

In [1]:

```
from sklearn.datasets import load_boston
from sklearn.model_selection import train_test_split
from sklearn.cluster import KMeans
from sklearn import preprocessing
from sklearn.preprocessing import MaxAbsScaler
from yellowbrick.cluster import KElbowVisualizer
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline

pd.set_option('display.max_columns', 500)
boston_data = load_boston()
df = pd.DataFrame(boston_data.data, columns=boston_data.feature_names)
df['target'] = pd.Series(boston_data.target)
```

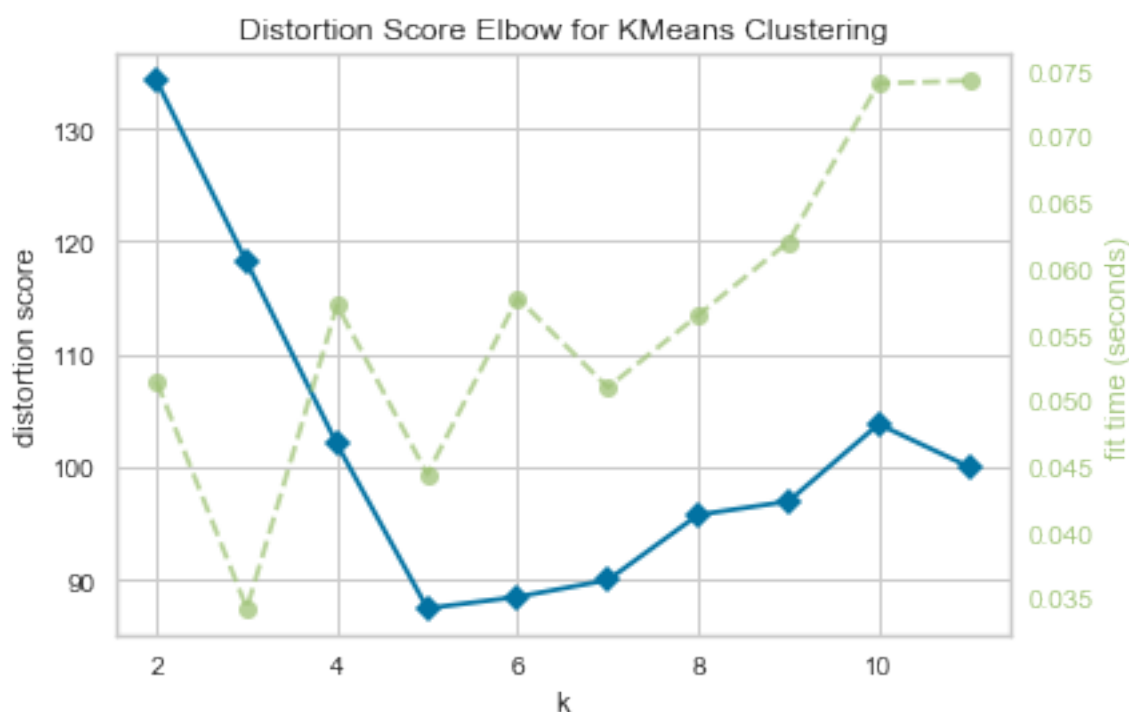
In [2]:

```
#plot the prediction to a cross validation set evaluate the difference in pred  
iction and reality add the average to the model.
```

In [3]:

```
L = df[['LSTAT', 'target']]
# L = preprocessing.scale(L)
mets = ['silhouette', 'calinski_harabaz', 'distortion']
km = KMeans()
visualize = KElbowVisualizer(km, k = (2, 12), metric=mets[2])

visualize.fit(L)
visualize.poof()
```



In [4]:

```
def centroid_sort(centroids):
    new_centroids = []
    c0s = np.sort(centroids[:,0])

    for j in c0s:
        for i in range(len(centroids)):
            if j == centroids[i][0]:
                new_centroids.append([j, centroids[i][1]])
            else:
                continue

    new_centroids = np.array(new_centroids)
    return new_centroids

def generate_points(poly,start,stop,number_of_points=100):
    y_points=[]
    dummies = np.linspace(start,stop,number_of_points)
    for i in range(len(dummies)):
        y_points.append(poly(dummies[i]))
    return dummies, y_points

def get_ave_coefficients(dimension_of_eq, scribbling):
    ave_coefficients = [0 for k in range(dimension_of_eq)]
    for j in range(dimension_of_eq):
        for i in range(len(scribbling)):
            ave_coefficients[j] += scribbling[f'run_{i}'][j]

        ave_coefficients[j] = ave_coefficients[j]/len(scribbling)
    return ave_coefficients
```

In [5]:

```
# # print(centroids)

# colors = ['azure', 'darkolivegreen', 'darkblue', 'm', 'y', 'green', 'r', 'pa
payawhip' ]

# labels = km.labels_

# plt.scatter(L[:,0], L[:,1])
# plt.scatter(centroids[:, 0], centroids[:, 1], marker = 'x',c='m', s=75, line
widths = 5, zorder=10)

# # thats the line of best fit
# x, y = generate_points(p, -1, 2.5, 100)

# plt.plot(x, y, 'go-')
```

In [16]:

```
model_testing = {}
```

```
# feeding through KMeans on sample, generating the centroids again and again with the best fit line
# Average over the result and return
# Assumptions data is in format already prepared for KMeans ie.: preprocess, operations, aggregations and so on.
```

```
def rigorous_test(Var_col, no_of_clusters, leading_order=3, iterations=23, preproc = False):
    ave_eq = []
    ave_centroids = []

    if not preproc:
        X = df[[Var_col]]
        y = df[['target']]

        for i in range(iterations):

            X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.35)
            data = np.column_stack((X_train, y_train))

            km = KMeans(n_clusters = no_of_clusters)
            km.fit(data)
            centroids = km.cluster_centers_
            centroids = centroid_sort(centroids) #sorting centroids in order of x coord

            # fitting current centroids

            cx = centroids[:,0]
            cy = centroids[:,1]

            ploy = np.polyfit(cx, cy, leading_order)

            ave_eq.append(ploy)
            ave_centroids.append(centroids)

    if preproc:
        testab = np.array(df[[Var_col, "target"]])
        transformer = MaxAbsScaler().fit(testab)
        H = transformer.transform(testab)

        for i in range(iterations):

            X_train, X_test, y_train, y_test = train_test_split(H[:,0], H[:,1], test_size = 0.35)
            data = np.column_stack((X_train, y_train))

            km = KMeans(n_clusters = no_of_clusters)
            km.fit(data)
            centroids = km.cluster_centers_
            centroids = centroid_sort(centroids) #sorting centroids in order of x coord

            # fitting current centroids
```

```

        cx = centroids[:,0]
        cy = centroids[:,1]

        ploy = np.polyfit(cx, cy, leading_order)

        ave_eq.append(ploy)
        ave_centroids.append(centroids)

    ave_coefficients = np.mean(ave_eq, axis = 0)

    average_line = np.polyld(ave_coefficients)

    ave_centroids = np.mean(ave_centroids, axis = 0)

    return average_line, ave_centroids

```

In [44]:

```

line, centroids = rigorous_test('PTRATIO', no_of_clusters=5, leading_order=2, iterations=45, preproc=True )
print('the line is: \n ', line)
print('\n \n' , 'Centroids are:\n',centroids)

plt.plot(centroids[:,0], centroids[:,1], 'ro')
xz, yz = generate_points(line ,0.6, 1, number_of_points=100)
plt.plot(xz, yz, 'g--')

```

the line is:

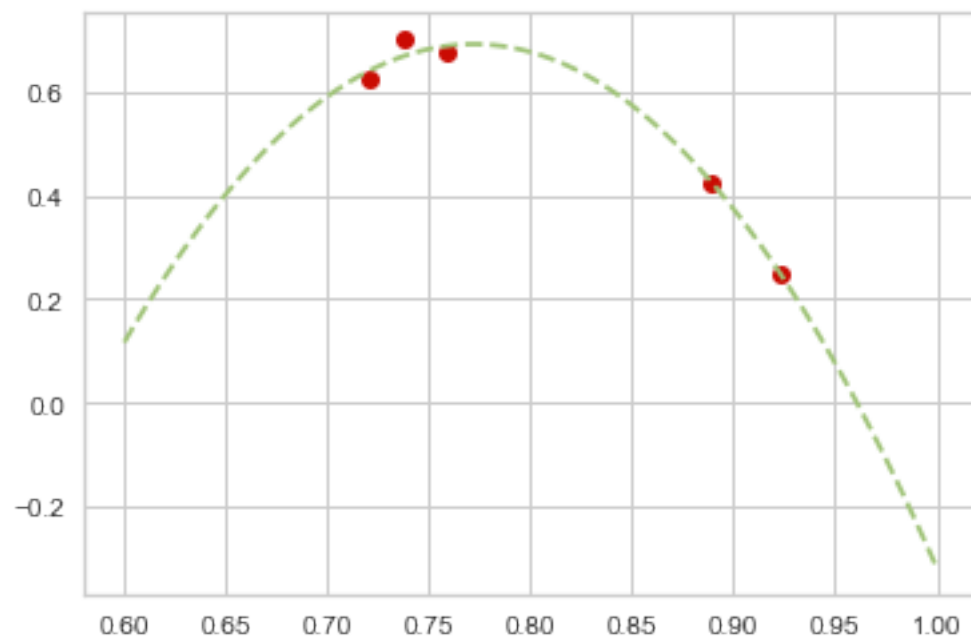
$$-19.47 x^2 + 30.07 x - 10.91$$

Centroids are:

```
[[0.72098668 0.62744311]  
[0.73749233 0.70206851]  
[0.75862239 0.67757107]  
[0.889471    0.42613279]  
[0.92317864 0.24920647]]
```

Out[44]:

[<matplotlib.lines.Line2D at 0x1239a1160>]



In [45]:

```
line.c
```

Out[45]:

```
array([-19.47162973,  30.06587534, -10.9131322  ])
```

In []: