

## K-means with IRIS data

In [7]:

```
from sklearn import datasets
import pandas as pd
iris = datasets.load_iris()
labels = pd.DataFrame(iris.target)
labels.columns=['labels']
data = pd.DataFrame(iris.data)
data.columns=['Sepal length', 'Sepal width', 'Petal length', 'Petal width']
data = pd.concat([data, labels], axis=1)

data.head()
```

Out[7]:

	Sepal length	Sepal width	Petal length	Petal width	labels
0	5.1	3.5	1.4	0.2	0
1	4.9	3.0	1.4	0.2	0
2	4.7	3.2	1.3	0.2	0
3	4.6	3.1	1.5	0.2	0
4	5.0	3.6	1.4	0.2	0

## Extract feature

In [14]:

```
feature = data[ ['Sepal length', 'Sepal width']]
feature.head()
```

Out[14]:

	Sepal length	Sepal width
0	5.1	3.5
1	4.9	3.0
2	4.7	3.2
3	4.6	3.1
4	5.0	3.6

## Create model, training & Prediction

In [15]:

```
from sklearn.cluster import KMeans
import matplotlib.pyplot as plt
import seaborn as sns

# create model and prediction
model = KMeans(n_clusters=3, algorithm='auto')
model.fit(feature)
predict = pd.DataFrame(model.predict(feature))
predict.columns=['predict']

# concatenate labels to df as a new column
r = pd.concat([feature, predict], axis=1)

print(r)
```

	Sepal length	Sepal width	predict
0	5.1	3.5	0
1	4.9	3.0	0
2	4.7	3.2	0
3	4.6	3.1	0
4	5.0	3.6	0
...	...	...	...
145	6.7	3.0	1
146	6.3	2.5	2
147	6.5	3.0	1
148	6.2	3.4	1
149	5.9	3.0	2

[150 rows x 3 columns]

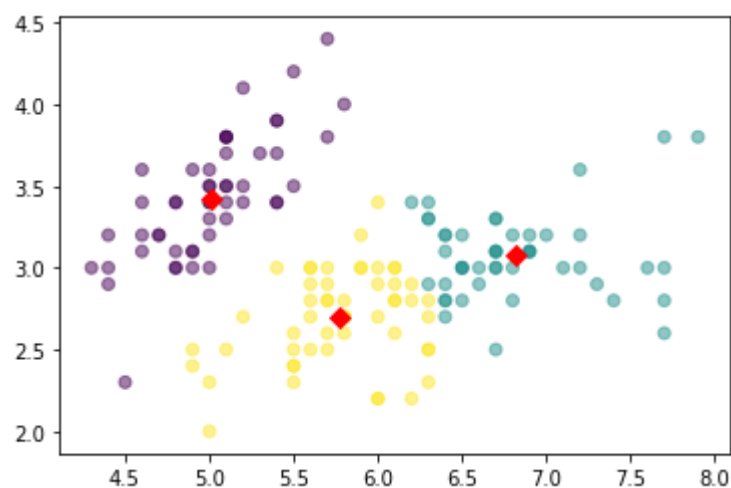
## visualize result

In [16]:

```
centers = pd.DataFrame(model.cluster_centers_, columns=['Sepal length', 'Sepal width'])
center_x = centers['Sepal length']
center_y = centers['Sepal width']
```

```
# scatter plot
```

```
plt.scatter(r['Sepal length'], r['Sepal width'], c=r['predict'], alpha=0.5)
plt.scatter(center_x, center_y, s=50, marker='D', c='r')
plt.show()
```



## Evaluate model with Cross tabuliazation

In [17]:

```
ct = pd.crosstab(data['labels'], r['predict'])
print(ct)
```

```
predict  0  1  2
labels
0         50  0  0
1          0 12 38
2          0 35 15
```

## Determine number of clusters with Inertia value

In [19]:

```
ks = range(1,10)
inertias = []

for k in ks:
    model = KMeans(n_clusters=k)
    model.fit(feature)
    inertias.append(model.inertia_)    #inertia : 군집화가된 후에, 각 중심점에서 군집의 데이터간의 거

# Plot ks vs inertias
plt.plot(ks, inertias, '-o')
plt.xlabel('number of clusters, k')
plt.ylabel('inertia')
plt.xticks(ks)
plt.show()    #3-5개가 적당
```

