

---

Python

# Matplotlib/빈도분석/kNN

---

김선녕(sykim.lecture@gmail.com)

참고문헌 : 파이썬을 이용한 빅데이터 수집, 분석과 시각화 - 비팬북스, 이원화

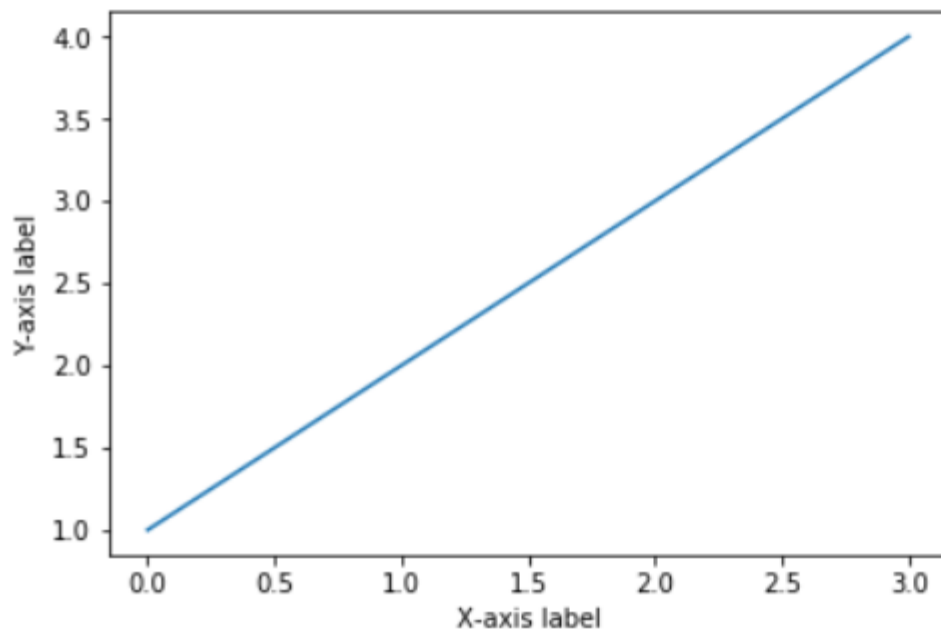
---

**Matplotlib**

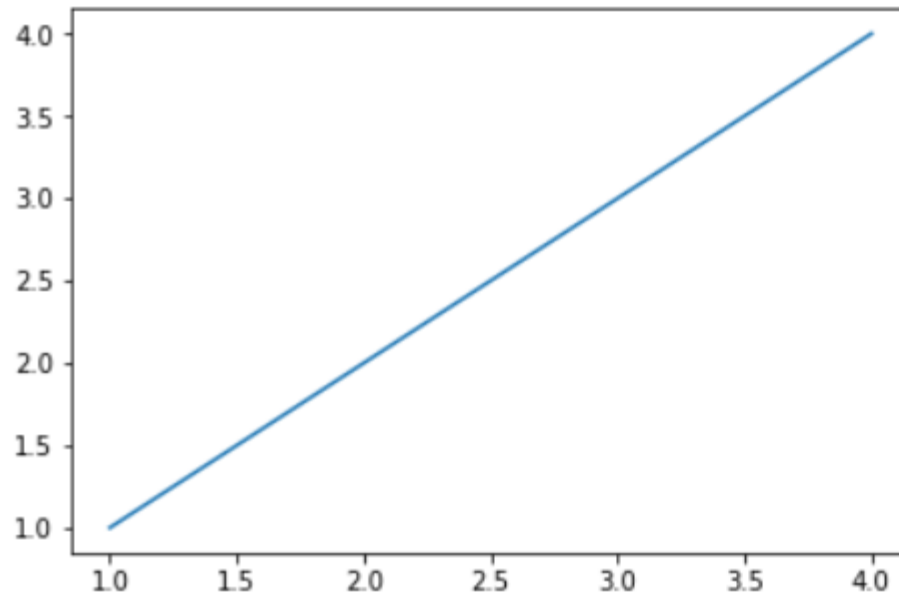
빈도분석

k-최근접 이웃 알고리즘(k-Nearest Neighbors)

```
In [2]: 1 import matplotlib.pyplot as plt
        2 plt.plot([1,2,3,4])
        3 plt.xlabel("X-axis label")
        4 plt.ylabel("Y-axis label")
        5 plt.show()
```

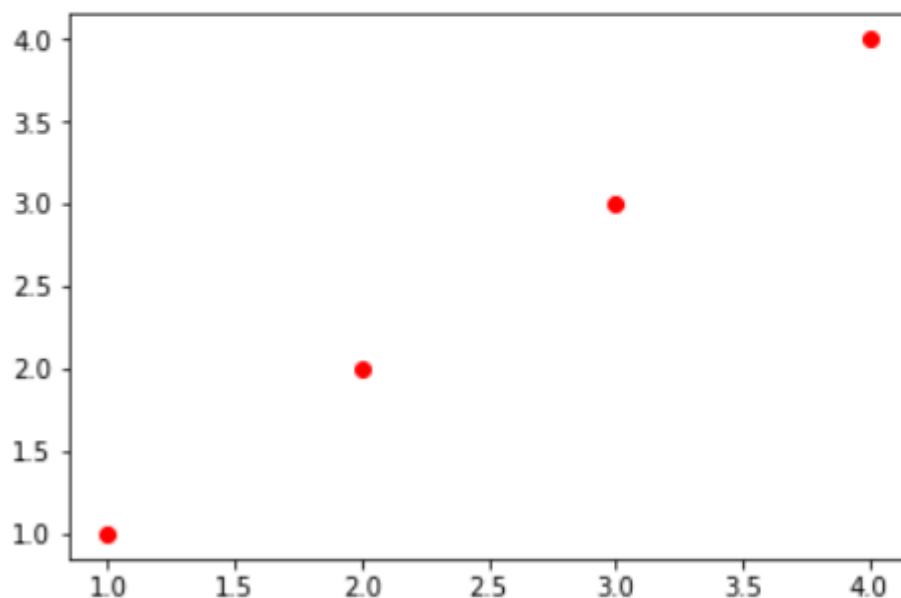


```
In [3]: 1 plt.plot([1,2,3,4],[1,2,3,4])  
        2 plt.show()
```



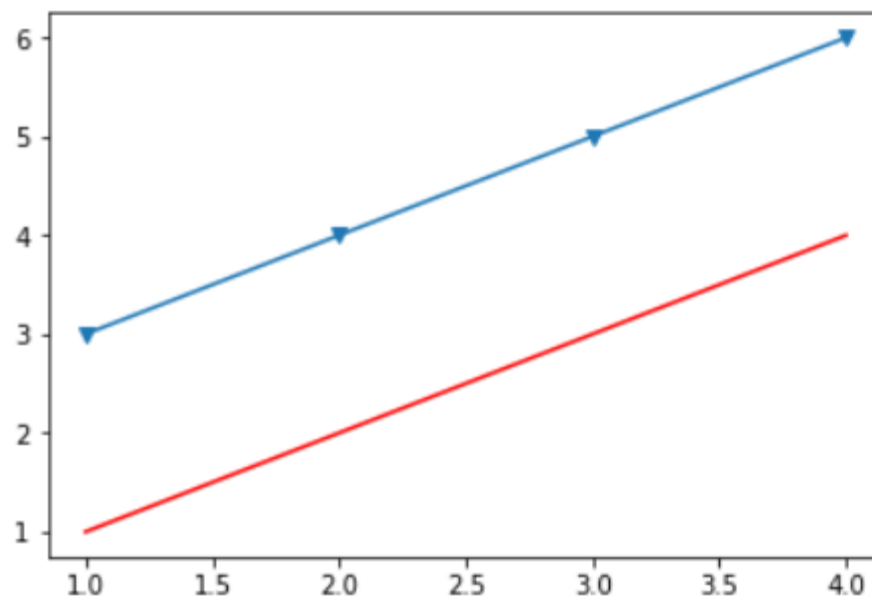
In [6]:

```
1 # 기본값 : 파란색(b) 라인(-)  
2 # ro : 적색 o, bv : 파란색 v 마크, 그외 matplotlib 공식사이트에서 확인  
3 plt.plot([1,2,3,4],[1,2,3,4],'ro')  
4 plt.show()
```



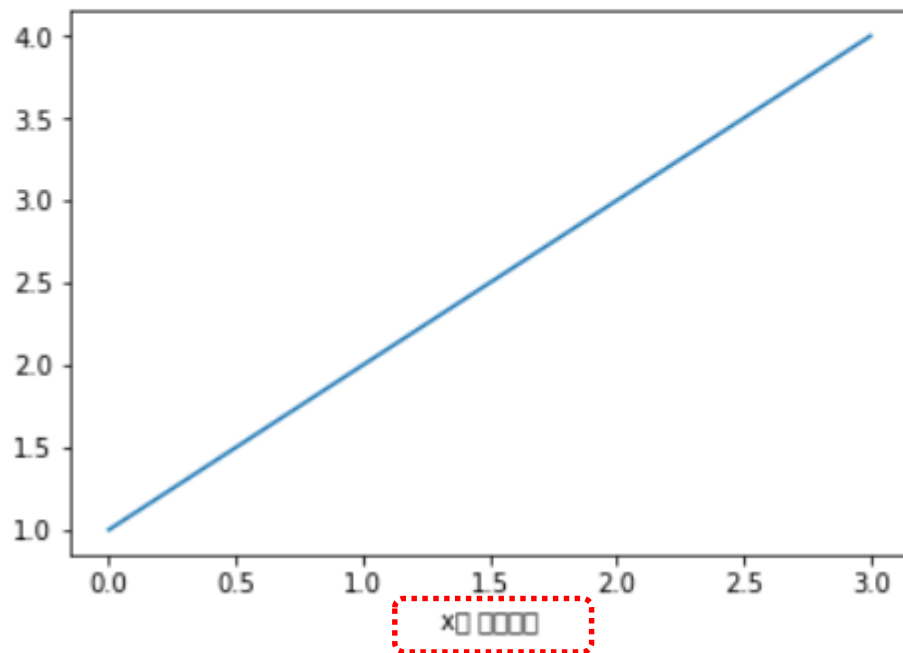
```
In [7]: 1 # 다수의 그래프 그리기  
2 plt.plot([1,2,3,4],[1,2,3,4], 'r-',[1,2,3,4],[3,4,5,6], 'v-')
```

```
Out[7]: [<matplotlib.lines.Line2D at 0x22df8386828>,  
<matplotlib.lines.Line2D at 0x22df83869e8>]
```



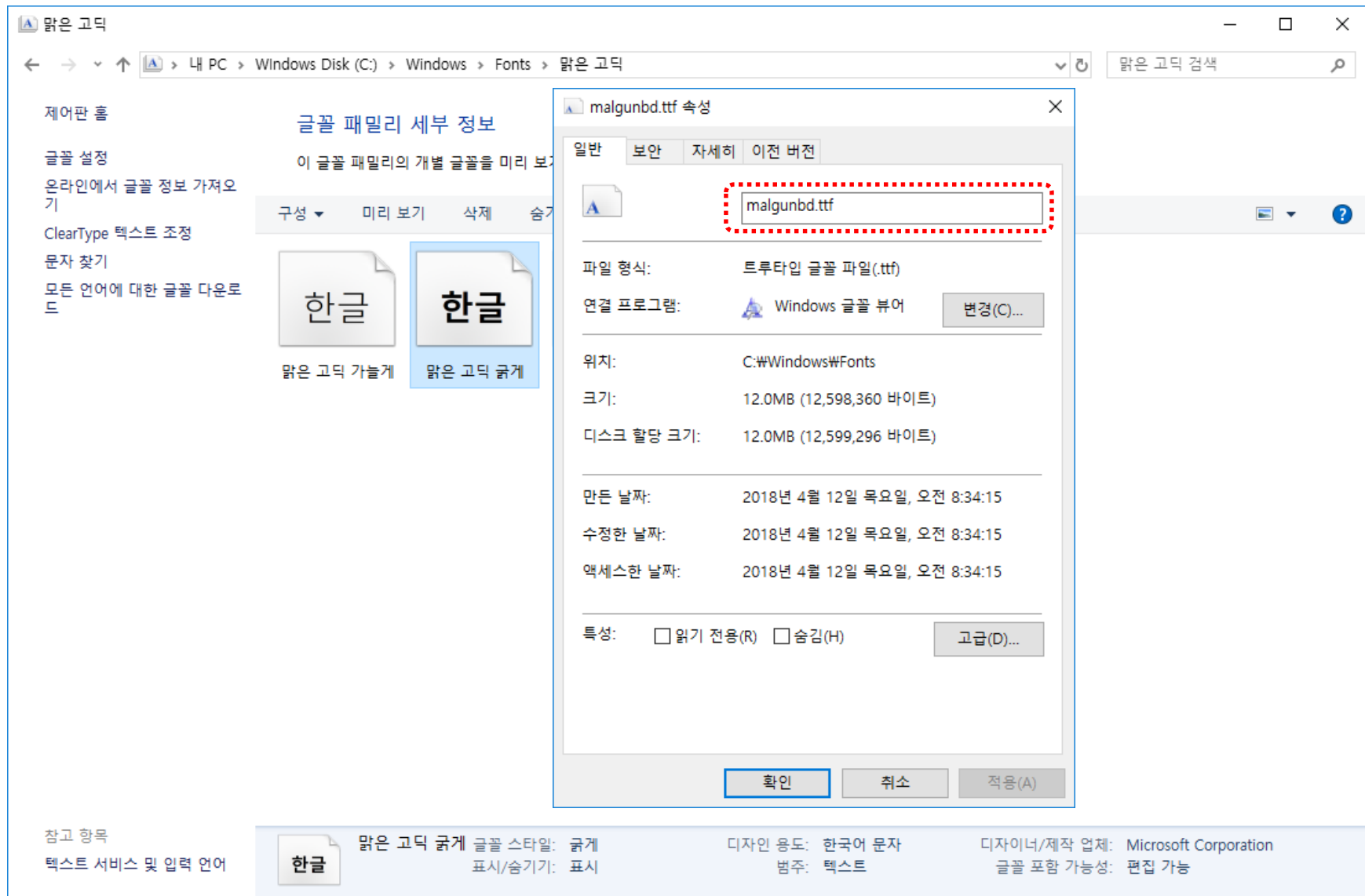
```
In [8]: 1 plt.plot([1,2,3,4])  
        2 plt.xlabel('x축 한글표시')  
        3 plt.show
```

Out[8]: <function matplotlib.pyplot.show(\*args, \*\*kw)>



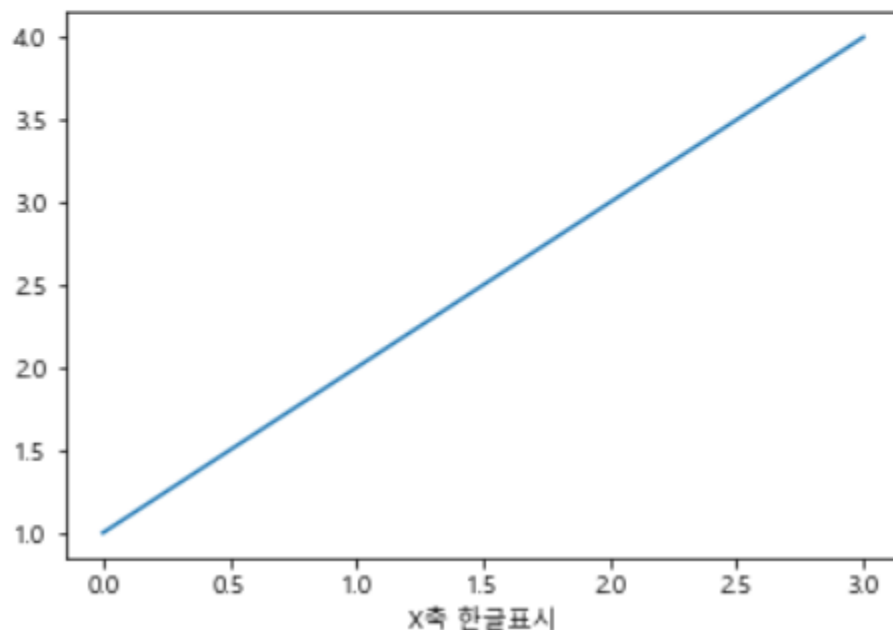
이름	글꼴 스타일	표시/숨기기	디자인 용도	범주	디자이너/제작 업체
궁서체 보통	보통	표시	한국어 문자	텍스트	Han Yang I & C Co. Ltd.
돋움 보통	보통	표시	한국어 문자	텍스트	Han Yang I & C Co. Ltd.
돋움체 보통	보통	표시	한국어 문자	텍스트	Han Yang I & C Co. Ltd.
<b>맑은 고딕</b>	<b>보통; 가늘게; 굵게</b>	<b>표시</b>	<b>한국어 문자</b>	<b>텍스트</b>	<b>Microsoft Corporation</b>
문체부 궁서체 정자체 보통	보통	표시			
문체부 궁서체 흘림체 보통	보통	표시			
문체부 돋움체 보통	보통	표시			
문체부 바탕체 보통	보통	표시			
문체부 쓰기 정체 보통	보통	표시			
문체부 쓰기 흘림체 보통	보통	표시			
문체부 제목 돋움체 보통	보통	표시			
문체부 제목 바탕체 보통	보통	표시			
문체부 훈민정음체 보통	보통	표시			
바탕 보통	보통	표시	한국어 문자	텍스트	Han Yang I & C Co. Ltd.
바탕체 보통	보통	표시	한국어 문자	텍스트	Han Yang I & C Co. Ltd.
새굴림 보통	보통	표시	한국어 문자	텍스트	HanYang System Co., LTD.
안상수2006가는 보통	보통	표시			
안상수2006굵은 보통	보통	표시			
안상수2006중간 보통	보통	표시			





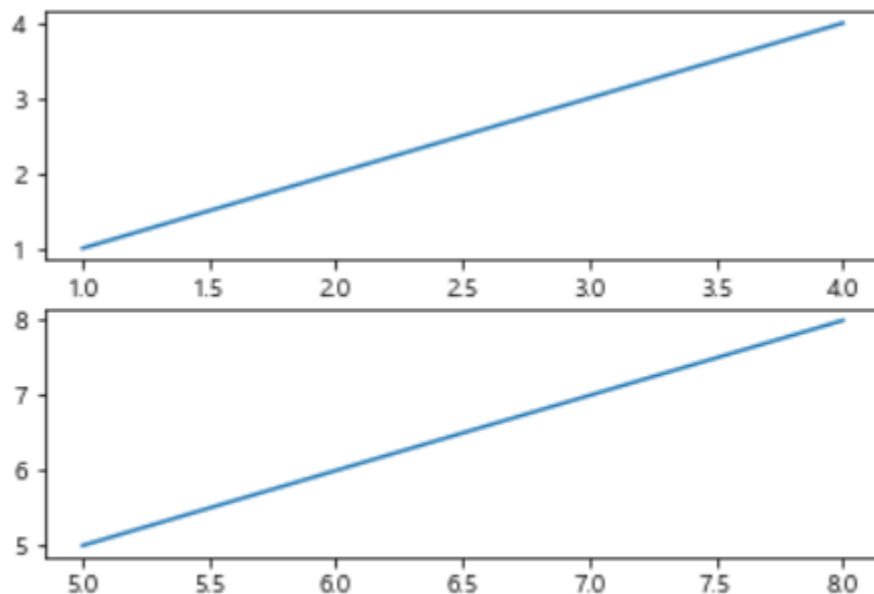
In [10]:

```
1 from matplotlib import font_manager, rc
2 import matplotlib
3 font_location="c:/Windows/fonts/malgunbd.ttf"
4 font_name = font_manager.FontProperties(fname=font_location).get_name()
5 matplotlib.rc('font', family=font_name)
6 plt.plot([1,2,3,4])
7 plt.xlabel("X축 한글표시")
8 plt.show()
```



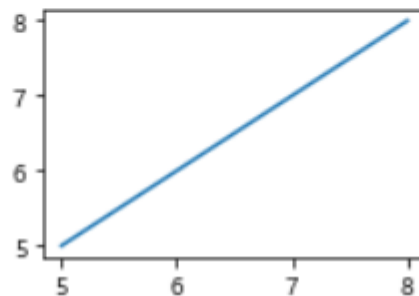
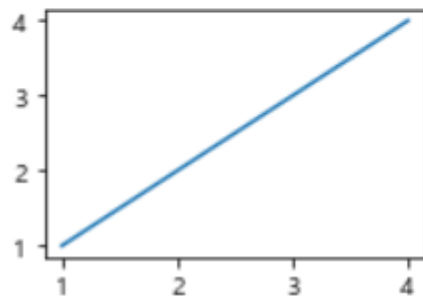
In [15]:

```
1 plt.figure() # 하나의 캔버스를 생성
2 # subplot(m,n,idx)
3 # 매트릭스 형태로 행2 열1개인 창을 의미. idx는 mn형태의 idx번째
4 plt.subplot(2,1,1)
5 plt.plot([1,2,3,4],[1,2,3,4])
6 plt.subplot(2,1,2)
7 plt.plot([5,6,7,8],[5,6,7,8])
8 plt.show()
```



In [18]:

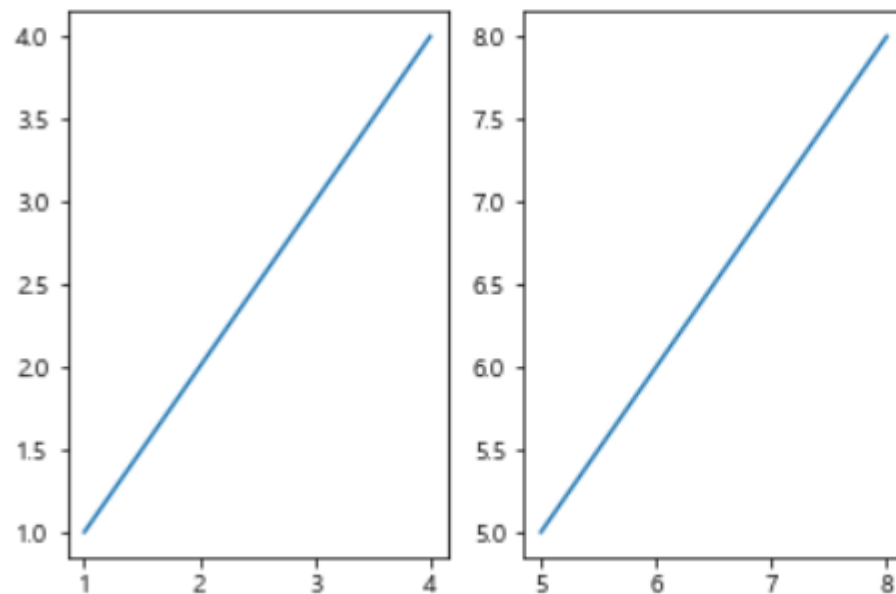
```
1 plt.figure()  
2 plt.subplot(2,2,1)  
3 plt.plot([1,2,3,4],[1,2,3,4])  
4 plt.subplot(2,2,2)  
5 plt.plot([5,6,7,8],[5,6,7,8])  
6 plt.show()
```



## 여러 개의 그래프 그리기

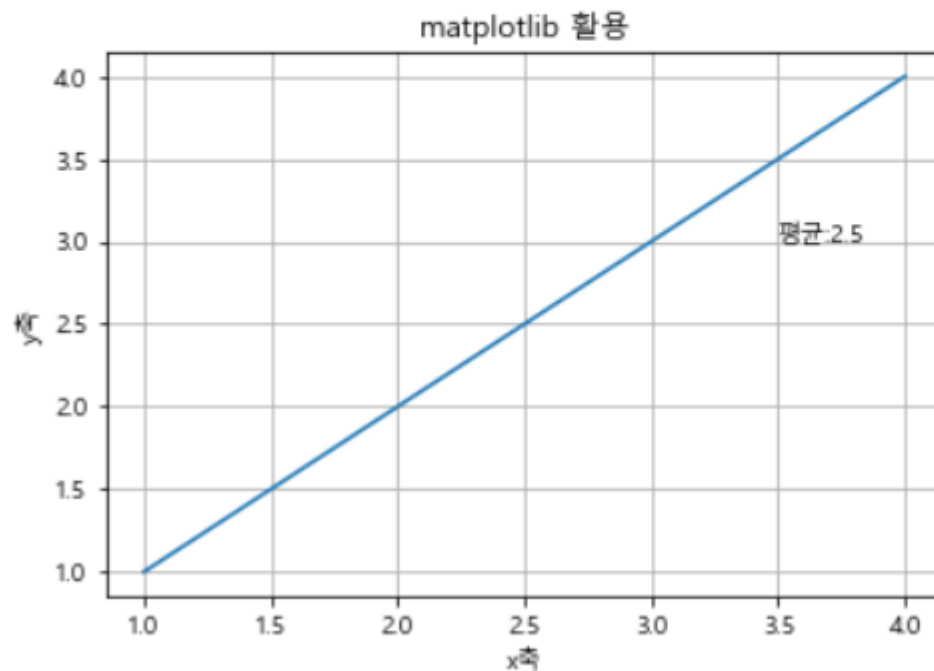
13

```
In [19]: 1 plt.figure()
          2 plt.subplot(1,2,1) # 1행의 첫 번째 컬럼
          3 plt.plot([1,2,3,4],[1,2,3,4])
          4 plt.subplot(1,2,2) # 1행의 두 번째 컬럼
          5 plt.plot([5,6,7,8],[5,6,7,8])
          6 plt.show()
```



In [22]:

```
1 plt.plot([1,2,3,4],[1,2,3,4])
2 plt.xlabel('x축')
3 plt.ylabel('y축')
4 plt.title('matplotlib 활용')
5 plt.text(3.5, 3.0, '평균:2.5')
6 plt.grid(True)
7 plt.show()
```

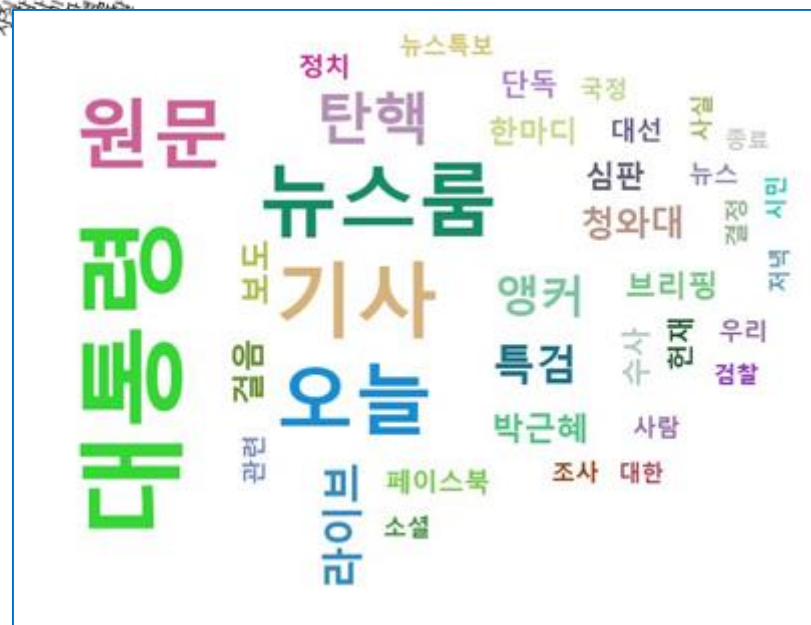
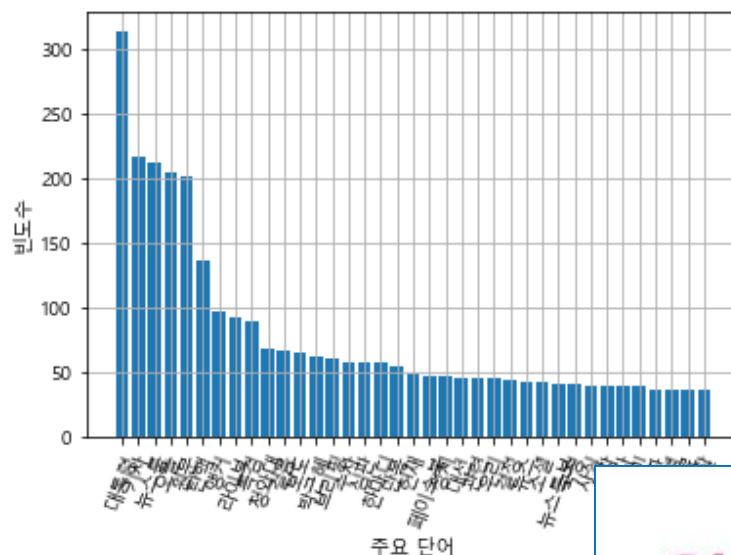


Matplotlib

**빈도분석**

k-최근접 이웃 알고리즘(k-Nearest Neighbors)

대통령 : 314  
 기사 : 217  
 뉴스룸 : 212  
 오늘 : 204  
 원문 : 202  
 탄핵 : 137  
 앵커 : 97  
 라이브 : 92  
 특검 : 90  
 청와대 : 69  
 결음 : 67  
 보도 : 65  
 박근혜 : 62  
 브리핑 : 61  
 수사 : 58  
 심판 : 58  
 한마디 : 57  
 단독 : 54  
 현재 : 49  
 페이스북 : 47  
 정치 : 47  
 대선 : 46  
 관련 : 45  
 우리 : 45  
 결정 : 44  
 뉴스 : 43  
 소셜 : 42  
 뉴스특보 : 41  
 국정 : 41  
 사실 : 40  
 사람 : 40  
 조사 : 40  
 시민 : 39



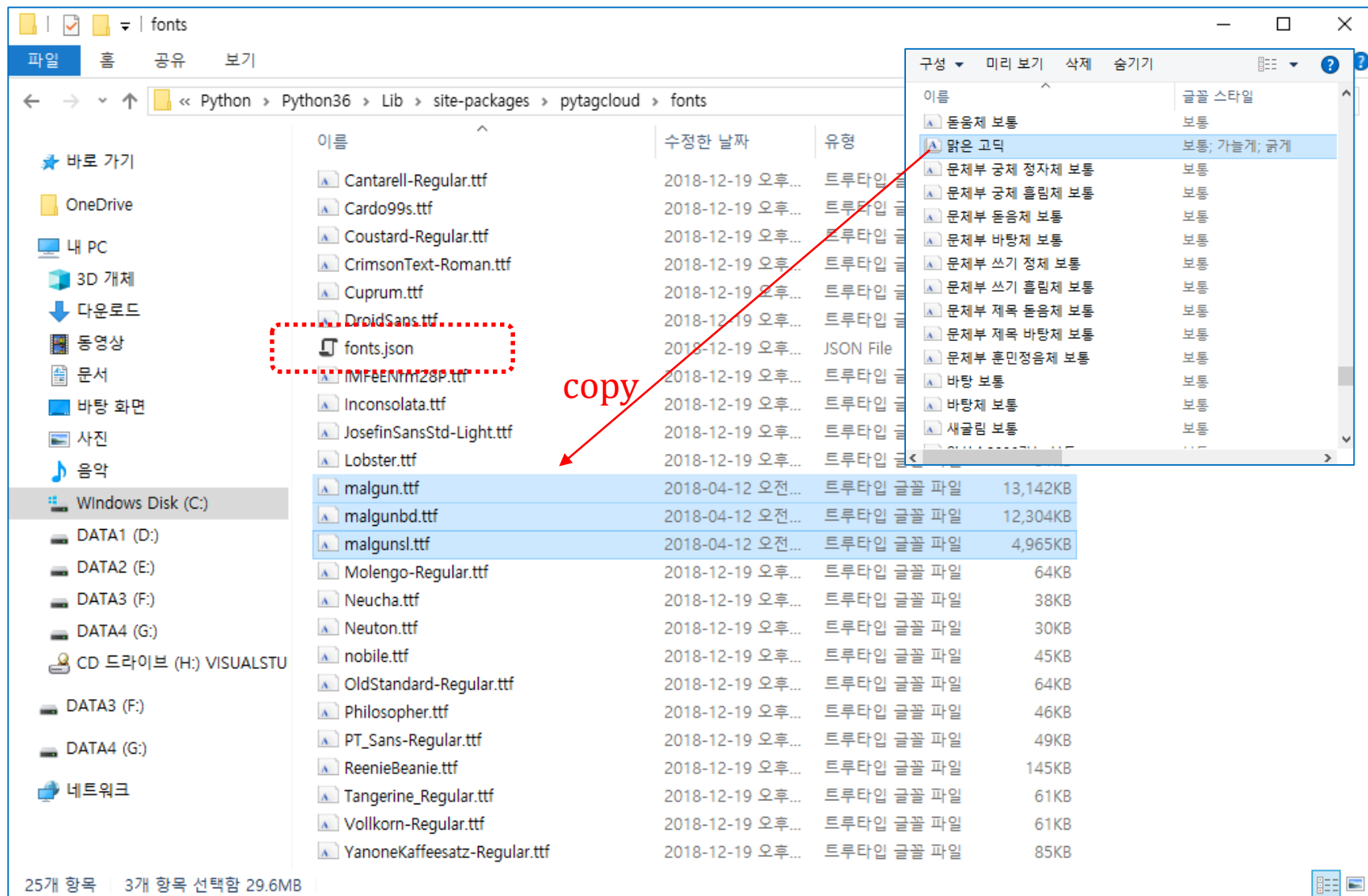


- 페이스북의 JTBC뉴스와 조선일보데이터 분석
  - 기간 : 2016-10-01~2017-03-12
- 설치 패키지
  - `pip install JPytype1`
  - `pip install KoNLPy`
  - `pip install pytagcloud`
  - `pip install pygame`
  - `pip install simplejson`

```
In [3]: 1 import json
        2 import re
        3
        4 from konlpy.tag import Okt
        5 from collections import Counter
        6
        7 import matplotlib.pyplot as plt
        8 import matplotlib
        9 from matplotlib import font_manager, rc
        10
        11 import pytagcloud
        12 import webbrowser
        13
        14 #[CODE 1]
        15 def showGraph(wordInfo):
        16
        17     font_location = "c:/Windows/fonts/malgunbd.ttf"
        18     font_name = font_manager.FontProperties(fname=font_location).get_name()
        19     matplotlib.rc('font', family=font_name)
        20
        21     plt.xlabel('주요 단어')
        22     plt.ylabel('빈도수')
        23     plt.grid(True)
        24     # 최대빈도수값과 최대빈도수 단어가 저장
        25     Sorted_Dict_Values = sorted(wordInfo.values(), reverse=True)
        26     Sorted_Dict_Keys = sorted(wordInfo, key=wordInfo.get, reverse=True)
        27     # 막대그래프 그리는 함수
        28     plt.bar(range(len(wordInfo)), Sorted_Dict_Values, align='center')
        29     # x축의 각 데이터 별 문자열(tick)을 지정
        30     plt.xticks(range(len(wordInfo)), list(Sorted_Dict_Keys), rotation='70')
        31
        32     plt.show()
```

```
In [4]: 1 #[CODE 2]
        2 def saveWordCloud(wordInfo, filename):
        3
        4     taglist = pytagcloud.make_tags(dict(wordInfo).items(), maxsize=80)
        5     pytagcloud.create_tag_image(taglist, filename, size=(640, 480), fontname='korean', rectangular=False)
        6     webbrowser.open(filename)
        7
```

```
In [11]: 1 def main():
2
3     openFileName = 'G:\python_workspace\jupyter\BigData\jtbcnews_facebook_2016-10-01_2017-03-12.json'
4     #openFileName = 'G:/python_workspace/jupyter/BigData/chosun_facebook_2016-10-01_2017-03-12.json'
5
6     cloudImagePath = openFileName + '.jpg'
7
8     rfile = open(openFileName, 'r', encoding='utf-8').read()
9
10    jsonData = json.loads(rfile)
11    message = ''
12
13    # jsonData의 개별 message를 합쳐서 하나의 문자열로 구성. 불필요한 \t, \n등의 문자 제거
14    for item in jsonData:
15        if 'message' in item.keys():
16            message = message + re.sub(r'[\t\n]', ' ', item['message']) + ' '
17
18    # 품사 클래스 - 명사만 추출하여 갯수를 세어 상위 50개만 가지고 온다.
19    nlp = Okt()
20    nouns = nlp.nouns(message)
21    count = Counter(nouns)
22
23    wordInfo = dict()
24    for tags, counts in count.most_common(50):
25        if (len(str(tags)) > 1):
26            wordInfo[tags] = counts
27            print ("%s : %d" % (tags, counts))
28
29    showGraph(wordInfo)
30    saveWordCloud(wordInfo, cloudImagePath)
31
32    if __name__ == "__main__":
33        main()
```



```
*C:\Python\Python36\Lib\site-packages\pytagcloud\fonts\fonts.json - Notepad++
파일(F) 편집(E) 찾기(S) 보기(V) 인코딩(N) 언어(L) 설정(T) 도구(O) 매크로 실행 플러그인 창 관리 ?
fonts.json
1 [
2   {
3     "name": "korean",
4     "ttf": "malgunbd.ttf",
5     "web": "http://fonts.googleapis.com/css?family=Nobile"
6   },
7   {
8     "name": "Old Standard TT",
9     "ttf": "OldStandard-Regular.ttf",
10    "web": "http://fonts.googleapis.com/css?family=Old+Standard+TT"
11  },
12  {
13    "name": "Cantarell",
14    "ttf": "Cantarell-Regular.ttf",
15    "web": "http://fonts.googleapis.com/css?family=Cantarell"
16  },
17  {
18    "name": "Reenie Beanie",
19    "ttf": "ReenieBeanie.ttf",
20    "web": "http://fonts.googleapis.com/css?family=Reenie+Beanie"
21  },
22  {
23    "name": "Cuprum",
24    "ttf": "Cuprum.ttf",
25    "web": "http://fonts.googleapis.com/css?family=Cuprum"
26  },
27  {
28    "name": "Molengo",
29    "ttf": "Molengo-Regular.ttf",
```

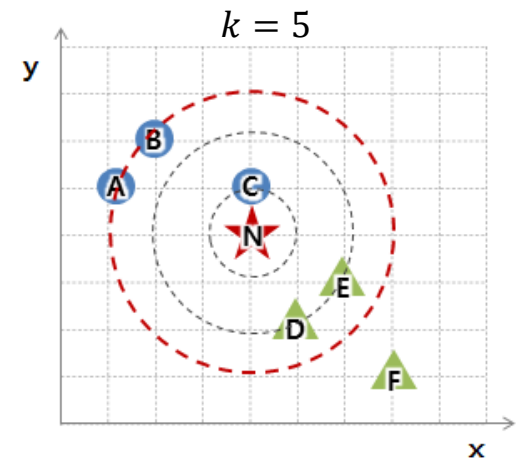
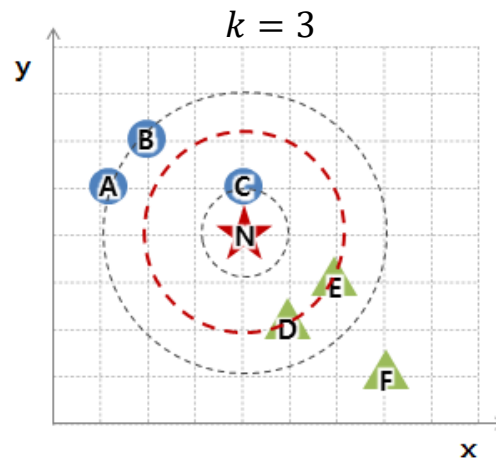
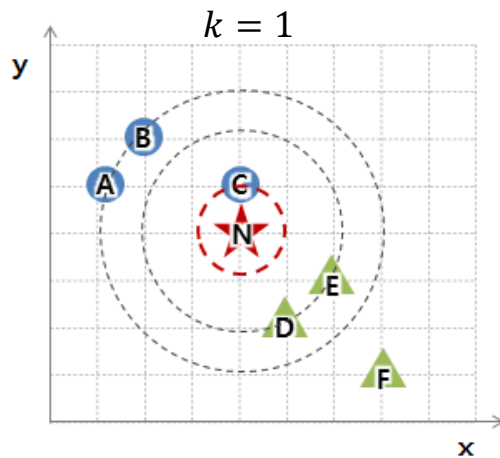
JSON file      length : 2,589    lines : 107    Ln : 4    Col : 1    Sel : 0 | 0    Unix (LF)    UTF-8    INS

Matplotlib

빈도분석

**k-최근접 이웃 알고리즘(k-Nearest Neighbors)**

- $k$ -NN( $k$  – *Nearest Neighbors*)은 머신러닝 분야에서 가장 간단한 알고리즘
- 훈련 데이터셋을 그냥 저장하는 것이 모델을 만드는 과정의 전부
- 새로운 데이터 포인트에 대해 예측할 땐 훈련 데이터셋에서 가장 가까운 데이터 포인트(최근접 이웃)을 찾는다.
- $k$ -NN은 새로 들어온 "★은 ○ (or △) 그룹의 데이터와 가장 가까우니 ★은 ○ (or △) 그룹이다." 라고 분류하는 알고리즘이다.

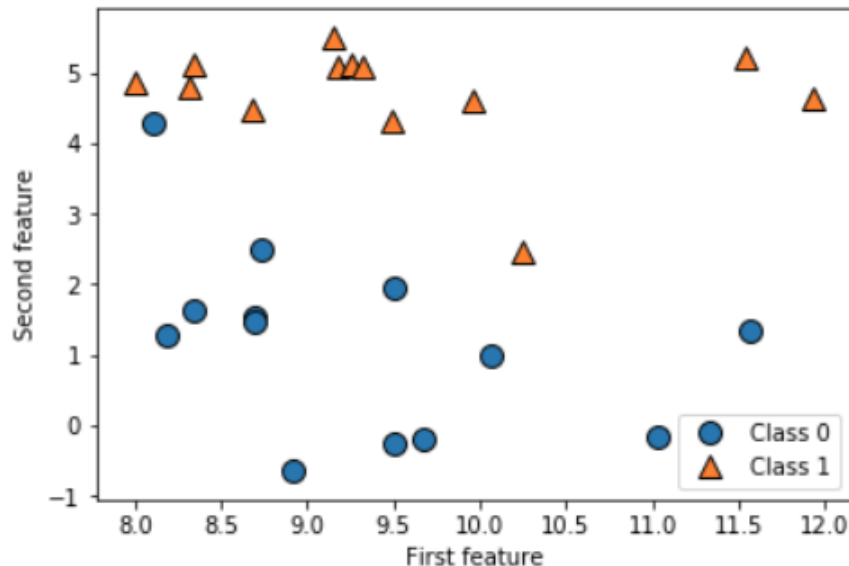




- k-최근접 이웃 알고리즘(k-Nearest Neighbors Algorithm)

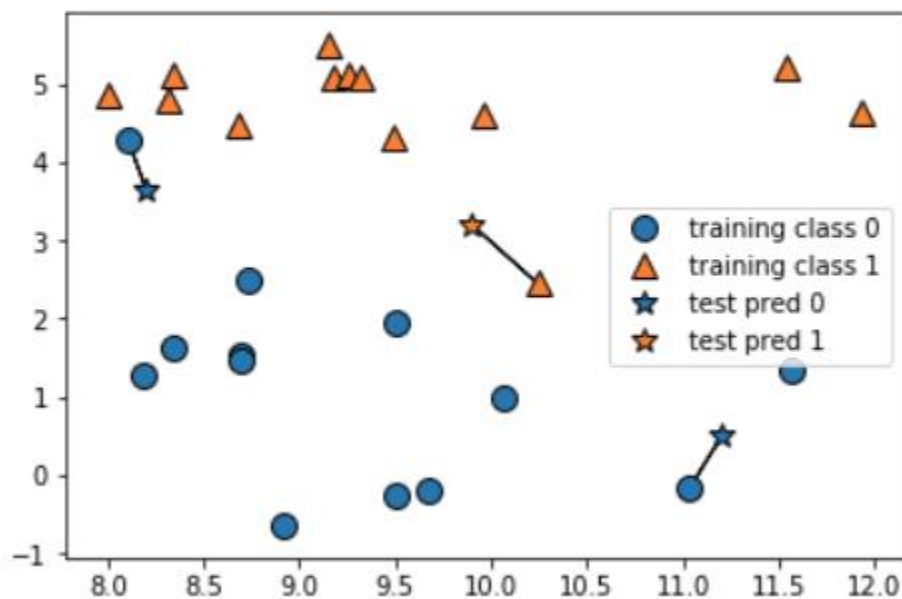
```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import mglearn as mglearn
4
5 # generate dataset
6 X, y = mglearn.datasets.make_forge()
7 # plot dataset
8 mglearn.discrete_scatter(X[:, 0], X[:, 1], y)
9 plt.legend(["Class 0", "Class 1"], loc=4)
10 plt.xlabel("First feature")
11 plt.ylabel("Second feature")
12 print("X.shape:", X.shape)
```

X.shape: (26, 2)



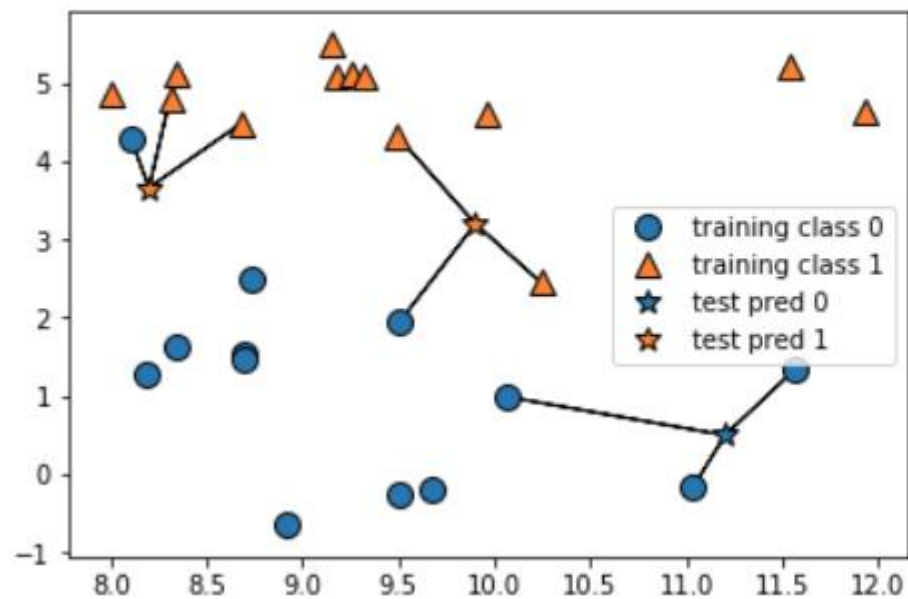
- 가장 가까운 훈련 데이터 포인트 하나를 최근접 이웃으로 찾아 예측에 사용
- 데이터 포인트 3개 추가: 1-최근접 이웃 알고리즘 예측

```
1 mglearn.plots.plot_knn_classification(n_neighbors=1)
```



- 3-최근접 이웃 알고리즘 예측

```
1 mglearn.plots.plot_knn_classification(n_neighbors=3)
```



- k-최근접 이웃 알고리즘(k-Nearest Neighbors Algorithm)

```
1 from sklearn.model_selection import train_test_split
2 X, y = mglearn.datasets.make_forge()
3
4 X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)
5 from sklearn.neighbors import KNeighborsClassifier
6 clf = KNeighborsClassifier(n_neighbors=3) # 이웃의 수 : 3
7
8 # 훈련 세트를 사용하여 분류 모델 학습
9 clf.fit(X_train, y_train)
```

```
KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
                     metric_params=None, n_jobs=None, n_neighbors=3, p=2,
                     weights='uniform')
```

```
1 # 테스트 데이터에 대해 predict 메서드를 호출해서 예측
2 print("Test set predictions:", clf.predict(X_test))
```

Test set predictions: [1 0 1 0 1 0 0]

```
1 # 모델이 얼마나 잘 일반화되었는지 평가
2 # 테스트 데이터와 테스트 레이블을 넣어 호출
3 print("Test set accuracy: {:.2f}".format(clf.score(X_test, y_test)))
```

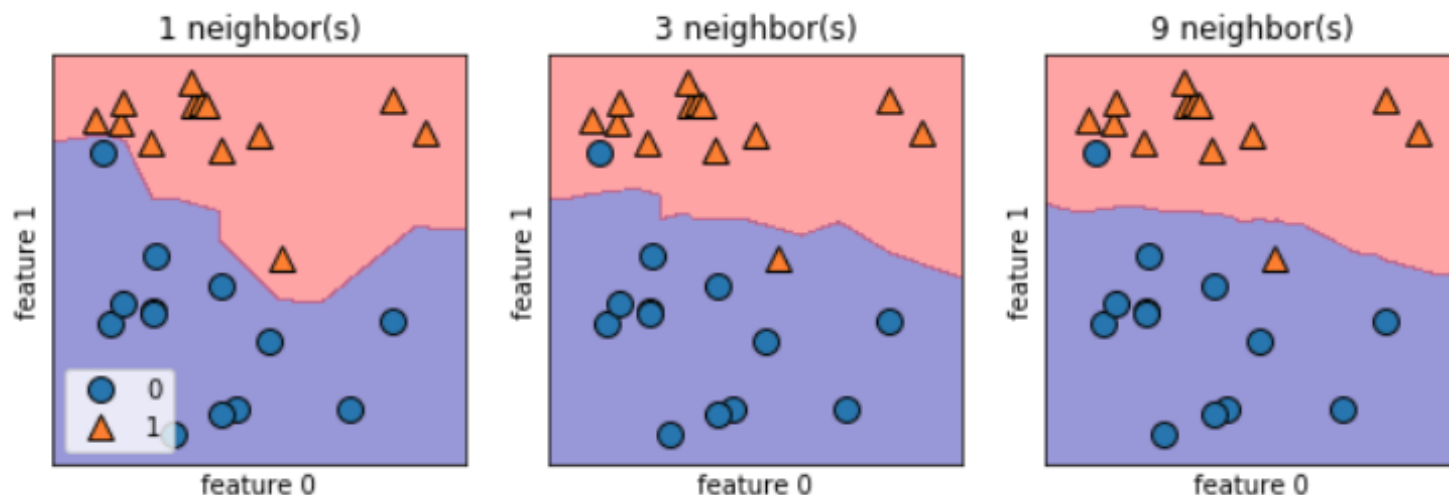
Test set accuracy: 0.86

```

1 fig, axes = plt.subplots(1, 3, figsize=(10, 3))
2
3 for n_neighbors, ax in zip([1, 3, 9], axes):
4     # the fit method returns the object self, so we can instantiate
5     # and fit in one line
6     clf = KNeighborsClassifier(n_neighbors=n_neighbors).fit(X, y)
7     mglearn.plots.plot_2d_separator(clf, X, fill=True, eps=0.5, ax=ax, alpha=.4)
8     mglearn.discrete_scatter(X[:, 0], X[:, 1], y, ax=ax)
9     ax.set_title("{} neighbor(s)".format(n_neighbors))
10    ax.set_xlabel("feature 0")
11    ax.set_ylabel("feature 1")
12 axes[0].legend(loc=3)

```

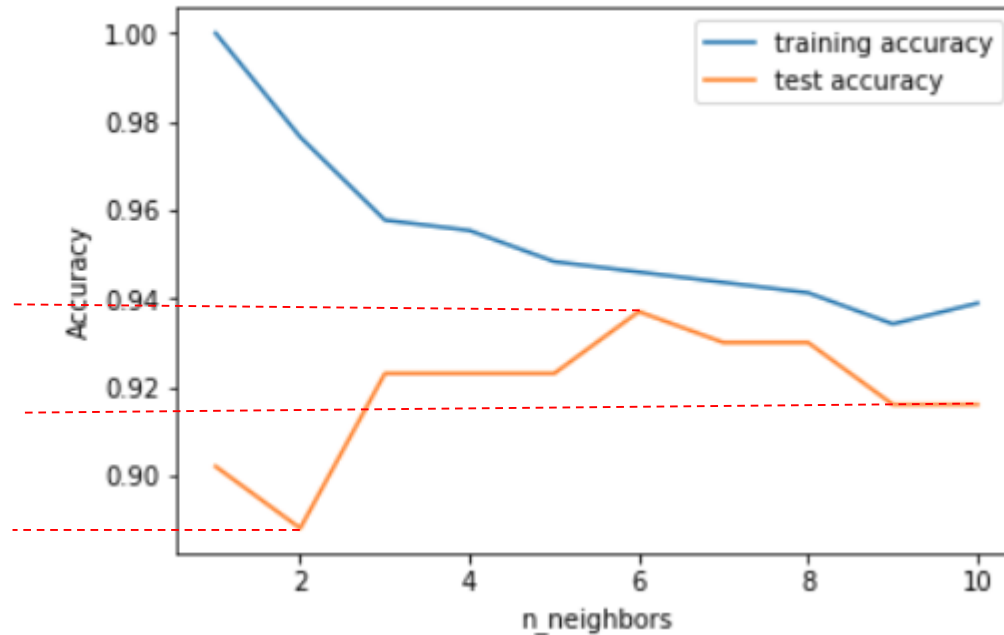
<matplotlib.legend.Legend at 0x117bf1ef160>



- 먼저 훈련 세트와 테스트 세트로 나눈 후 이웃의 수를 달리하여 훈련 세트와 테스트 세트의 성능 평가

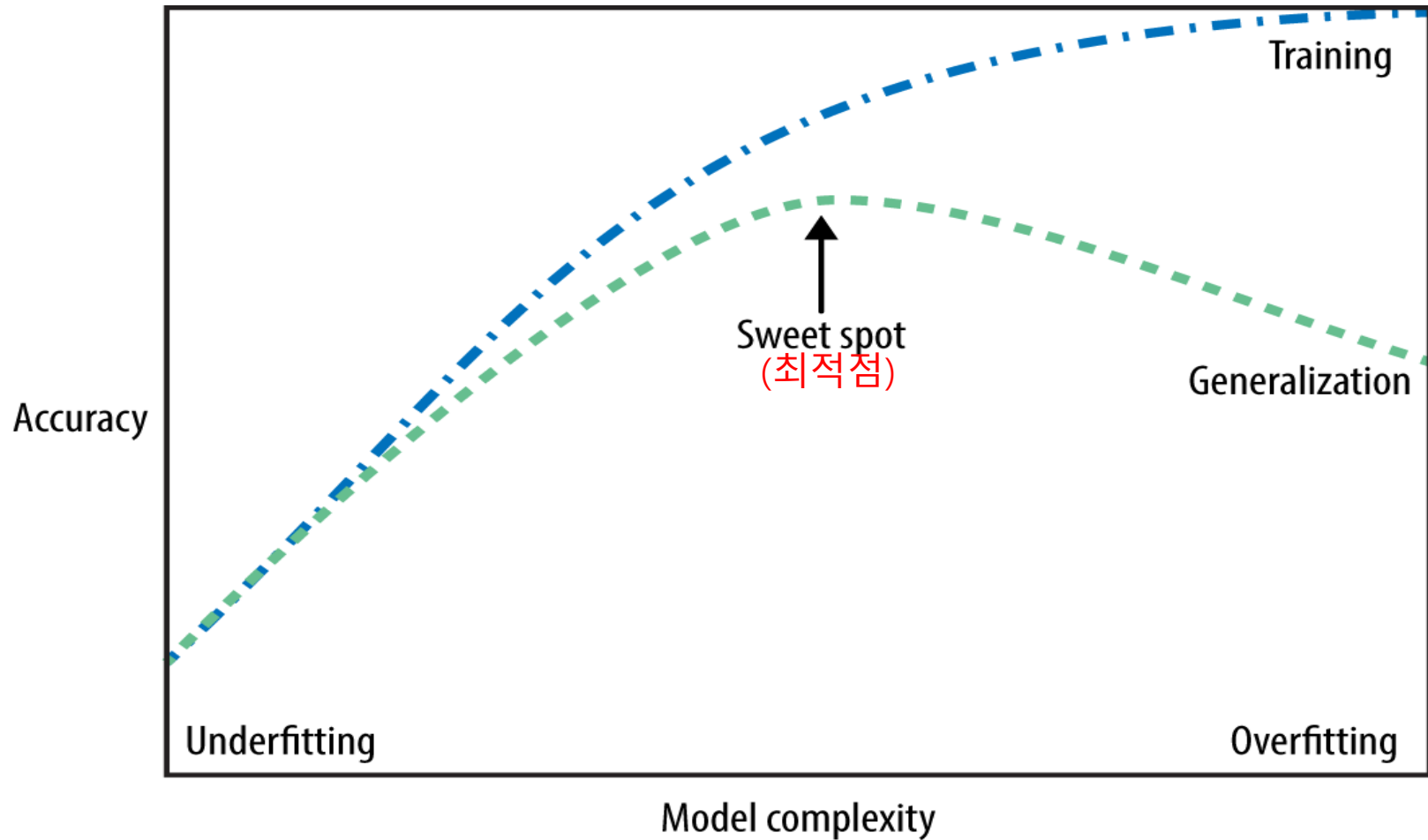
```
In [34]: 1 from sklearn.datasets import load_breast_cancer
2
3 cancer = load_breast_cancer()
4 X_train, X_test, y_train, y_test = train_test_split(
5     cancer.data, cancer.target, stratify=cancer.target, random_state=66)
6
7 training_accuracy = []
8 test_accuracy = []
9 # try n_neighbors from 1 to 10
10 neighbors_settings = range(1, 11)
11
12 for n_neighbors in neighbors_settings:
13     # build the model
14     clf = KNeighborsClassifier(n_neighbors=n_neighbors)
15     clf.fit(X_train, y_train)
16     # record training set accuracy
17     training_accuracy.append(clf.score(X_train, y_train))
18     # record generalization accuracy
19     test_accuracy.append(clf.score(X_test, y_test))
20
21 plt.plot(neighbors_settings, training_accuracy, label="training accuracy")
22 plt.plot(neighbors_settings, test_accuracy, label="test accuracy")
23 plt.ylabel("Accuracy")
24 plt.xlabel("n_neighbors")
25 plt.legend()
```

Out[34]: <matplotlib.legend.Legend at 0x117bf2ced68>



- 이웃을 하나 사용한 테스트 세트의 정확도는 이웃을 많이 사용했을 때보다 낮다.
- 그러나 10개 사용했을 때는 모델이 너무 단순해서 정확도는 더 나빠진다.
- 가장 좋을 때는 중간 정도인 여섯 개를 사용한 경우이다.
- 가장 나쁜 정확도도 88%이므로 수긍할 만하다.

- k-최근접 이웃 알고리즘(k-Nearest Neighbors Algorithm)





- 장점
  - 이해하기 매우 쉬운 모델
  - 많이 조정하지 않아도 자주 좋은 성능 발휘
  - 더 복잡한 알고리즘을 적용해 보기 전에 시도해 볼 수 있는 좋은 시작점
- 단점
  - 훈련세트가 매우 크면 예측이 느려진다.
  - (수백 개 이상의)많은 특성을 가진 데이터 셋에는 잘 동작하지 않으며,
  - 대부분이 0인(희소)데이터 셋은 특히 잘 작동되지 않는다.
  - 예측이 느리고 많은 특성을 처리하는 능력이 부족해 잘 사용하지 않는다.