

Data Science Fundamentals 5

Basic introduction on how to perform typical machine learning tasks with Python.

Prepared by Mykhailo Vladymyrov & Aris Marcolongo, Science IT Support, University Of Bern, 2020

This work is licensed under [CC0](https://creativecommons.org/share-your-work/public-domain/cc0/) (<https://creativecommons.org/share-your-work/public-domain/cc0/>).

Solutions to Part 1.

```
In [0]: from sklearn import linear_model

        from sklearn.datasets import make_blobs
        from sklearn.model_selection import train_test_split
        from sklearn import metrics

        from matplotlib import pyplot as plt
        import numpy as np
        import os
        from imageio import imread
        import pandas as pd
        from time import time as timer

        import tensorflow as tf

        %matplotlib inline
        from matplotlib import animation
        from IPython.display import HTML

In [2]: if not os.path.exists('data'):
        path = os.path.abspath('.') + '/colab_material.tgz'
        tf.keras.utils.get_file(path, 'https://github.com/newworldemancer/DSF5/raw/master/colab_material.tgz')
        !tar -xvzf colab_material.tgz > /dev/null 2>&1

Downloading data from https://github.com/newworldemancer/DSF5/raw/master/colab_material.tgz
98304/96847 [=====] - 0s 0us/step
```

Datasets

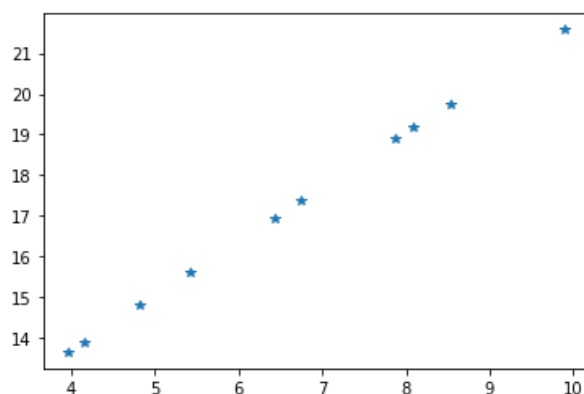
In this course we will use several synthetic and real-world datasets to illustrate the behavior of the models and exercise our skills.

1. Synthetic linear

```
In [0]: def get_linear(n_d=1, n_points=10, w=None, b=None, sigma=5):  
        x = np.random.uniform(0, 10, size=(n_points, n_d))  
  
        w = w or np.random.uniform(0.1, 10, n_d)  
        b = b or np.random.uniform(-10, 10)  
        y = np.dot(x, w) + b + np.random.normal(0, sigma, size=n_points)  
  
        print('true w =', w, ' ;  b =', b)  
  
        return x, y
```

```
In [4]: x, y = get_linear(n_d=1, sigma=0)  
        plt.plot(x[:, 0], y, '*')  
  
true w = [1.34032066] ;  b = 8.326857960042354
```

```
Out[4]: [<matplotlib.lines.Line2D at 0x7f08a22aef28>]
```

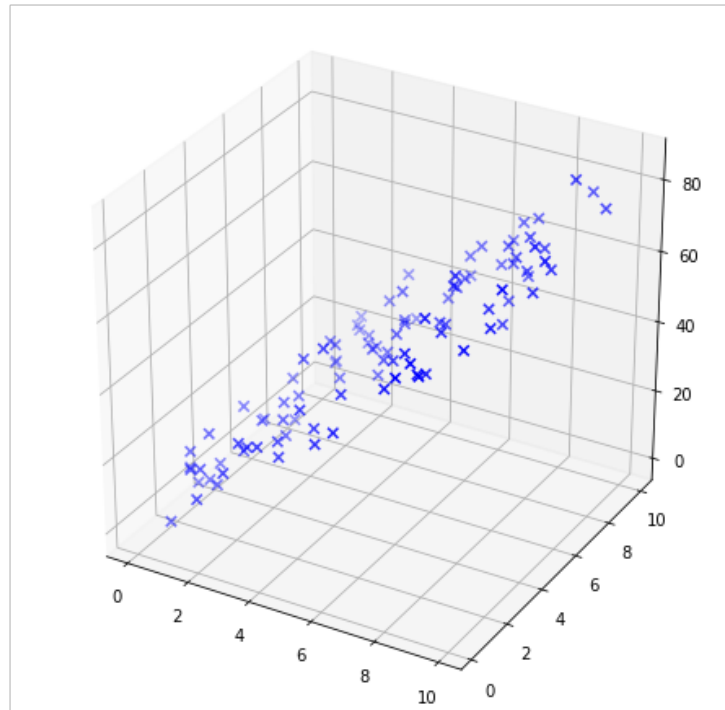


```
In [5]: n_d = 2
x, y = get_linear(n_d=n_d, n_points=100)

fig = plt.figure(figsize=(8,8))
ax = fig.add_subplot(111, projection='3d')
ax.scatter(x[:,0], x[:,1], y, marker='x', color='b', s=40)

true w = [7.03409766 0.83697333] ; b = 2.7928040906788194
```

Out[5]: <mpl_toolkits.mplot3d.art3d.Path3DCollection at 0x7f08a1dab198>



2. House prices

Subset of the the hous pricess kaggle dataset: <https://www.kaggle.com/c/house-prices-advanced-regression-techniques>
(<https://www.kaggle.com/c/house-prices-advanced-regression-techniques>)

```

In [0]: def house_prices_dataset(return_df=False, price_max=400000, area_max=400
00):
    path = 'data/train.csv'

    df = pd.read_csv(path, na_values="NaN", keep_default_na=False)

    useful_fields = ['LotArea',
                     'Utilities', 'OverallQual', 'OverallCond',
                     'YearBuilt', 'YearRemodAdd', 'ExterQual', 'ExterCond',
                     'HeatingQC', 'CentralAir', 'Electrical',
                     '1stFlrSF', '2ndFlrSF', 'GrLivArea',
                     'FullBath', 'HalfBath',
                     'BedroomAbvGr', 'KitchenAbvGr', 'KitchenQual', 'TotRms
AbvGrd',
                     'Functional', 'PoolArea',
                     'YrSold', 'MoSold'
                    ]
    target_field = 'SalePrice'

    cleanup_nums = {"Street":      {"GrvL": 0, "Pave": 1},
                    "LotFrontage": {"NA":0},
                    "Alley":       {"NA":0, "GrvL": 1, "Pave": 2},
                    "LotShape":    {"IR3":0, "IR2": 1, "IR1": 2, "Reg":3},
                    "Utilities":   {"ELO":0, "NoSeWa": 1, "NoSewr": 2, "Al
lPub": 3},
                    "LandSlope":   {"Sev":0, "Mod": 1, "Gtl": 3},
                    "ExterQual":   {"Po":0, "Fa": 1, "TA": 2, "Gd": 3, "E
x":4},
                    "ExterCond":   {"Po":0, "Fa": 1, "TA": 2, "Gd": 3, "E
x":4},
                    "BsmtQual":    {"NA":0, "Po":1, "Fa": 2, "TA": 3, "G
d": 4, "Ex":5},
                    "BsmtCond":    {"NA":0, "Po":1, "Fa": 2, "TA": 3, "G
d": 4, "Ex":5},
                    "BsmtExposure":{"NA":0, "No":1, "Mn": 2, "Av": 3, "G
d": 4},
                    "BsmtFinType1":{"NA":0, "Unf":1, "LwQ": 2, "Rec": 3, "
BLQ": 4, "ALQ":5, "GLQ":6},
                    "BsmtFinType2":{"NA":0, "Unf":1, "LwQ": 2, "Rec": 3, "
BLQ": 4, "ALQ":5, "GLQ":6},
                    "HeatingQC":   {"Po":0, "Fa": 1, "TA": 2, "Gd": 3, "E
x":4},
                    "CentralAir":  {"N":0, "Y": 1},
                    "Electrical":  {"NA":0, "Mix":1, "FuseP":2, "FuseF":
3, "FuseA": 4, "SBrkr": 5},
                    "KitchenQual": {"Po":0, "Fa": 1, "TA": 2, "Gd": 3, "E
x":4},
                    "Functional":  {"Sal":0, "Sev":1, "Maj2": 2, "Maj1":
3, "Mod": 4, "Min2":5, "Min1":6, 'Typ':7},
                    "FireplaceQu": {"NA":0, "Po":1, "Fa": 2, "TA": 3, "G
d": 4, "Ex":5},
                    "PoolQC":      {"NA":0, "Fa": 1, "TA": 2, "Gd": 3, "E
x":4},
                    "Fence":       {"NA":0, "MnWw": 1, "GdWo": 2, "MnPrv":
3, "GdPrv":4},
                    }

    df_X = df[useful_fields].copy()
    df_X.replace(cleanup_nums, inplace=True) # convert continous categori
al variables to numerical
    df_Y = df[target_field].copy()

    x = df_X.to_numpy().astype(np.float32)
    y = df_Y.to_numpy().astype(np.float32)

    if price_max>0:
        idxs = y<price_max
        x = x[idxs]

```

```
In [7]: x, y, df = house_prices_dataset(return_df=True)
print(x.shape, y.shape)
df.head()
```

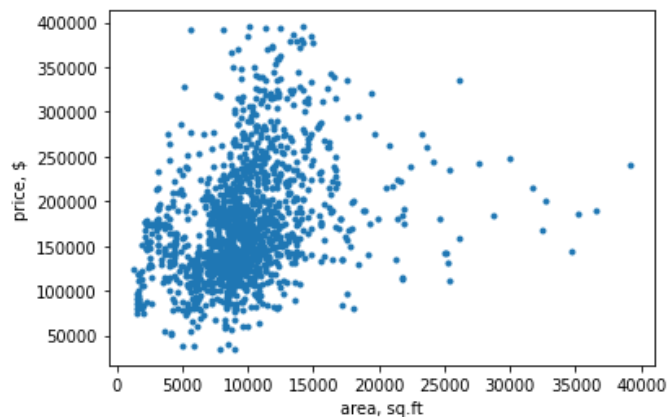
```
(1420, 24) (1420,)
```

```
Out[7]:
```

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour	Util
0	1	60	RL	65	8450	Pave	NA	Reg	Lvl	AllF
1	2	20	RL	80	9600	Pave	NA	Reg	Lvl	AllF
2	3	60	RL	68	11250	Pave	NA	IR1	Lvl	AllF
3	4	70	RL	60	9550	Pave	NA	IR1	Lvl	AllF
4	5	60	RL	84	14260	Pave	NA	IR1	Lvl	AllF

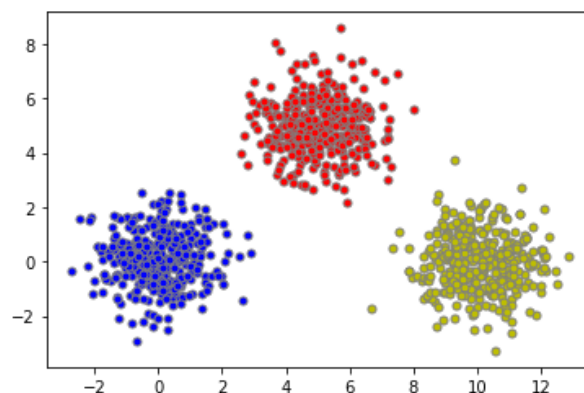
```
5 rows × 81 columns
```

```
In [8]: plt.plot(x[:, 0], y, '.')
plt.xlabel('area, sq.ft')
plt.ylabel('price, $');
```



3. Blobs

```
In [9]: x, y = make_blobs(n_samples=1000, centers=[[0,0], [5,5], [10, 0]])
colors = "bry"
for i, color in enumerate(colors):
    idx = y == i
    plt.scatter(x[idx, 0], x[idx, 1], c=color, edgecolor='gray', s=25)
```



4. Fashion MNIST

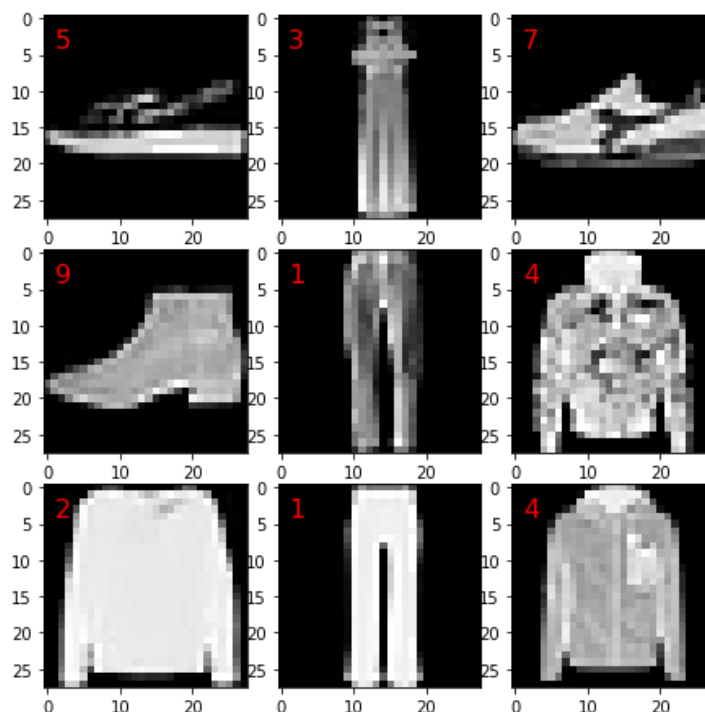
Fashion-MNIST is a dataset of Zalando's article images—consisting of a training set of 60,000 examples and a test set of 10,000 examples. Each example is a 28x28 grayscale image, associated with a label from 10 classes. (from <https://github.com/zalandoresearch/fashion-mnist> (<https://github.com/zalandoresearch/fashion-mnist>))

```
In [10]: fashion_mnist = tf.keras.datasets.fashion_mnist
(train_images, train_labels), (test_images, test_labels) = fashion_mnist.load_data()
```

```
Downloading data from https://storage.googleapis.com/tensorflow/tf-keras-
datasets/train-labels-idx1-ubyte.gz
32768/29515 [=====] - 0s 0us/step
Downloading data from https://storage.googleapis.com/tensorflow/tf-keras-
datasets/train-images-idx3-ubyte.gz
26427392/26421880 [=====] - 0s 0us/step
Downloading data from https://storage.googleapis.com/tensorflow/tf-keras-
datasets/t10k-labels-idx1-ubyte.gz
8192/5148 [=====] - 0s 0us/step
Downloading data from https://storage.googleapis.com/tensorflow/tf-keras-
datasets/t10k-images-idx3-ubyte.gz
4423680/4422102 [=====] - 0s 0us/step
```

Let's check few samples:

```
In [11]: n = 3
fig, ax = plt.subplots(n, n, figsize=(2*n, 2*n))
ax = [ax_xy for ax_xy in ax for ax_xy in ax_xy]
for axi, im_idx in zip(ax, np.random.choice(len(train_images), n**2)):
    im = train_images[im_idx]
    im_class = train_labels[im_idx]
    axi.imshow(im, cmap='gray')
    axi.text(1, 4, f'{im_class}', color='r', size=16)
plt.tight_layout(0,0,0)
```



Each training and test example is assigned to one of the following labels:

Label	Description
0	T-shirt/top
1	Trouser
2	Pullover
3	Dress
4	Coat
5	Sandal
6	Shirt
7	Sneaker
8	Bag
9	Ankle boot

EXERCISE 1.

```

In [12]: # Solution:
x, y = house_prices_dataset()

# 1. make train/test split
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)

# 2. fit the model
reg = linear_model.LinearRegression()
reg.fit(x_train, y_train)

# 3. evaluate MSE, MAD, and R2 on train and test datasets
#prediction:
y_p_train = reg.predict(x_train)
y_p_test = reg.predict(x_test)

# mse
print('train mse =', np.std(y_train - y_p_train))
print('test mse =', np.std(y_test - y_p_test))
# mae
print('train mae =', np.mean(np.abs(y_train - y_p_train)))
print('test mae =', np.mean(np.abs(y_test - y_p_test)))
# R2
print('train R2 =', reg.score(x_train, y_train))
print('test R2 =', reg.score(x_test, y_test))

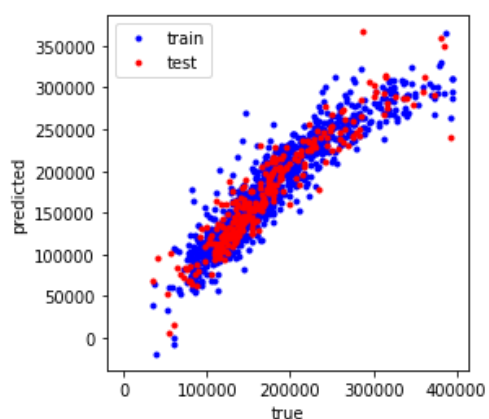
# 4. plot y vs predicted y for test and train parts
plt.plot(y_train, y_p_train, 'b.', label='train')
plt.plot(y_test, y_p_test, 'r.', label='test')

plt.plot([0], [0], 'w.') # dummy to have origin
plt.xlabel('true')
plt.ylabel('predicted')
plt.gca().set_aspect('equal')
plt.legend()

train mse = 24834.885
test mse = 24293.932
train mae = 18221.04
test mae = 17279.08
train R2 = 0.8534532342221349
test R2 = 0.8706671727120681

```

Out[12]: <matplotlib.legend.Legend at 0x7f089d87e278>



EXERCISE 2.


```
In [0]: fashion_mnist = tf.keras.datasets.fashion_mnist
(train_images, train_labels), (test_images, test_labels) = fashion_mnist.load_data()
```

We will reshape 2-d images to 1-d arrays for use in scikit-learn:

```
In [0]: n_train = len(train_labels)
x_train = train_images.reshape((n_train, -1))
y_train = train_labels

n_test = len(test_labels)
x_test = test_images.reshape((n_test, -1))
y_test = test_labels
```

Now use a multinomial logistic regression classifier, and measure the accuracy:

```
In [15]: #solution
# 1. Create classifier
multi_class = 'multinomial'
clf = linear_model.LogisticRegression(solver='sag', max_iter=20,
                                     multi_class=multi_class)

# 2. fit the model
t1 = timer()
clf.fit(x_train, y_train)
t2 = timer()
print ('training time: %.1fs'%(t2-t1))

# 3. evaluate accuracy on train and test datasets
print("training score : %.3f" % (clf.score(x_train, y_train)))
print("test score : %.3f" % (clf.score(x_test, y_test)))

/usr/local/lib/python3.6/dist-packages/sklearn/linear_model/_sag.py:330:
ConvergenceWarning: The max_iter was reached which means the coef_ did not
converge
  "the coef_ did not converge", ConvergenceWarning)

training time: 40.6s
training score : 0.874
test score : 0.843
```