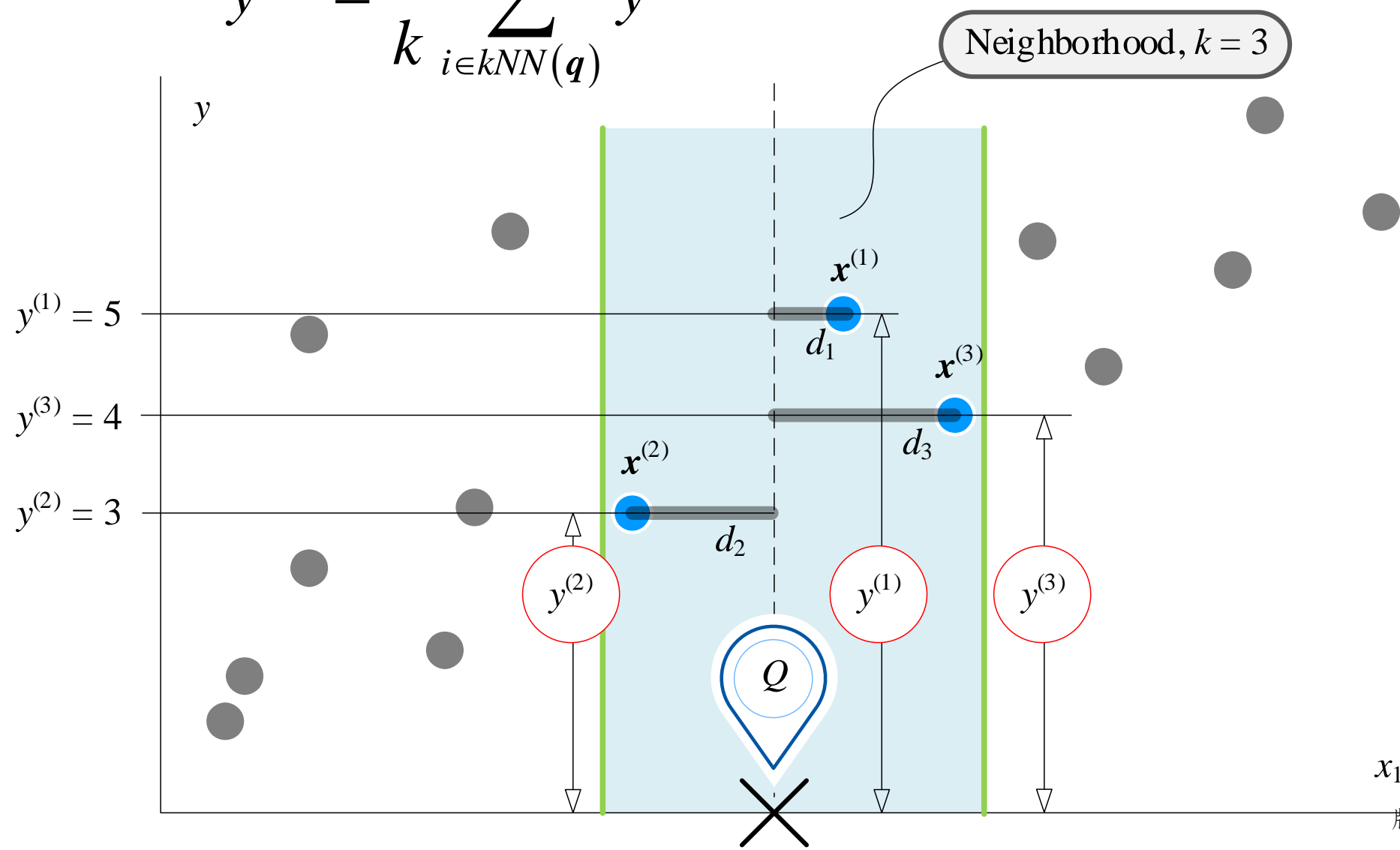


知乎



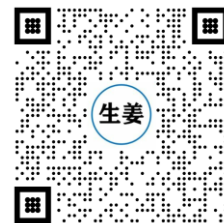
$$y^{(q)} = \frac{1}{k} \sum_{i \in kNN(q)} y^{(i)}$$

$$y^{(q)} = \frac{1}{3} \left(y^{(1)} + y^{(2)} + y^{(3)} \right) = \frac{1}{3} (5 + 3 + 4) = 4$$



bilibili

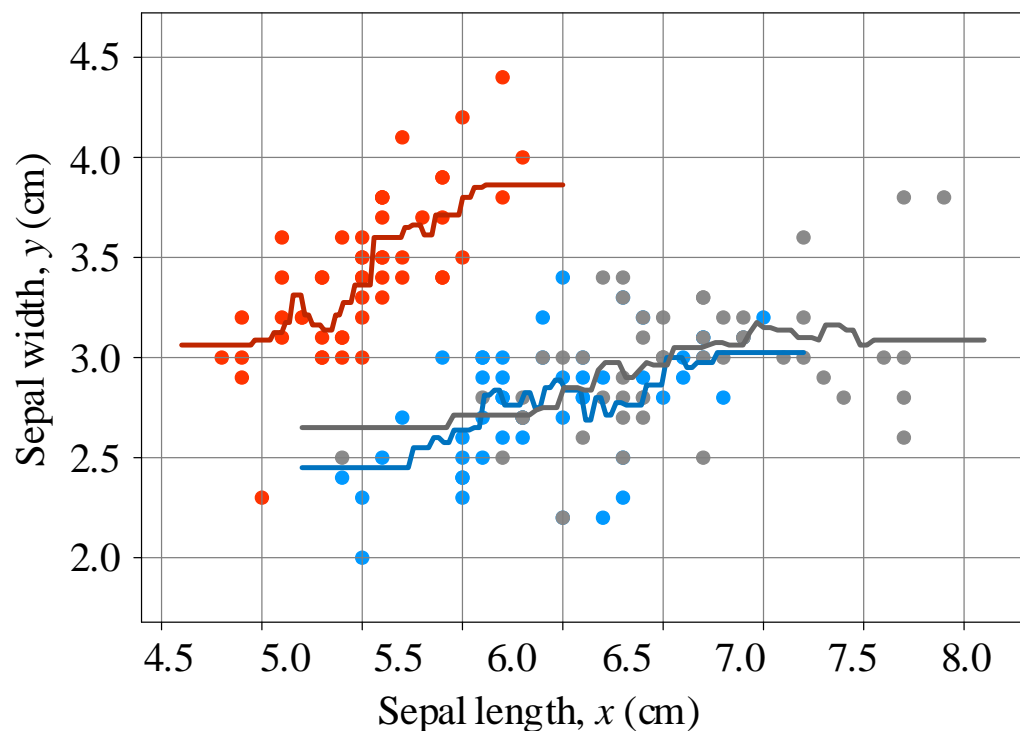
知乎



k NN回归，不同种类鸢尾花花萼长度 x_1 和花萼宽度 x_2 回归关系

● Setosa, C_1 ● Versicolor, C_2 ● Virginica, C_3

(a) k NN regressor, $k = 8$, weights = 'uniform'



(b) k NN regressor, $k = 8$, weights = 'distance'

