

# LetsFaceIt\_S06Ex1

December 14, 2021

## Solution for Sheet 06 - Exercise 1 (slightly adapted based on Hani's comments)

Contributors:

- Annika Bätz: annika.baetz@uni-potsdam.de
- Max Serra: maximilia-manuel.serra.lasierra@uni-potsdam.de
- Adrián de Miguel: adrian.de.miguel.palacio@uni-potsdam.de
- Ignacio Llorca: llorcarodriguez@uni-potsdam.de

```
[ ]: import numpy as np
import pandas as pd
from numpy.linalg import inv
from sklearn.model_selection import train_test_split

import matplotlib.pyplot as plt
plt.style.use('seaborn')
```

### 0.0.1 Load data

```
[ ]: X_potential_drivers = np.array(pd.read_csv('X_potential_drivers.txt', ↴
                                               header=None))
Y_snow = np.array(pd.read_csv('Y_snow.txt', header=None))
```

## 1 1 Preprocessing

Target analysis ( $Y_{snow}$ ):

```
[ ]: df_Y_snow = pd.DataFrame(Y_snow)

print("- First 5 target values:")
print(df_Y_snow.head())
print("\n- Type:")
print(df_Y_snow.dtypes)
print("\n- Main statistics:")
print(df_Y_snow.describe())
```

- First 5 target values:

0

0 44.747

1 43.015

```

2 43.397
3 41.650
4 41.919

- Type:
0    float64
dtype: object

- Main statistics:
   0
count 201.000000
mean -13.905716
std 33.601792
min -70.328000
25% -44.056000
50% -13.604000
75% 16.012000
max 44.747000

```

#### Features analysis ( $x_1, x_2, x_3, x_4, x_5$ ):

```
[ ]: feature_names = ['1.feature', '2.feature', '3.feature', '4.feature', '5.feature']
df_X_potential_drivers = pd.read_csv('X_potential_drivers.txt', ↵
                                         names=feature_names)

print("- First 5 realisations:")
print(df_X_potential_drivers.head())
print("\n- Type:")
print(df_X_potential_drivers.dtypes)
print("\n- Main statistics:")
print(df_X_potential_drivers.describe())
```

- First 5 realisations:

	1.feature	2.feature	3.feature	4.feature	5.feature
0	11.728	-0.64836	-5.00	2.3026+3.1416i	-23.0
1	11.771	-0.51455	-4.95	2.2925+3.1416i	-22.7
2	11.817	-0.39388	-4.90	2.2824+3.1416i	-22.4
3	11.864	-0.28239	-4.85	2.2721+3.1416i	-22.1
4	11.913	-0.17704	-4.80	2.2618+3.1416i	-21.8

- Type:

1.feature	float64
2.feature	float64
3.feature	float64
4.feature	object
5.feature	float64

dtype: object

- Main statistics:

	1.feature	2.feature	3.feature	5.feature
count	201.000000	2.010000e+02	201.000000	201.000000
mean	12.000000	-7.070077e-17	0.000000	7.000000
std	0.345947	9.658693e+00	2.908393	17.450358
min	11.500000	-8.071300e+01	-5.000000	-23.000000
25%	11.669000	-8.855600e-01	-2.500000	-8.000000
50%	12.000000	0.000000e+00	0.000000	7.000000
75%	12.331000	8.855600e-01	2.500000	22.000000
max	12.500000	8.071300e+01	5.000000	37.000000

### 1.0.1 Deal with complex values

There might be many different alternatives to the one we have followed below to deal with complex values, for instance using Complex Matrices. This approach has also been covered in this topic and whose summary can be found in the section [Annex: Complex Matrices](#)

**Our idea here:** Split the fourth feature into a real and imaginary part. Additionally, we can compute the absolute value of the complex value and can use it as a third additional feature.

```
[ ]: # lists to save our extracted features:
complex_feature=[] # original values
real_feature=[] # real part
imag_feature=[] # imaginary part
abs_feature=[] # absolute value

for s in X_potential_drivers[:,3]:
    s = s.replace('i', 'j')
    c = complex(s)
    complex_feature.append(c)
    real_feature.append(np.real(c))
    imag_feature.append(np.imag(c))
    abs_feature.append(np.abs(c))

[ ]: print('complex feature:')
print(complex_feature[:5])
print('...')
print(complex_feature[-5:])

complex feature:
[(2.3026+3.1416j), (2.2925+3.1416j), (2.2824+3.1416j), (2.2721+3.1416j),
(2.2618+3.1416j)]
...
[(2.2618+0j), (2.2721+0j), (2.2824+0j), (2.2925+0j), (2.3026+0j)]

[ ]: print('real feature:')
print(real_feature[:5])
print('...')
print(real_feature[-5:])
```

```

real feature:
[2.3026, 2.2925, 2.2824, 2.2721, 2.2618]
...
[2.2618, 2.2721, 2.2824, 2.2925, 2.3026]

[ ]: print('imaginary feature:')
print(imag_feature[:5])
print('...')
print(imag_feature[-5:])

imaginary feature:
[3.1416, 3.1416, 3.1416, 3.1416, 3.1416]
...
[0.0, 0.0, 0.0, 0.0, 0.0]

[ ]: print('absolute feature:')
print(abs_feature[:5])
print('...')
print(abs_feature[-5:])

absolute feature:
[3.8950760352013667, 3.8891138849357447, 3.883168850307697, 3.877123801221725,
3.8710967179857443]
...
[2.2618, 2.2721, 2.2824, 2.2925, 2.3026]

[ ]: np.unique(imag_feature)

[ ]: array([0.      , 3.1416])

```

Because the imaginary feature only takes on two unique values, i.e. it is a binary feature, we could also reformulate it as: \* 0: no imaginary part \* 1: imaginary part (instead of 3.1416) But we will leave it with the original implementation for now.

## 1.0.2 Dealing with infinity

We have one instance were the value of the original 4.feature is -inf.

```

[ ]: np.unique(real_feature)

[ ]: array([-inf, -2.3026 , -1.6094 , -1.204 , -0.91629, -0.69315,
-0.51083, -0.35667, -0.22314, -0.10536, 0.      , 0.09531,
0.18232, 0.26236, 0.33647, 0.40547, 0.47    , 0.53063,
0.58779, 0.64185, 0.69315, 0.74194, 0.78846, 0.83291,
0.87547, 0.91629, 0.95551, 0.99325, 1.0296 , 1.0647 ,
1.0986 , 1.1314 , 1.1632 , 1.1939 , 1.2238 , 1.2528 ,
1.2809 , 1.3083 , 1.335  , 1.361  , 1.3863 , 1.411  ,
1.4351 , 1.4586 , 1.4816 , 1.5041 , 1.5261 , 1.5476 ,
1.5686 , 1.5892 , 1.6094 , 1.6292 , 1.6487 , 1.6677 ,
1.6864 , 1.7047 , 1.7228 , 1.7405 , 1.7579 , 1.775  ,
1.7932 , 1.8105 , 1.8278 , 1.8451 , 1.8624 , 1.8797 ,
1.897  , 1.9143 , 1.9316 , 1.9489 , 1.9662 , 1.9835 ,
1.9998 , 2.0171 , 2.0344 , 2.0517 , 2.069  , 2.0863 ,
2.1033 , 2.1206 , 2.1379 , 2.1552 , 2.1725 , 2.1898 ,
2.2071 , 2.2244 , 2.2417 , 2.259  , 2.2763 , 2.2936 ,
2.3109 , 2.3282 , 2.3455 , 2.3628 , 2.3801 , 2.3974 ,
2.4147 , 2.432  , 2.4493 , 2.4666 , 2.4839 , 2.5012 ,
2.5185 , 2.5358 , 2.5531 , 2.5704 , 2.5877 , 2.605  ,
2.6223 , 2.6396 , 2.6569 , 2.6742 , 2.6915 , 2.7088 ,
2.7261 , 2.7434 , 2.7607 , 2.778  , 2.7953 , 2.8126 ,
2.8299 , 2.8472 , 2.8645 , 2.8818 , 2.8991 , 2.9164 ,
2.9337 , 2.951  , 2.9683 , 2.9856 , 3.0029 , 3.0202 ,
3.0375 , 3.0548 , 3.0721 , 3.0894 , 3.1067 , 3.124  ,
3.1413 , 3.1586 , 3.1759 , 3.1932 , 3.2105 , 3.2278 ,
3.2451 , 3.2624 , 3.2797 , 3.297  , 3.3143 , 3.3316 ,
3.3489 , 3.3662 , 3.3835 , 3.4008 , 3.4181 , 3.4354 ,
3.4527 , 3.47  , 3.4873 , 3.5046 , 3.5219 , 3.5392 ,
3.5565 , 3.5738 , 3.5911 , 3.6084 , 3.6257 , 3.643  ,
3.6603 , 3.6776 , 3.6949 , 3.7122 , 3.7295 , 3.7468 ,
3.7641 , 3.7814 , 3.7987 , 3.816  , 3.8333 , 3.8506 ,
3.8679 , 3.8852 , 3.9025 , 3.92  , 3.9373 , 3.9546 ,
3.9719 , 3.9892 , 3.9998 , 4.0105 , 4.0212 , 4.0319 ,
4.0426 , 4.0533 , 4.064  , 4.0747 , 4.0854 , 4.0961 ,
4.1068 , 4.1175 , 4.1282 , 4.1389 , 4.1496 , 4.1603 ,
4.171  , 4.1817 , 4.1924 , 4.2031 , 4.2138 , 4.2245 ,
4.2352 , 4.2459 , 4.2566 , 4.2673 , 4.278  , 4.2887 ,
4.2994 , 4.3101 , 4.3208 , 4.3315 , 4.3422 , 4.3529 ,
4.3636 , 4.3743 , 4.385  , 4.3957 , 4.4064 , 4.4171 ,
4.4278 , 4.4385 , 4.4492 , 4.4599 , 4.4706 , 4.4813 ,
4.492  , 4.5027 , 4.5134 , 4.5241 , 4.5348 , 4.5455 ,
4.5562 , 4.5669 , 4.5776 , 4.5883 , 4.599  , 4.6097 ,
4.6204 , 4.6311 , 4.6418 , 4.6525 , 4.6632 , 4.6739 ,
4.6846 , 4.6953 , 4.706  , 4.7167 , 4.7274 , 4.7381 ,
4.7488 , 4.7595 , 4.7602 , 4.7709 , 4.7816 , 4.7923 ,
4.793  , 4.7937 , 4.7944 , 4.7951 , 4.7958 , 4.7965 ,
4.7972 , 4.7979 , 4.7986 , 4.7993 , 4.7998 , 4.8005 ,
4.8012 , 4.8019 , 4.8026 , 4.8033 , 4.804  , 4.8047 ,
4.8054 , 4.8061 , 4.8068 , 4.8075 , 4.8082 , 4.8089 ,
4.8096 , 4.8103 , 4.811  , 4.8117 , 4.8124 , 4.8131 ,
4.8138 , 4.8145 , 4.8152 , 4.8159 , 4.8166 , 4.8173 ,
4.818  , 4.8187 , 4.8194 , 4.8197 , 4.8198 , 4.8199 ,
4.82  , 4.8201 , 4.8202 , 4.8203 , 4.8204 , 4.8205 ,
4.8206 , 4.8207 , 4.8208 , 4.8209 , 4.821  , 4.8211 ,
4.8212 , 4.8213 , 4.8214 , 4.8215 , 4.8216 , 4.8217 ,
4.8218 , 4.8219 , 4.822  , 4.8221 , 4.8222 , 4.8223 ,
4.8224 , 4.8225 , 4.8226 , 4.8227 , 4.8228 , 4.8229 ,
4.8229 , 4.823  , 4.8231 , 4.8232 , 4.8233 , 4.8234 ,
4.8235 , 4.8236 , 4.8237 , 4.8238 , 4.8239 , 4.824  ,
4.8241 , 4.8242 , 4.8243 , 4.8244 , 4.8245 , 4.8246 ,
4.8247 , 4.8248 , 4.8249 , 4.825  , 4.8251 , 4.8252 ,
4.8253 , 4.8254 , 4.8255 , 4.8256 , 4.8257 , 4.8258 ,
4.8259 , 4.826  , 4.8261 , 4.8262 , 4.8263 , 4.8264 ,
4.8265 , 4.8266 , 4.8267 , 4.8268 , 4.8269 , 4.827  ,
4.8271 , 4.8272 , 4.8273 , 4.8274 , 4.8275 , 4.8276 ,
4.8277 , 4.8278 , 4.8279 , 4.828  , 4.8281 , 4.8282 ,
4.8283 , 4.8284 , 4.8285 , 4.8286 , 4.8287 , 4.8288 ,
4.8289 , 4.829  , 4.8291 , 4.8292 , 4.8293 , 4.8294 ,
4.8295 , 4.8296 , 4.8297 , 4.8298 , 4.8299 , 4.83  ,
4.8301 , 4.8302 , 4.8303 , 4.8304 , 4.8305 , 4.8306 ,
4.8307 , 4.8308 , 4.8309 , 4.831  , 4.8311 , 4.8312 ,
4.8313 , 4.8314 , 4.8315 , 4.8316 , 4.8317 , 4.8318 ,
4.8319 , 4.832  , 4.8321 , 4.8322 , 4.8323 , 4.8324 ,
4.8325 , 4.8326 , 4.8327 , 4.8328 , 4.8329 , 4.833  ,
4.8331 , 4.8332 , 4.8333 , 4.8334 , 4.8335 , 4.8336 ,
4.8337 , 4.8338 , 4.8339 , 4.834  , 4.8341 , 4.8342 ,
4.8343 , 4.8344 , 4.8345 , 4.8346 , 4.8347 , 4.8348 ,
4.8349 , 4.835  , 4.8351 , 4.8352 , 4.8353 , 4.8354 ,
4.8355 , 4.8356 , 4.8357 , 4.8358 , 4.8359 , 4.836  ,
4.8361 , 4.8362 , 4.8363 , 4.8364 , 4.8365 , 4.8366 ,
4.8367 , 4.8368 , 4.8369 , 4.837  , 4.8371 , 4.8372 ,
4.8373 , 4.8374 , 4.8375 , 4.8376 , 4.8377 , 4.8378 ,
4.8379 , 4.838  , 4.8381 , 4.8382 , 4.8383 , 4.8384 ,
4.8385 , 4.8386 , 4.8387 , 4.8388 , 4.8389 , 4.839  ,
4.8391 , 4.8392 , 4.8393 , 4.8394 , 4.8395 , 4.8396 ,
4.8397 , 4.8398 , 4.8399 , 4.84  , 4.8401 , 4.8402 ,
4.8403 , 4.8404 , 4.8405 , 4.8406 , 4.8407 , 4.8408 ,
4.8409 , 4.841  , 4.8411 , 4.8412 , 4.8413 , 4.8414 ,
4.8415 , 4.8416 , 4.8417 , 4.8418 , 4.8419 , 4.842  ,
4.8421 , 4.8422 , 4.8423 , 4.8424 , 4.8425 , 4.8426 ,
4.8427 , 4.8428 , 4.8429 , 4.843  , 4.8431 , 4.8432 ,
4.8433 , 4.8434 , 4.8435 , 4.8436 , 4.8437 , 4.8438 ,
4.8439 , 4.844  , 4.8441 , 4.8442 , 4.8443 , 4.8444 ,
4.8445 , 4.8446 , 4.8447 , 4.8448 , 4.8449 , 4.845  ,
4.8451 , 4.8452 , 4.8453 , 4.8454 , 4.8455 , 4.8456 ,
4.8457 , 4.8458 , 4.8459 , 4.846  , 4.8461 , 4.8462 ,
4.8463 , 4.8464 , 4.8465 , 4.8466 , 4.8467 , 4.8468 ,
4.8469 , 4.847  , 4.8471 , 4.8472 , 4.8473 , 4.8474 ,
4.8475 , 4.8476 , 4.8477 , 4.8478 , 4.8479 , 4.848  ,
4.8481 , 4.8482 , 4.8483 , 4.8484 , 4.8485 , 4.8486 ,
4.8487 , 4.8488 , 4.8489 , 4.849  , 4.8491 , 4.8492 ,
4.8493 , 4.8494 , 4.8495 , 4.8496 , 4.8497 , 4.8498 ,
4.8499 , 4.85  , 4.8501 , 4.8502 , 4.8503 , 4.8504 ,
4.8505 , 4.8506 , 4.8507 , 4.8508 , 4.8509 , 4.851  ,
4.8511 , 4.8512 , 4.8513 , 4.8514 , 4.8515 , 4.8516 ,
4.8517 , 4.8518 , 4.8519 , 4.852  , 4.8521 , 4.8522 ,
4.8523 , 4.8524 , 4.8525 , 4.8526 , 4.8527 , 4.8528 ,
4.8529 , 4.853  , 4.8531 , 4.8532 , 4.8533 , 4.8534 ,
4.8535 , 4.8536 , 4.8537 , 4.8538 , 4.8539 , 4.854  ,
4.8541 , 4.8542 , 4.8543 , 4.8544 , 4.8545 , 4.8546 ,
4.8547 , 4.8548 , 4.8549 , 4.855  , 4.8551 , 4.8552 ,
4.8553 , 4.8554 , 4.8555 , 4.8556 , 4.8557 , 4.8558 ,
4.8559 , 4.856  , 4.8561 , 4.8562 , 4.8563 , 4.8564 ,
4.8565 , 4.8566 , 4.8567 , 4.8568 , 4.8569 , 4.857  ,
4.8571 , 4.8572 , 4.8573 , 4.8574 , 4.8575 , 4.8576 ,
4.8577 , 4.8578 , 4.8579 , 4.858  , 4.8581 , 4.8582 ,
4.8583 , 4.8584 , 4.8585 , 4.8586 , 4.8587 , 4.8588 ,
4.8589 , 4.859  , 4.8591 , 4.8592 , 4.8593 , 4.8594 ,
4.8595 , 4.8596 , 4.8597 , 4.8598 , 4.8599 , 4.86  ,
4.8601 , 4.8602 , 4.8603 , 4.8604 , 4.8605 , 4.8606 ,
4.8607 , 4.8608 , 4.8609 , 4.861  , 4.8611 , 4.8612 ,
4.8613 , 4.8614 , 4.8615 , 4.8616 , 4.8617 , 4.8618 ,
4.8619 , 4.862  , 4.8621 , 4.8622 , 4.8623 , 4.8624 ,
4.8625 , 4.8626 , 4.8627 , 4.8628 , 4.8629 , 4.863  ,
4.8631 , 4.8632 , 4.8633 , 4.8634 , 4.8635 , 4.8636 ,
4.8637 , 4.8638 , 4.8639 , 4.864  , 4.8641 , 4.8642 ,
4.8643 , 4.8644 , 4.8645 , 4.8646 , 4.8647 , 4.8648 ,
4.8649 , 4.865  , 4.8651 , 4.8652 , 4.8653 , 4.8654 ,
4.8655 , 4.8656 , 4.8657 , 4.8658 , 4.8659 , 4.866  ,
4.8661 , 4.8662 , 4.8663 , 4.8664 , 4.8665 , 4.8666 ,
4.8667 , 4.8668 , 4.8669 , 4.867  , 4.8671 , 4.8672 ,
4.8673 , 4.8674 , 4.8675 , 4.8676 , 4.8677 , 4.8678 ,
4.8679 , 4.868  , 4.8681 , 4.8682 , 4.8683 , 4.8684 ,
4.8685 , 4.8686 , 4.8687 , 4.8688 , 4.8689 , 4.869  ,
4.8691 , 4.8692 , 4.8693 , 4.8694 , 4.8695 , 4.8696 ,
4.8697 , 4.8698 , 4.8699 , 4.87  , 4.8701 , 4.8702 ,
4.8703 , 4.8704 , 4.8705 , 4.8706 , 4.8707 , 4.8708 ,
4.8709 , 4.871  , 4.8711 , 4.8712 , 4.8713 , 4.8714 ,
4.8715 , 4.8716 , 4.8717 , 4.8718 , 4.8719 , 4.872  ,
4.8721 , 4.8722 , 4.8723 , 4.8724 , 4.8725 , 4.8726 ,
4.8727 , 4.8728 , 4.8729 , 4.873  , 4.8731 , 4.8732 ,
4.8733 , 4.8734 , 4.8735 , 4.8736 , 4.8737 , 4.8738 ,
4.8739 , 4.874  , 4.8741 , 4.8742 , 4.8743 , 4.8744 ,
4.8745 , 4.8746 , 4.8747 , 4.8748 , 4.8749 , 4.875  ,
4.8751 , 4.8752 , 4.8753 , 4.8754 , 4.8755 , 4.8756 ,
4.8757 , 4.8758 , 4.8759 , 4.876  , 4.8761 , 4.8762 ,
4.8763 , 4.8764 , 4.8765 , 4.8766 , 4.8767 , 4.8768 ,
4.8769 , 4.877  , 4.8771 , 4.8772 , 4.8773 , 4.8774 ,
4.8775 , 4.8776 , 4.8777 , 4.8778 , 4.8779 , 4.878  ,
4.8781 , 4.8782 , 4.8783 , 4.8784 , 4.8785 , 4.8786 ,
4.8787 , 4.8788 , 4.8789 , 4.879  , 4.8791 , 4.8792 ,
4.8793 , 4.8794 , 4.8795 , 4.8796 , 4.8797 , 4.8798 ,
4.8799 , 4.88  , 4.8801 , 4.8802 , 4.8803 , 4.8804 ,
4.8805 , 4.8806 , 4.8807 , 4.8808 , 4.8809 , 4.881  ,
4.8811 , 4.8812 , 4.8813 , 4.8814 , 4.8815 , 4.8816 ,
4.8817 , 4.8818 , 4.8819 , 4.882  , 4.8821 , 4.8822 ,
4.8823 , 4.8824 , 4.8825 , 4.8826 , 4.8827 , 4.8828 ,
4.8829 , 4.883  , 4.8831 , 4.8832 , 4.8833 , 4.8834 ,
4.8835 , 4.8836 , 4.8837 , 4.8838 , 4.8839 , 4.884  ,
4.8841 , 4.8842 , 4.8843 , 4.8844 , 4.8845 , 4.8846 ,
4.8847 , 4.8848 , 4.8849 , 4.885  , 4.8851 , 4.8852 ,
4.8853 , 4.8854 , 4.8855 , 4.8856 , 4.8857 , 4.8858 ,
4.8859 , 4.886  , 4.8861 , 4.8862 , 4.8863 , 4.8864 ,
4.8865 , 4.8866 , 4.8867 , 4.8868 , 4.8869 , 4.887  ,
4.8871 , 4.8872 , 4.8873 , 4.8874 , 4.8875 , 4.8876 ,
4.8877 , 4.8878 , 4.8879 , 4.888  , 4.8881 , 4.8882 ,
4.8883 , 4.8884 , 4.8885 , 4.8886 , 4.8887 , 4.8888 ,
4.8889 , 4.889  , 4.8891 , 4.8892 , 4.8893 , 4.8894 ,
4.8895 , 4.8896 , 4.8897 , 4.8898 , 4.8899 , 4.89  ,
4.8901 , 4.8902 , 4.8903 , 4.8904 , 4.8905 , 4.8906 ,
4.8907 , 4.8908 , 4.8909 , 4.891  , 4.8911 , 4.8912 ,
4.8913 , 4.8914 , 4.8915 , 4.8916 , 4.8917 , 4.8918 ,
4.8919 , 4.892  , 4.8921 , 4.8922 , 4.8923 , 4.8924 ,
4.8925 , 4.8926 , 4.8927 , 4.8928 , 4.8929 , 4.893  ,
4.8931 , 4.8932 , 4.8933 , 4.8934 , 4.8935 , 4.8936 ,
4.8937 , 4.8938 , 4.8939 , 4.894  , 4.8941 , 4.8942 ,
4.8943 , 4.8944 , 4.8945 , 4.8946 , 4.8947 , 4.8948 ,
4.8949 , 4.895  , 4.8951 , 4.8952 , 4.8953 , 4.8954 ,
4.8955 , 4.8956 , 4.8957 , 4.8958 , 4.8959 , 4.896  ,
4.8961 , 4.8962 , 4.8963 , 4.8964 , 4.8965 , 4.8966 ,
4.8967 , 4.8968 , 4.8969 , 4.897  , 4.8971 , 4.8972 ,
4.8973 , 4.8974 , 4.8975 , 4.8976 , 4.8977 , 4.8978 ,
4.8979 , 4.898  , 4.8981 , 4.8982 , 4.8983 , 4.8984 ,
4.8985 , 4.8986 , 4.8987 , 4.8988 , 4.8989 , 4.899  ,
4.8991 , 4.8992 , 4.8993 , 4.8994 , 4.8995 , 4.8996 ,
4.8997 , 4.8998 , 4.8999 , 4.9  , 4.9001 , 4.9002 ,
4.9003 , 4.9004 , 4.9005 , 4.9006 , 4.9007 , 4.9008 ,
4.9009 , 4.901  , 4.9011 , 4.9012 , 4.9013 , 4.9014 ,
4.9015 , 4.9016 , 4.9017 , 4.9018 , 4.9019 , 4.902  ,
4.9021 , 4.9022 , 4.9023 , 4.9024 , 4.9025 , 4.9026 ,
4.9027 , 4.9028 , 4.9029 , 4.903  , 4.9031 , 4.9032 ,
4.9033 , 4.9034 , 4.9035 , 4.9036 , 4.9037 , 4.9038 ,
4.9039 , 4.904  , 4.9041 , 4.9042 , 4.9043 , 4.9044 ,
4.9045 , 4.9046 , 4.9047 , 4.9048 , 4.9049 , 4.905  ,
4.9051 , 4.9052 , 4.9053 , 4.9054 , 4.9055 , 4.9056 ,
4.9057 , 4.9058 , 4.9059 , 4.906  , 4.9061 , 4.9062 ,
4.9063 , 4.9064 , 4.9065 , 4.9066 , 4.9067 , 4.9068 ,
4.9069 , 4.907  , 4.9071 , 4.9072 , 4.9073 , 4.9074 ,
4.9075 , 4.9076 , 4.9077 , 4.9078 , 4.9079 , 4.908  ,
4.9081 , 4.9082 , 4.9083 , 4.9084 , 4.9085 , 4.9086 ,
4.9087 , 4.9088 , 4.9089 , 4.909  , 4.9091 , 4.9092 ,
4.9093 , 4.9094 , 4.9095 , 4.9096 , 4.9097 , 4.9098 ,
4.9099 , 4.91  , 4.9101 , 4.9102 , 4.9103 , 4.9104 ,
4.9105 , 4.9106 , 4.9107 , 4.9108 , 4.9109 , 4.911  ,
4.9111 , 4.9112 , 4.9113 , 4.9114 , 4.9115 , 4.9116 ,
4.9117 , 4.9118 , 4.9119 , 4.912  , 4.9121 , 4.9122 ,
4.9123 , 4.9124 , 4.9125 , 4.9126 , 4.9127 , 4.9128 ,
4.9129 , 4.913  , 4.9131 , 4.9132 , 4.9133 , 4.9134 ,
4.9135 , 4.9136 , 4.9137 , 4.9138 , 4.9139 , 4.914  ,
4.9141 , 4.9142 , 4.9143 , 4.9144 , 4.9145 , 4.9146 ,
4.9147 , 4.9148 , 4.9149 , 4.915  , 4.9151 , 4.9152 ,
4.9153 , 4.9154 , 4.9155 , 4.9156 , 4.9157 , 4.9158 ,
4.9159 , 4.916  , 4.9161 , 4.9162 , 4.9163 , 4.9164 ,
4.9165 , 4.9166 , 4.9167 , 4.9168 , 4.9169 , 4.917  ,
4.9171 , 4.9172 , 4.9173 , 4.9174 , 4.9175 , 4.9176 ,
4.9177 , 4.9178 , 4.9179 , 4.918  , 4.9181 , 4.9182 ,
4.9183 , 4.9184 , 4.9185 , 4.9186 , 4.9187 , 4.9188 ,
4.9189 , 4.919  , 4.9191 , 4.9192 , 4.9193 , 4.9194 ,
4.9195 , 4.9196 , 4.9197 , 4.9198 , 4.9199 , 4.92  ,
4.9201 , 4.9202 , 4.9203 , 4.9204 , 4.9205 , 4.9206 ,
4.9207 , 4.9208 , 4.9209 , 4.921  , 4.9211 , 4.9212 ,
4.9213 , 4.9214 , 4.9215 , 4.9216 , 4.9217 , 4.9218 ,
4.9219 , 4.922  , 4.9221 , 4.9222 , 4.9223 , 4.9224 ,
4.9225 , 4.9226 , 4.9227 , 4.9228 , 4.9229 , 4.923  ,
4.9231 , 4.9232 , 4.9233 , 4.9234 , 4.9235 , 4.9236 ,
4.9237 , 4.9238 , 4
```

```

1.7918 , 1.8083 , 1.8245 , 1.8405 , 1.8563 , 1.8718 ,
1.8871 , 1.9021 , 1.9169 , 1.9315 , 1.9459 , 1.9601 ,
1.9741 , 1.9879 , 2.0015 , 2.0149 , 2.0281 , 2.0412 ,
2.0541 , 2.0669 , 2.0794 , 2.0919 , 2.1041 , 2.1163 ,
2.1282 , 2.1401 , 2.1518 , 2.1633 , 2.1748 , 2.1861 ,
2.1972 , 2.2083 , 2.2192 , 2.23 , 2.2407 , 2.2513 ,
2.2618 , 2.2721 , 2.2824 , 2.2925 , 2.3026 ])

```

```
[ ]: instance_to_drop = np.where(np.isinf(real_feature))[0][0]
print(instance_to_drop)
```

100

Because we have no information given about these features we cannot infer what -inf represents. If this was clear, we could find an appropriate solution to deal with it, e.g. imputation. Moreover, there is just one instance with -inf value for this feature, thus dropping it from the dataset makes more sense since it will not have a significative impact on our final results.

Nevertheless, we have also tried mean imputation to handle this value during our experiments with the Complex Matrices' approach. Results and comparisons can be found in [Annex: Complex Matrices](#).

### 1.0.3 Replace fourth feature by new features and delete instance with infinity

```
[ ]: # delete original feature column
X = np.delete(X_potential_drivers, 3, 1)
# add new feature columns
n = len(Y_snow)
X = np.hstack((X[:,[0,1,2]],
               np.array(real_feature).reshape(n,1),
               np.array(imag_feature).reshape(n,1),
               np.array(abs_feature).reshape(n,1),
               X[:,3].reshape(n,1)))

# drop the feature with -inf
X = np.delete(X,instance_to_drop, axis=0)
X = X.astype('float64')

X[:3]
```

```
[ ]: array([[ 11.728      , -0.64836      , -5.          , 2.3026      ,
            3.1416      , 3.89507604, -23.          ],
           [ 11.771      , -0.51455      , -4.95        , 2.2925      ,
            3.1416      , 3.88911388, -22.7        ],
           [ 11.817      , -0.39388      , -4.9          , 2.2824      ,
            3.1416      , 3.88316885, -22.4        ]])
```

```
[ ]: # we also need to delete the instance in the target array  
y = np.delete(Y_snow, instance_to_drop)
```

```
[ ]: # number of instances  
n = X.shape[0]  
n
```

```
[ ]: 200
```

```
[ ]: # number of features:  
p = X.shape[1]  
p
```

```
[ ]: 7
```

#### 1.0.4 Data Inspection

```
[ ]: feature_names = ['1.feature', '2.feature', '3.feature', '4.feature: real', '4.  
→feature: imaginary', '4.feature: absolute', '5.feature']
```

```
[ ]: df_X = pd.DataFrame(X, columns= feature_names)
```

```
print("- First 5 realisations:")  
print(df_X.head())  
print("\n- Type:")  
print(df_X.dtypes)  
print("\n- Main statistics:")  
print(df_X.describe())
```

- First 5 realisations:

	1.feature	2.feature	3.feature	4.feature: real	4.feature: imaginary	\
0	11.728	-0.64836	-5.00	2.3026	3.1416	
1	11.771	-0.51455	-4.95	2.2925	3.1416	
2	11.817	-0.39388	-4.90	2.2824	3.1416	
3	11.864	-0.28239	-4.85	2.2721	3.1416	
4	11.913	-0.17704	-4.80	2.2618	3.1416	

```
4.feature: absolute 5.feature  
0 3.895076 -23.0  
1 3.889114 -22.7  
2 3.883169 -22.4  
3 3.877124 -22.1  
4 3.871097 -21.8
```

- Type:

1.feature	float64
2.feature	float64
3.feature	float64

```

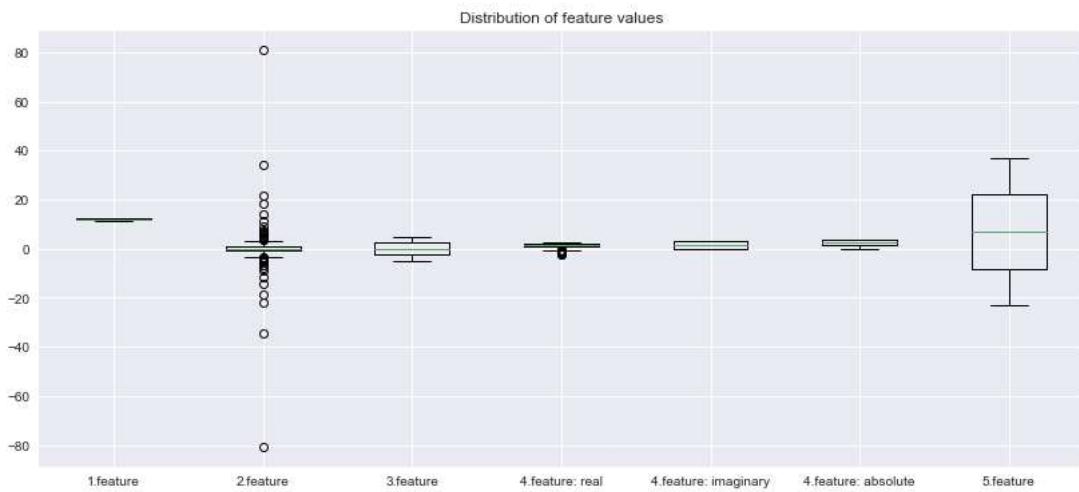
4.feature: real          float64
4.feature: imaginary     float64
4.feature: absolute      float64
5.feature                 float64
dtype: object

- Main statistics:
    1.feature   2.feature   3.feature  4.feature: real \
count  200.000000  2.000000e+02  200.000000  200.000000
mean   12.000000 -7.105427e-17  0.000000  1.334810
std    0.346815  9.682931e+00  2.915691  0.925720
min   11.500000 -8.071300e+01 -5.000000 -2.302600
25%   11.667250 -8.931725e-01 -2.512500  0.945705
50%   12.000000  0.000000e+00  0.000000  1.619300
75%   12.332750  8.931725e-01  2.512500  2.018200
max   12.500000  8.071300e+01  5.000000  2.302600

        4.feature: imaginary  4.feature: absolute  5.feature
count          200.000000  200.000000  200.000000
mean          1.570800   2.510633   7.000000
std           1.574742   1.127963  17.494148
min          0.000000   0.000000 -23.000000
25%          0.000000   1.643825 -8.075000
50%          1.570800   2.722100   7.000000
75%          3.141600   3.541173  22.075000
max          3.141600   3.895076  37.000000

```

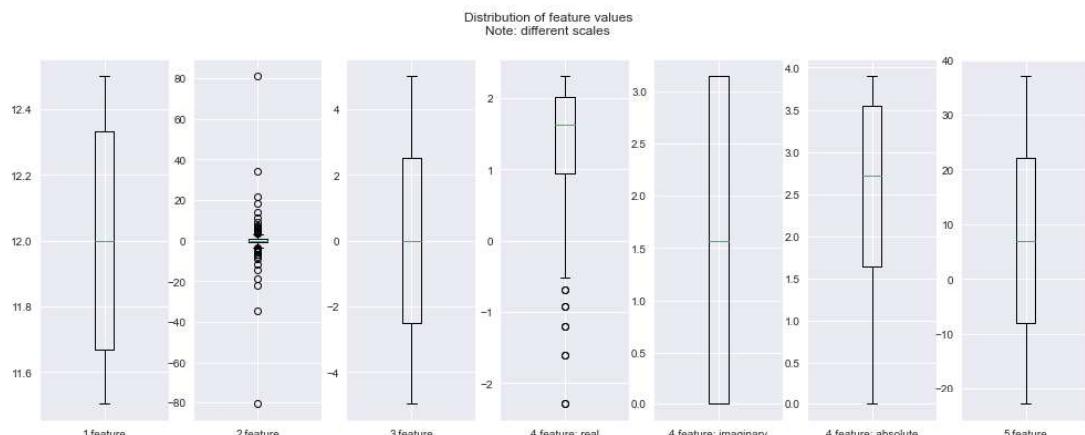
```
[ ]: fig, ax = plt.subplots(figsize=(14,6))
ax.boxplot(X)
ax.set_title('Distribution of feature values')
ax.set_xticklabels(feature_names)
plt.show()
```



In order to get a better impression of the feature values we create individual boxplots for each feature using **optimal yaxis scales for each feature**.

**Note:** In the above plot, we have drawn all the features together to compare their distributions using the same scale value for the y-axis, so we can have a general idea about the differences between their value ranges. On the other hand, the next plot just wants to assess the singularities of each distribution individually (i.e. kind of distribution, skeweness, mean,...). Thus, we need a suitable scale value for each of them in order to illustrate these characteristics properly.

```
[ ]: fig, ax = plt.subplots(1,p,figsize=(17,6))
for i in range(p):
    ax[i].boxplot(X[:,i])
    name = str(i+1) + '.feature'
    ax[i].set_xticklabels([feature_names[i]])
plt.suptitle('Distribution of feature values\nNote: different scales')
plt.show()
```



As you can see, each feature has a quite unique distribution. The first feature has a mean of approximately 12, the third feature of approximately 0. Some features have values in a quite small range, e.g. the 4.feature, while others have a huge range, e.g. the 5.feature.

In the following we use the following schema:  
 \* triangles: original 4.feature has an imaginary part  
 \* dots: original 4.feature has no imaginary part

and the color of the triangles/dots represents the y value

```
[ ]: # resource: https://stackoverflow.com/questions/55271471/
      ↵plotting-some-third-variable-against-x-and-y-in-matplotlib-scatter
import matplotlib.colors as colors
import matplotlib.cm as cmx
```

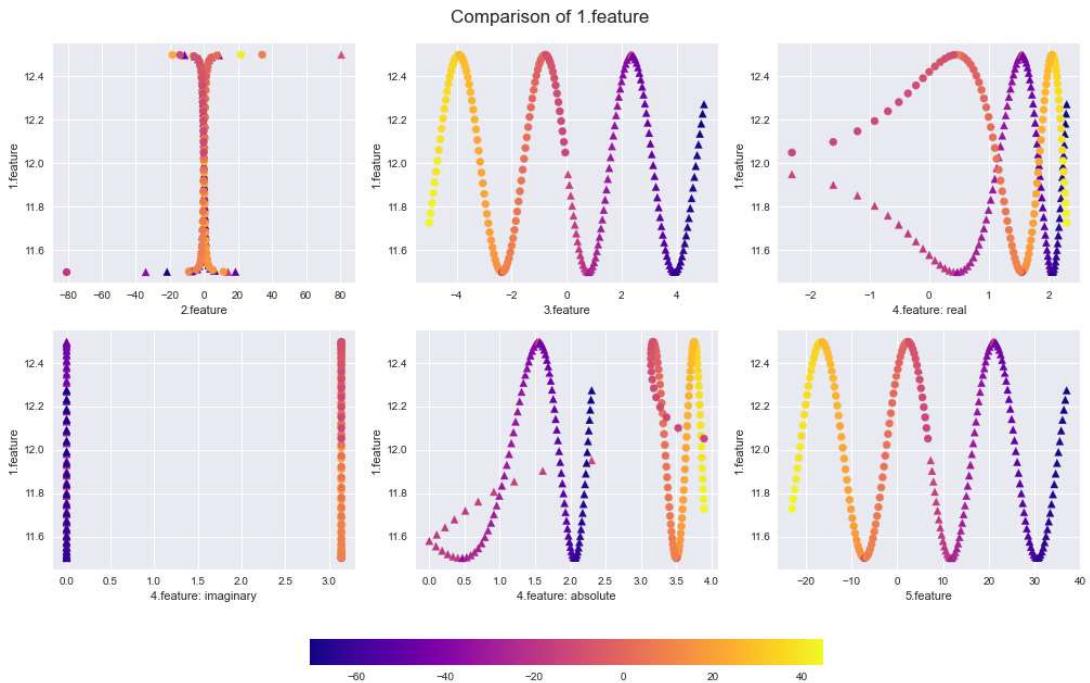
```
[ ]: def compare_specific_feature(index_feature, title, feature_names):
    fig, axis = plt.subplots(2, 3, figsize=(17,11))
    cmap = plt.get_cmap('plasma')
    cNorm = colors.Normalize(vmin=min(y), vmax=max(y))
    scalarMap = cmx.ScalarMappable(norm=cNorm, cmap=cmap)
    fig.subplots_adjust(top=0.94)
    f = 0

    for i,ax in enumerate(axis.flatten()):
        index_to_compare = i
        if i>=index_feature:
            index_to_compare = i+1
        else:
            index_to_compare = i
        ax.scatter(X[100:,index_to_compare],X[100:,index_feature], c=scalarMap.
        ↪to_rgba(y[100:]), marker='^', cmap=cmx.plasma)
        ax.scatter(X[:100,index_to_compare],X[:100,index_feature], c=scalarMap.
        ↪to_rgba(y[:100]), marker='o', cmap=cmx.plasma)
        ax.set_xlabel(feature_names[index_to_compare])
        ax.set_ylabel(feature_names[index_feature])

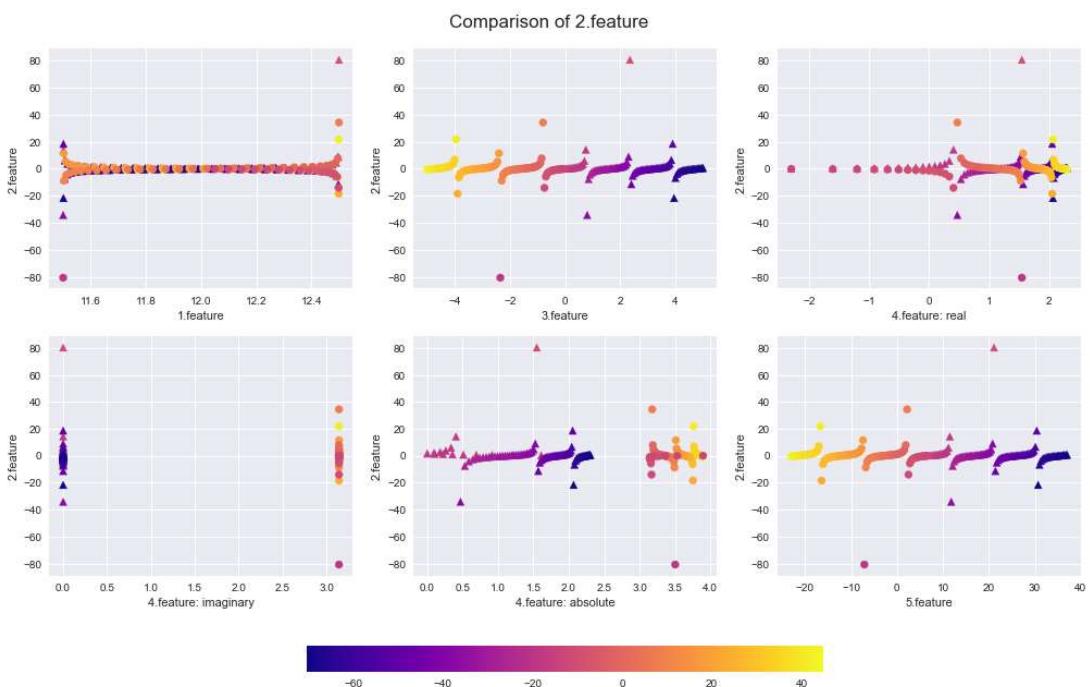
    plt.suptitle(title, fontsize=17, y=0.98)
    fig.colorbar(scalarMap,ax=axis[:, :],location='bottom', shrink=0.5,pad=0.1)
    plt.show()
```

```
[ ]: feature_names = ['1.feature','2.feature','3.feature','4.feature: real','4.
    ↪feature: imaginary','4.feature: absolute','5.feature']
```

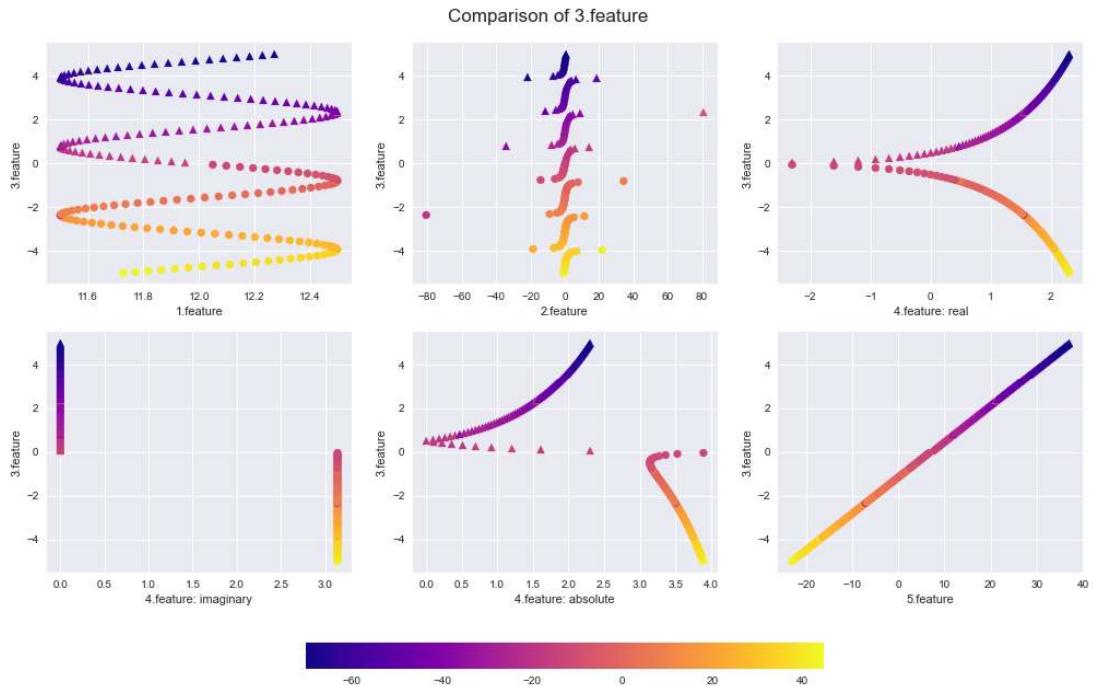
```
[ ]: compare_specific_feature(index_feature=0, title='Comparison of 1.feature', ↪
    ↪feature_names=feature_names)
```



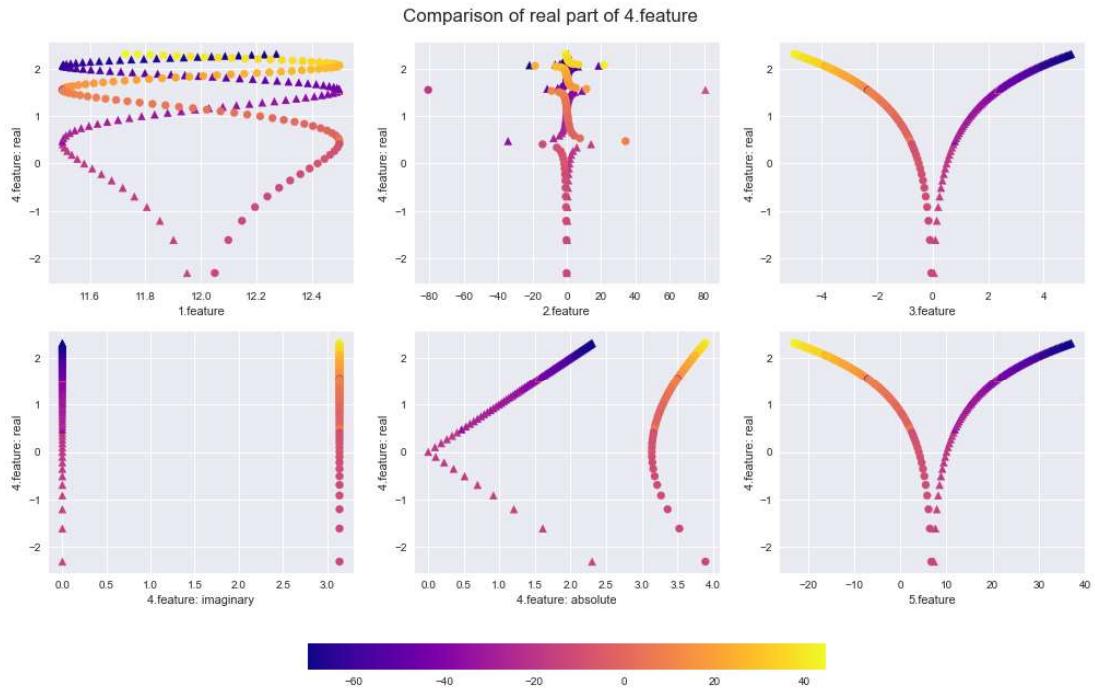
```
[ ]: compare_specific_feature(index_feature=1, title='Comparison of 2.feature',
                           ↪feature_names=feature_names)
```



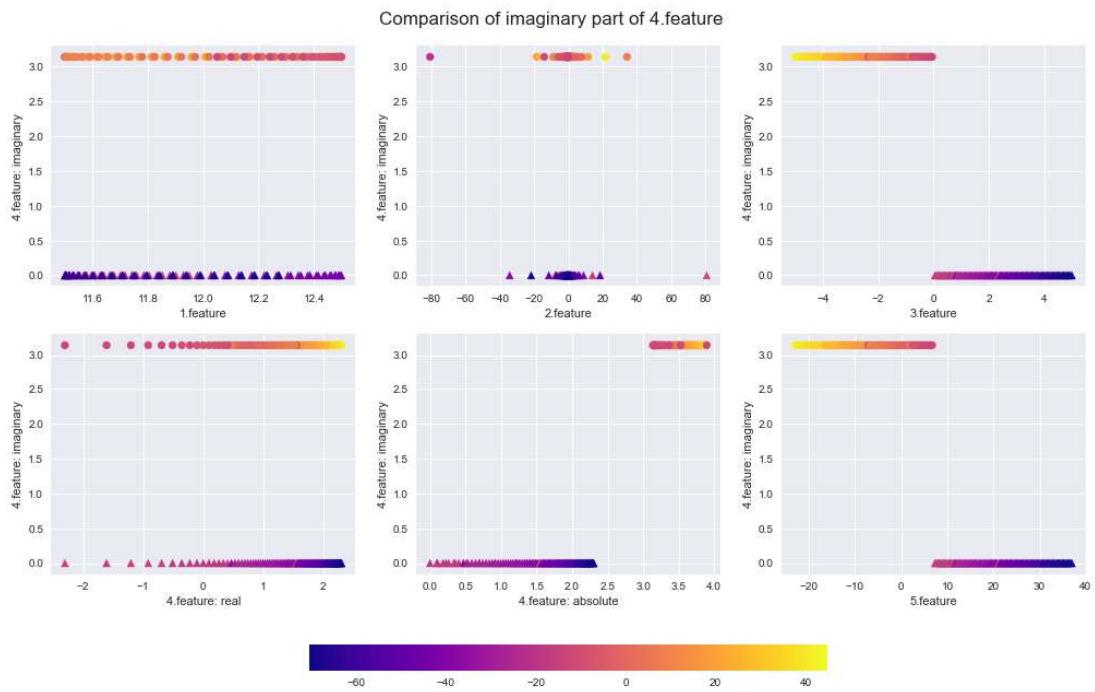
```
[ ]: compare_specific_feature(index_feature=2, title='Comparison of 3.feature',  
                             ↪feature_names=feature_names)
```



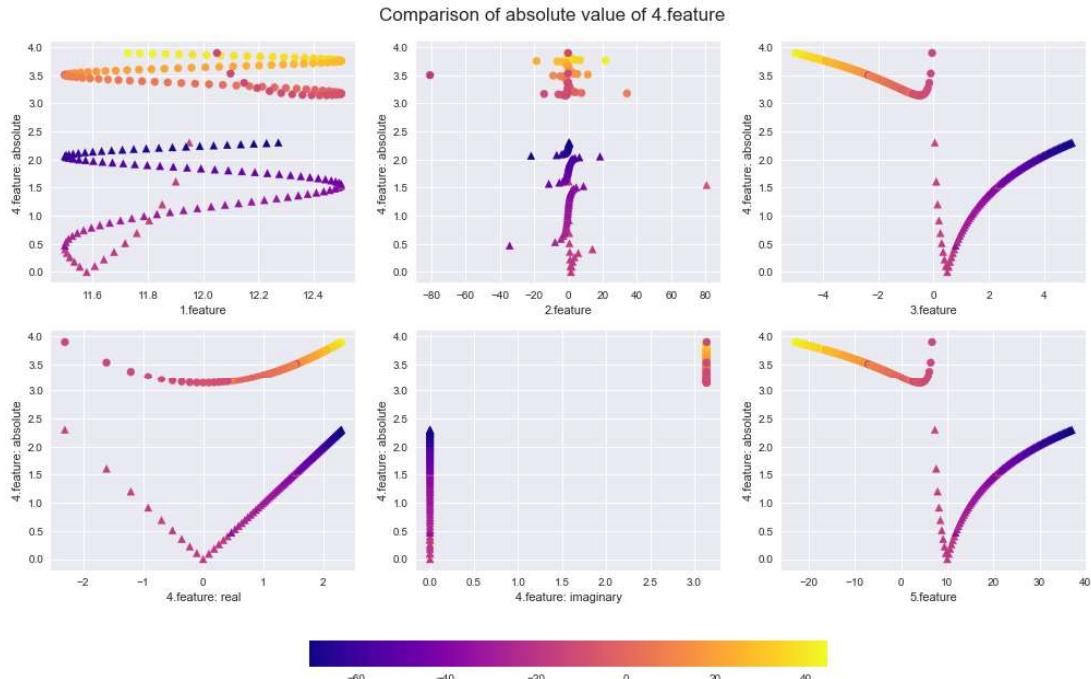
```
[ ]: compare_specific_feature(index_feature=3, title='Comparison of real part of 4.  
                             ↪feature', feature_names=feature_names)
```



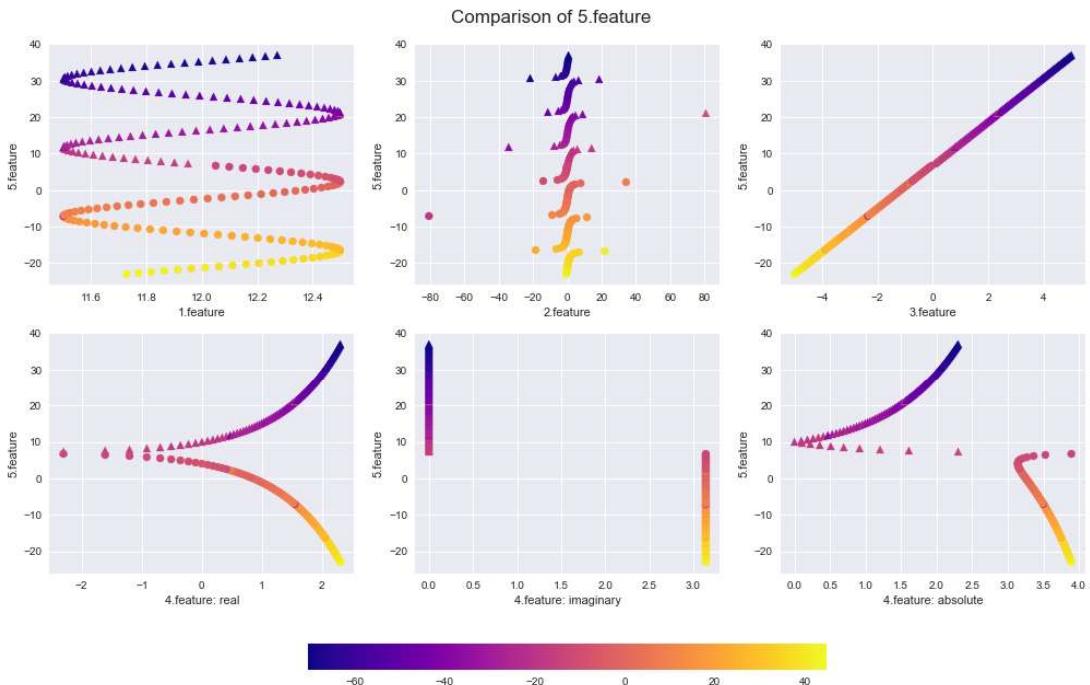
```
[ ]: compare_specific_feature(index_feature=4, title='Comparison of imaginary part of 4.feature', feature_names=feature_names)
```



```
[ ]: compare_specific_feature(index_feature=5, title='Comparison of absolute value of 4.feature', feature_names=feature_names)
```



```
[ ]: compare_specific_feature(index_feature=6, title='Comparison of 5.feature', feature_names=feature_names)
```



One important observation can be done in these plots: The 3. and the 5.feature seem correlated. Let's compute the Pearson correlation coefficient:

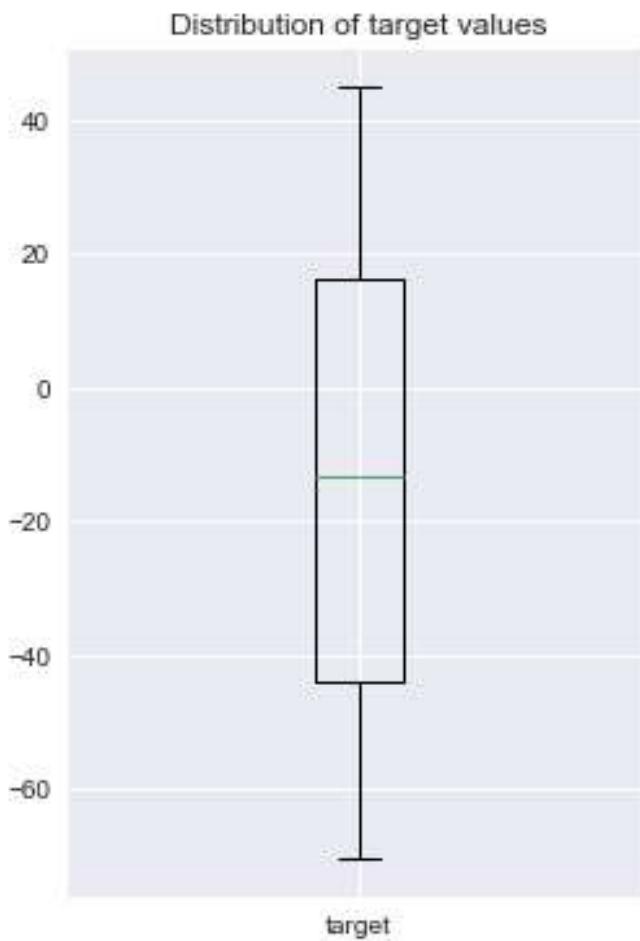
```
[ ]: from scipy.stats import pearsonr
corr, _ = pearsonr(X[:,2],X[:,6])
print('Pearson correlation: %.3f' % corr)
```

Pearson correlation: 1.000

1 is the highest possible Pearson correlation score and indicates that the 3. and the 7.feature are perfectly linear correlated.

### Target

```
[ ]: fig, ax = plt.subplots(figsize=(4,6))
ax.boxplot(y)
ax.set_title('Distribution of target values')
ax.set_xticklabels(['target'])
plt.show()
```



### 1.0.5 Should we extend the feature matrix for an additional column of 1's to allow a bias?

```
[ ]: p = X.shape[1]
r = np.linalg.matrix_rank(X)
print('Number of features:', n)
print('Column rank of matrix:', r)
```

Number of features: 200  
 Column rank of matrix: 7

```
[ ]: X_extended = np.hstack((np.ones((n,1)),X))
r_extended = np.linalg.matrix_rank(X_extended)
print('Number of columns in matrix:', p+1)
print('Column rank of matrix:', r_extended)
```

Number of columns in matrix: 8  
 Column rank of matrix: 7

If  $X$  is not of full column rank,  $X^T X$  will also not be of full column rank, i.e. be a singular matrix. This means, it is not invertible (Using the python function `linalg.inv()` we will probably be able to compute an output, but it is questionable how valuable this is.) This is probably due to the correlation of the 3. and the 5.feature. In the following we will thus use two types of matrixes. See **3.2 Second approach: Select features based on corrected coefficient of determination** for more details on this.

## 2 2 Prepare functions for computation and evaluation of model

### 2.0.1 Define the function that can compute the LS-estimator

```
[ ]: def beta_hat(X,y):
    return inv(X.T@X)@X.T@y
```

### 2.0.2 Coefficient of Determination

```
[ ]: def r_sq(true, pred):
    mean_true = np.mean(true)
    nominator = np.sum(np.power(true-pred,2))
    denominator = np.sum(np.power(true-mean_true,2))
    return 1 - (nominator/denominator)
```

### 2.0.3 The corrected coefficient of determination

Motivation for corrected/adjusted coefficient of determination:

In contrast to the regular coefficient of determination this statistics takes the number of features  $p$  into account. The coefficient of determination has the characteristics that increasing the number of features will **never** decrease it, it can only increase. Since we want to compare different feature combinations we don't want that the number of features itself to affect our evaluation. The corrected/adjusted coefficient of determination is supposed to provide exactly this.

```
[ ]: def corr_r_sq(true, pred, p):
    n = len(true)
    return 1 - ((n-1)/(n-p-1))*(1-r_sq(true, pred))
```

## 3 3 Feature Selection

### 3.1 3.1 First approach: Select features based on observations in plots

During the data analysis we noticed dependencies of our labels  $y$  to the 3. and 5.feature:

```
[ ]: fig, ax = plt.subplots(1, 2, figsize=(10,5))
cmap = plt.get_cmap('plasma')
cNorm = colors.Normalize(vmin=min(y), vmax=max(y))
scalarMap = cmx.ScalarMappable(norm=cNorm, cmap=cmap)
ax[0].scatter(X[100:,2],X[100:,0], c=scalarMap.to_rgba(y[100:]), marker='^', cmap=cmap)
```

```

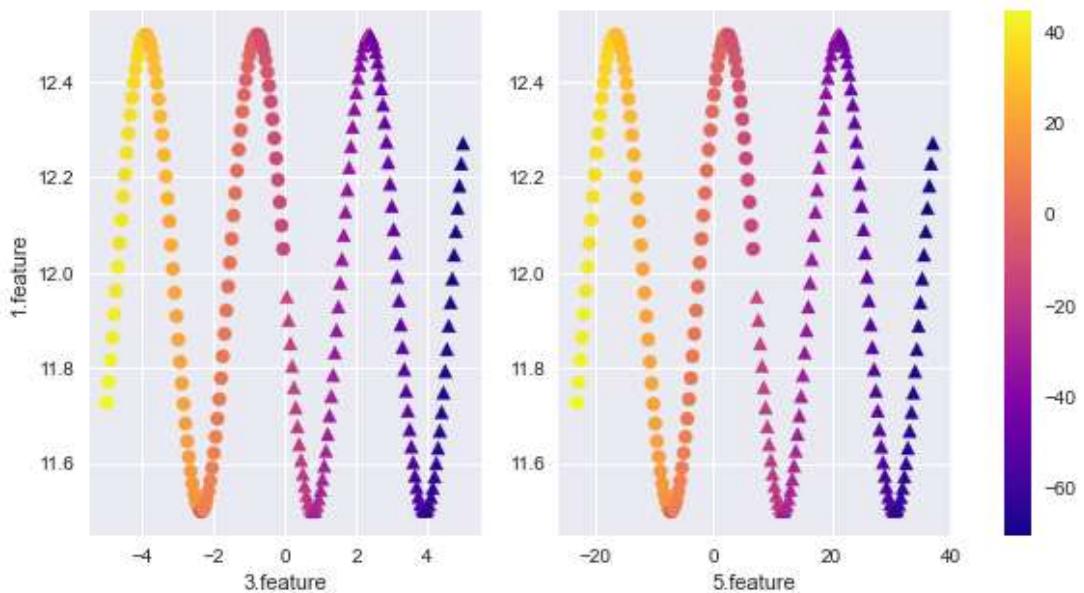
ax[0].scatter(X[:100,2],X[:100,0], c=scalarMap.to_rgba(y[:100]), marker='o',  

              cmap=cmx.plasma)
ax[0].set_xlabel('3.feature')
ax[0].set_ylabel('1.feature')
ax[1].scatter(X[100:,6],X[100:,0], c=scalarMap.to_rgba(y[100:]), marker='^',  

              cmap=cmx.plasma)
ax[1].scatter(X[:100,6],X[:100,0], c=scalarMap.to_rgba(y[:100]), marker='o',  

              cmap=cmx.plasma)
ax[1].set_xlabel('5.feature')
#ax[1].set_ylabel('1.feature')
fig.colorbar(scalarMap,ax=ax[:])
plt.show()

```

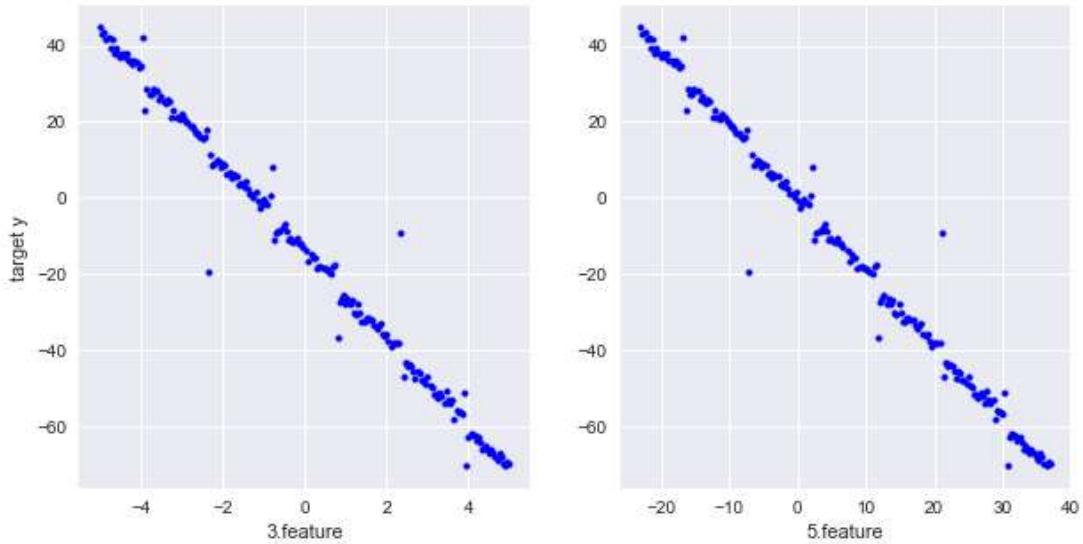


These two plots, that already appeared in our data analysis part, show that the target (represented by the color) can be described pretty well both by the 3. and the 5.feature. In the following two plots this relationship gets even more obvious:

```

[ ]: fig, ax = plt.subplots(1, 2, figsize=(10,5))
ax[0].scatter(X[:,2],y, s=12, facecolors='blue')
ax[0].set_xlabel('3.feature')
ax[0].set_ylabel('target y')
ax[1].scatter(X[:,6],y, s=12, facecolors='blue')
ax[1].set_xlabel('5.feature')
plt.show()

```



Based on this observation, it seems realistic that we can already get a good model just using the 3. or the 5.feature. Here we can include a bias.

```
[ ]: # Check the rank of the matrix we will use for our first model:
np.linalg.matrix_rank(X_extended[:,[0,3]]) == 2

[ ]: True

[ ]: # Check the rank of the matrix we will use for our second model:
np.linalg.matrix_rank(X_extended[:,[0,7]]) == 2

[ ]: True

[ ]: beta1 = beta_hat(X_extended[:,[0,3]],y)
pred1 = X_extended[:,[0,3]]@beta1
r_sq1 = r_sq(y, pred1)
corr_r_sq1 = corr_r_sq(y, pred1, 2)
print('Results of model using 3.feature and bias:')
print('- coefficient of determination: {:.6f}'.format(r_sq1))
print('- corrected coefficient of determination: {:.6f}'.format(corr_r_sq1))

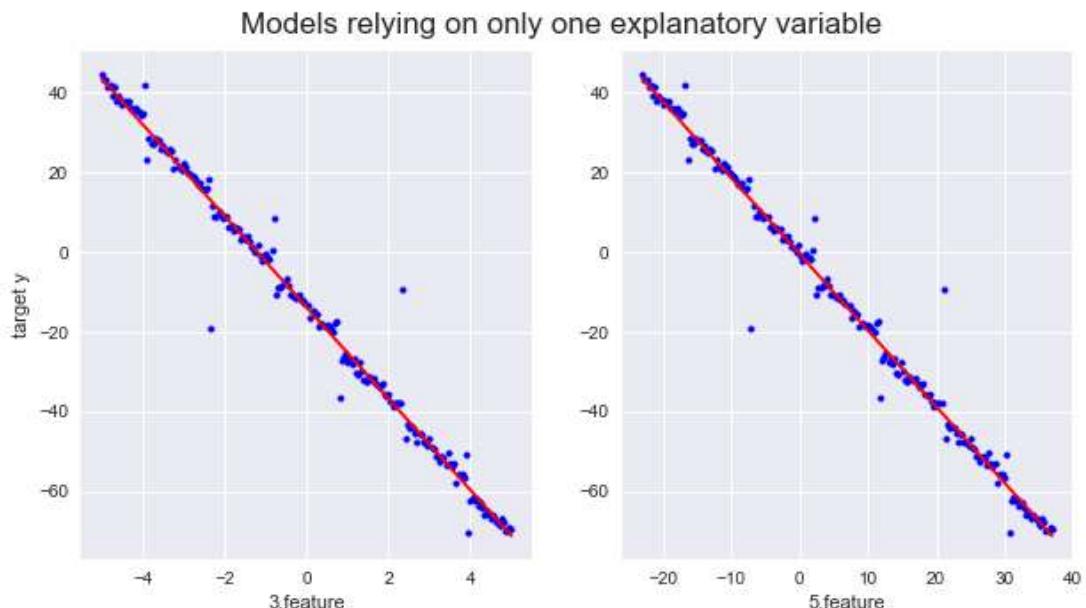
beta2 = beta_hat(X_extended[:,[0,7]],y)
pred2 = X_extended[:,[0,7]]@beta2
r_sq2 = r_sq(y, pred2)
corr_r_sq2 = corr_r_sq(y, pred2, 2)
print('\nResults of model using 5.feature and bias:')
print('- coefficient of determination: {:.6f}'.format(r_sq2))
print('- corrected coefficient of determination: {:.6f}'.format(corr_r_sq2))
```

```
Results of model using 3.feature and bias:
- coefficient of determination: 0.985808
- corrected coefficient of determination: 0.985664
```

```
Results of model using 5.feature and bias:
- coefficient of determination: 0.985808
- corrected coefficient of determination: 0.985664
```

```
[ ]: x_plot_1model = np.linspace(X[:,2].min(),X[:,2].max(),5)
x_plot_2model = np.linspace(X[:,6].min(),X[:,6].max(),5)
```

```
[ ]: fig, ax = plt.subplots(1, 2, figsize=(10,5))
fig.subplots_adjust(top=0.92)
ax[0].scatter(X[:,2],y, s=12, facecolors='blue')
ax[0].set_xlabel('3.feature')
ax[0].set_ylabel('target y')
ax[0].plot(x_plot_1model,beta1[0] + beta1[1]*x_plot_1model, color='red')
ax[1].scatter(X[:,6],y, s=12, facecolors='blue')
ax[1].set_xlabel('5.feature')
ax[1].plot(x_plot_2model,beta2[0] + beta2[1]*x_plot_2model, color='red')
plt.suptitle('Models relying on only one explanatory variable', fontsize=17, y=0.
             ↪98)
plt.show()
```



Both our models performed pretty good and achieved a high (corrected) coefficient of determination. Since the 3. and 5.feature are correlated, their scores are equal and no model is better than the other.

In the following we will try to further improve our model.

### 3.2 3.2 Second approach: Select features based on corrected coefficient of determination

Note: From now on we will use the expressions "1.,2., ..., 7.feature" to refer to the features. Otherwise we will explicitly use the term "original feature(s)".

example: \* original 3.feature -> 3.feature \* original 5.feature -> 7.feature

In the following we have used the corrected coefficient of determination to determine which features to use to build a good model. This is subdivided into two parts: 1. part: no bias, using all features 2. part: bias included, leaving out the 3.feature to ensure full column rank

Looking back on this work, the 2.part is the really important one. However, we also wanted to document here how our work progressed over time and we already made some useful observations during the 1.part, so we decided to not drop it.

#### 3.2.1 In the following we want to try each possible feature combination. Let's compute the powerset of our feature indexes:

```
[ ]: # the following function is copied from:  
# https://stackoverflow.com/questions/1482308/  
#how-to-get-all-subsets-of-a-set-powerset  
def powerset(s):  
    result = []  
    x = len(s)  
    masks = [1 << i for i in range(x)]  
    for i in range(1, 1 << x):  
        yield [ss for mask, ss in zip(masks, s) if i & mask]
```

```
[ ]: feature_comb = list(powerset(range(p)))  
print(feature_comb[:5])  
print(feature_comb[-5:])  
  
[[0], [1], [0, 1], [2], [0, 2]]  
[[0, 1, 3, 4, 5, 6], [2, 3, 4, 5, 6], [0, 2, 3, 4, 5, 6], [1, 2, 3, 4, 5, 6],  
[0, 1, 2, 3, 4, 5, 6]]
```

The list feature\_comb describes our features using the index (starting at 0), that's useful for the computation part. However for interpretation it is useful to start the counting at 1 for the 1.feature and so on. Therefore we create a new list:

```
[ ]: feature_comb_interp = []  
for f in feature_comb:  
    feature_comb_interp.append([i+1 for i in f])  
  
[ ]: print(feature_comb_interp[:5])  
print(feature_comb_interp[-5:])
```

```
[[1], [2], [1, 2], [3], [1, 3]]
[[1, 2, 4, 5, 6, 7], [3, 4, 5, 6, 7], [1, 3, 4, 5, 6, 7], [2, 3, 4, 5, 6, 7],
[1, 2, 3, 4, 5, 6, 7]]
```

### 3.2.2 1. Part: No bias, using all features

```
[ ]: r_sq_results = []
corr_r_sq_results = []
for f in feature_comb:
    beta = beta_hat(X[:,f],y)
    pred = X[:,f]@beta
    r_sq_results.append(r_sq(y, pred))
    corr_r_sq_results.append(corr_r_sq(y, pred, len(f)))
```

The list `r_sq_results` and `corr_r_sq_results` contain the coefficient of determination  $R^2$  and the corrected coefficient of determination for every feature combination model. The models are trained on the training data and evaluated on the test data.

We combine the results in a data frame. For this we use the notation where we start the counting of the features at 1 in order to fit to the plots created at the beginning.

```
[ ]: results = pd.DataFrame(data = np.vstack((np.array(r_sq_results), np.
array(corr_r_sq_results))).T,
                           index = [str(x) for x in feature_comb_interp],
                           columns=['R^2', 'corr R^2'])
```

Let's take a look at the feature combinations that performed the worst and at the feature combinations that performed the best.

1. Sorted according to the coefficient of determination
2. Sorted according to the corrected coefficient of determination

Worst feature combinations according to regular coefficient of determination:

```
[ ]: results.sort_values('R^2').head(10)
```

```
[ ]:          R^2  corr R^2
[2] -0.165994 -0.171883
[6] -0.163031 -0.168904
[2, 6] -0.158094 -0.169851
[5] -0.071306 -0.076717
[2, 5] -0.064315 -0.075120
[4] -0.055286 -0.060616
[2, 4] -0.050126 -0.060788
[1] -0.005058 -0.010134
[1, 4] -0.004969 -0.015171
[1, 2] 0.000372 -0.009776
```

Worst feature combinations according to corrected coefficient of determination:

```
[ ]: results.sort_values('corr R^2').head(10)
```

```
[ ]:          R^2  corr R^2
[2]      -0.165994 -0.171883
[2, 6]   -0.158094 -0.169851
[6]      -0.163031 -0.168904
[5]      -0.071306 -0.076717
[2, 5]   -0.064315 -0.075120
[2, 4]   -0.050126 -0.060788
[4]      -0.055286 -0.060616
[1, 4]   -0.004969 -0.015171
[1, 2, 4] 0.000452 -0.014847
[1]      -0.005058 -0.010134
```

Best feature combinations according to regular coefficient of determination:

```
[ ]: results.sort_values('R^2').tail(10)
```

Best feature combinations according to corrected coefficient of determination:

```
[ ]: results.sort_values('corr R^2').tail(10)
```

```
[ ]:          R^2  corr R^2
[1, 2, 6, 7] 0.999187 0.999170
[2, 3, 4, 7] 0.999188 0.999171
[1, 2, 3, 5, 7] 0.999192 0.999171
[1, 2, 3, 4, 7] 0.999192 0.999171
[2, 3, 5, 7] 0.999188 0.999171
[1, 2, 3, 6, 7] 0.999192 0.999172
[2, 3, 6, 7] 0.999189 0.999172
[1, 2, 7] 0.999186 0.999174
[2, 3, 7] 0.999188 0.999175
[1, 2, 3, 7] 0.999192 0.999175
```

```
[ ]: np.sum(results['R^2']>=0.9)
```

```
[ ]: 94
```

```
[ ]: np.sum(results['corr R^2']>=0.9)
```

```
[ ]: 94
```

Which feature combinations got a corrected coefficient of determination of at least 0.999 ?

```
[ ]: results[results['corr R^2']>=0.999].sort_values('corr R^2')
```

```
[ ]:          R^2  corr R^2
[2, 7]      0.999034 0.999024
[2, 4, 7]   0.999136 0.999123
```

[2, 4, 5, 6, 7]	0.999169	0.999148
[2, 5, 7]	0.999163	0.999150
[2, 5, 6, 7]	0.999168	0.999151
[2, 4, 6, 7]	0.999168	0.999151
[2, 4, 5, 7]	0.999169	0.999152
[2, 6, 7]	0.999168	0.999155
[1, 2, 4, 5, 6, 7]	0.999190	0.999165
[1, 2, 4, 5, 7]	0.999186	0.999165
[1, 2, 3, 4, 5, 6, 7]	0.999195	0.999166
[1, 2, 5, 6, 7]	0.999187	0.999167
[1, 2, 4, 6, 7]	0.999188	0.999167
[2, 3, 4, 5, 6, 7]	0.999192	0.999167
[1, 2, 3, 4, 5, 7]	0.999192	0.999167
[2, 3, 4, 5, 7]	0.999188	0.999167
[1, 2, 3, 4, 6, 7]	0.999193	0.999168
[1, 2, 3, 5, 6, 7]	0.999193	0.999168
[2, 3, 5, 6, 7]	0.999189	0.999169
[2, 3, 4, 6, 7]	0.999190	0.999169
[1, 2, 5, 7]	0.999186	0.999170
[1, 2, 4, 7]	0.999186	0.999170
[1, 2, 6, 7]	0.999187	0.999170
[2, 3, 4, 7]	0.999188	0.999171
[1, 2, 3, 5, 7]	0.999192	0.999171
[1, 2, 3, 4, 7]	0.999192	0.999171
[2, 3, 5, 7]	0.999188	0.999171
[1, 2, 3, 6, 7]	0.999192	0.999172
[2, 3, 6, 7]	0.999189	0.999172
[1, 2, 7]	0.999186	0.999174
[2, 3, 7]	0.999188	0.999175
[1, 2, 3, 7]	0.999192	0.999175

### 3.2.3 Discussion/Interpretation:

In order to simplify the explanations, we will sometimes use the word “metrics” if we refer to both the coefficient of determination and the adjusted coefficient of determination.

- As expected the models that used only one or two features didn’t perform well.
- From 127 feature combinations 94 have a coefficient of determination and an adjusted coefficient of determination of at least 0.9. Note: the best possible score is 1. We can interpret this as: Many feature combinations already include a lot of valuable information to build a good model. It is not necessary to use all features.
- The top “worst” feature combinations are quite similar for both metrics.
- In the top “best” feature combinations the results look different for the two metrics. The coefficient of determination is the highest when we use all features. The adjusted coefficient of determination suggests to use either the feature [1, 2, 3, 7] or [2, 3, 7] (i.e. without 1).
- This different observations for the two metrics fit their characteristics described above: The coefficient of determination is affected by the number of features used. By adding new features, the score will never decrease. The adjusted coefficient of determination is not pe-

nalizing a more complex model (i.e. a model using more features) and is thus more suitable for this kind of model comparison.

- If we take a closer look at the top 10 list (based on the adjusted coefficient of determination), we can see that the 2. and 7.feature (the 7.feature is in the original dataset the 5.feature) appear in all feature combinations. We can see in the last table that those two features are the only pair that got a corrected coefficient of determination above 0.999. Thus, using those features is for sure a good choice.
- To further improve the model it might be worth to include at least one additional feature, e.g. feature 1 or 3. (The combinations [1, 2, 7] and [2, 3, 7] can be found quite low in the results tables).
- It is interesting to see that even if the 3.feature and the 7.feature (original 5.feature) are correlated, the 7.feature appears more often than the 3.feature. A reason might be that the 7.feature has a bigger range of values than the 3.feature.

### 3.2.4 2. Part: Including bias, excluding 3.feature

We have seen above that our matrix including the extra column of 1's is not of full column rank. This can be explained by the correlation of the 3. and 7. feature. Leaving out either one of them and including the bias we get a matrix with full column rank.

```
[ ]: # This matrix doesn't have full column rank!
np.linalg.matrix_rank(X_extended[:,[0,3,7]])
```

```
[ ]: 2
```

```
[ ]: # leave out the 3.feature -> full column rank.
np.linalg.matrix_rank(X_extended[:,[0,1,2,4,5,6,7]])
```

```
[ ]: 7
```

```
[ ]: # leave out the 7.feature -> full column rank.
np.linalg.matrix_rank(X_extended[:,[0,1,2,3,4,5,6]])
```

```
[ ]: 7
```

Since we have seen above that the feature 7 was quite valuable for our models, we will now drop the 3.feature, add the column of 1's and repeat the evaluation done above:

```
[ ]: X_new = X_extended[:,[0,1,2,4,5,6,7]]
```

Again, we create a new list that corresponds to feature\_comb (index based) but which will allow for easier interpretation of the result table.

```
[ ]: feature_comb_interp_new = []
for f in feature_comb:
    feature_comb_interp_new.append([i if i<3 else i+1 for i in f])
```

```
[ ]: r_sq_results_new = []
corr_r_sq_results_new = []
for f in feature_comb:
    beta = beta_hat(X_new[:,f],y)
    pred = X_new[:,f]@beta
    r_sq_results_new.append(r_sq(y, pred))
    corr_r_sq_results_new.append(corr_r_sq(y, pred, len(f)))

[ ]: results_new = pd.DataFrame(data = np.vstack((np.array(r_sq_results_new), np.
    array(corr_r_sq_results_new))).T,
                                index = [str(x) for x in feature_comb_interp_new],
                                columns=['R^2', 'corr R^2'])
```

The table uses the same index as the first results table. 0 indicates that a bias was used, the digits from 1 to 7 indicate the feature (1. -> first feature).

```
[ ]: results_new.sort_values('corr R^2').tail(10)
```

```
[ ]:          R^2  corr R^2
[1, 2, 6, 7]  0.999187  0.999170
[0, 2, 4, 7]  0.999188  0.999171
[0, 1, 2, 5, 7]  0.999192  0.999171
[0, 1, 2, 4, 7]  0.999192  0.999171
[0, 2, 5, 7]  0.999188  0.999171
[0, 1, 2, 6, 7]  0.999192  0.999172
[0, 2, 6, 7]  0.999189  0.999172
[1, 2, 7]      0.999186  0.999174
[0, 2, 7]      0.999188  0.999175
[0, 1, 2, 7]  0.999192  0.999175
```

### 3.2.5 Discussion/Interpretation:

Similar to our first results the 2. and 7.feature appear in each feature combination in the top 10. Most interesting is probably the feature combination [0,2,7] (same score as [0,1,2,7]). This model is based on just the two features 2 and 7 plus a bias. To use it at application time we would need to collect only two values.

Because of the correlation of the 3. and the 7.feature, we can expect that replacing the 7. by the 3.feature (and remaining a bias!) should result in an equally good model:

```
[ ]: beta = beta_hat(X_extended[:,[0,2,3]],y)
pred = X_extended[:,[0,2,3]]@beta
print('Model: 2. and 3. feature including bias:')
print(r_sq(y, pred))
print(corr_r_sq(y, pred, 3))

beta = beta_hat(X_extended[:,[0,1,2,3]],y)
pred = X_extended[:,[0,1,2,3]]@beta
```

```

print('\nModel: 1., 2. and 3. feature including bias:')
print(r_sq(y, pred))
print(corr_r_sq(y, pred, 4))

```

Model: 2. and 3. feature including bias:  
0.9991878139546391  
0.9991753825355775

Model: 1., 2. and 3. feature including bias:  
0.999192060398189  
0.9991754872781519

As expected we get the same results for the (adjusted) coefficient of determination if we do the replacement.

### 3.3 Third approach: Hypothesis Testing

In the following we will apply our test to the model trained on all features except for the 3.feature and we will use a bias term.

References:

[http://reliawiki.org/index.php/Multiple\\_Linear\\_Regression\\_Analysis](http://reliawiki.org/index.php/Multiple_Linear_Regression_Analysis)

<https://online.stat.psu.edu/stat462/node/131/>

Given our multiple linear regression model trained with all the available features (without the bias term) we might ask our selves: “Is the  $i$ -th feature contributing much to  $y$ ? Or is it just correcting minor noise?”

To test it, we can use hypothesis testing on the individual factors  $\hat{\beta}_i$  with the following test:

$$H_0 : \hat{\beta}_i = 0 \quad H_A : \hat{\beta}_i \neq 0$$

The go-to test statistic for this setup (one sample, normally distributed errors and “big” sample size of 200 points) is the **one-sample Student's t-test**. The [history of this name](#) is interesting and involves Guinness, so check it if you like beer.

The general formulation of the t-test statistic is:

$$t = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}}$$

For our case we don't have samples of  $\beta_i$ , but only one single occurrence and its variance. And we have  $\mu = 0$ , so we can reformulate the test statistic as follows:

$$t_i = \frac{\hat{\beta}_i}{se(\hat{\beta}_i)}$$

Where  $se(\hat{\beta}_i)$  is the standard error of the factor  $\hat{\beta}_i$ . This is the square root of the  $i$ -th diagonal element of the covariance matrix  $C$ :

$$C = \hat{\sigma}(X^T X)^{-1}$$

Finally  $\hat{\sigma}$  is the variance of the predictions,  $\hat{\sigma} = \frac{MSE}{(n-p)}$ . Note that we divide by  $(n - p)$  and not  $(n - p - 1)$ , because our  $p$  is here defined as the number of columns in our input matrix (including the column of 1's) and not by the number of features.

With this, we can proceed to test which features are significant. We work with the standard significance level of  $\alpha = 0.05$

Step 1: Define constants

```
[ ]: n = X_new.shape[0]
      p = X_new.shape[1] # <- p is not defined by the number of features, but by the
      ↪number of columns
      # significance level
      alpha = 0.05
```

Step 2: Train model

```
[ ]: beta = beta_hat(X_new, y)
      pred = X_new@beta
      print(beta)

[-2.73818683  0.18198005  0.4016157   -0.06941017 -0.14394245  0.17572647
 -1.92418342]
```

Step 3: Estimate the variance and compute the variance-covariance matrix of the estimated regression coefficients

```
[ ]: est_var = np.sum((y - pred)**2) / (n - p)
      # variance-covariance matrix of the estimated regression coefficients
      C = est_var*inv(X_new.T@X_new)
```

Step 4: Define t-test function for our scenario

```
[ ]: from scipy.stats import t

def t_test(beta, se_beta, alpha, df) -> str:

    t_stat = beta / se_beta
    t_bound = t.ppf(1. - alpha, df=df)

    if abs(t_stat) > t_bound:
        print(f"""
              t-test result:
              REJECT null hypothesis (mean 0) with t-value {t_stat}
              """)
        return "REJECT"
    else:
        print(f"""
              t-test result:
              FAIL TO REJECT null hypothesis (mean 0) with t-value {t_stat}
              """)
        return "FAIL TO REJECT"
```

```
[ ]: for i in range(0,7):
    if i>=3:
        f = i+1
    else:
        f = i
    if i != 0:
        print('Test of ' + str(f) + '.feature:')
        t_test(beta[i], np.sqrt(C[i,i]), alpha, (n - p))
```

Test of 1.feature:  
t-test result:  
FAIL TO REJECT null hypothesis (mean 0) with t-value 0.8599102751778964

Test of 2.feature:  
t-test result:  
REJECT null hypothesis (mean 0) with t-value 55.53469904349027

Test of 4.feature:  
t-test result:  
FAIL TO REJECT null hypothesis (mean 0) with t-value -0.6712261591596251

Test of 5.feature:  
t-test result:  
FAIL TO REJECT null hypothesis (mean 0) with t-value -0.7227015207511889

Test of 6.feature:  
t-test result:  
FAIL TO REJECT null hypothesis (mean 0) with t-value 0.8202076221553543

Test of 7.feature:  
t-test result:  
REJECT null hypothesis (mean 0) with t-value -213.8776281826563

So we would only reject the  $H_0 : \hat{\beta}_i = 0$  for feature 2 and 7. Therefore, the regression model should take into account at least these features and might be able to omit the rest. (Note: the 3.feature and the bias were excluded from this analysis)

$$y = \hat{\beta}_0 + \hat{\beta}_2 X_2 + \hat{\beta}_7 X_7$$

We can go one step further and look at the (absolute) value of the t-statistic. The meaning of it is: "the lower, the less significant the feature". So we can say that the 4.feature is the least significant, closely followed by the 5.feature.

## 4 4 Conclusion

We have seen in our first approach that both the 3. and the original 5.feature can already be used to build a good model. The second approach has showed that we should continue either with the 3. or the 5.feature and combine it with the 2.feature and add a bias. During hypothesis testing we have seen that the only regression coefficients for which we could reject the null hypothesis that these coefficients are zero were those of the 2. and the 5. feature. (3.feature was not included during test).

All these findings support our final conclusion that the best model would be based on the 2. and 5. feature (or alternatively the 3.feature) combined with bias. The respective corrected coefficient of correlation is 0.999175.

We would like to point out that the feature selection cannot be based solely on these results. What these features actually are should also be taken into account. For example assume that it is expensive or time intensive to measure the 2.feature. Then, it might better to just use the 3. or 5. feature like we have seen in our 1.approach and to tolerate a model that is slightly worse than the best possible one.

### 4.1 Annex: Complex Matrices

As commented at the very beginning of this notebook, there are other takes on how to deal with the complex feature. Worth of special mention is that in which all features are casted into complex format and the Least Squares estimation is performed on the Complex Matrix. See below a summary of the reasoning, the differences to simple Least Squares and the resulting coefficients of determination. To see the full implementation please refer to the individual submissions from Max Serra and Ignacio Llorca.

#### 4.1.1 Reasoning

Let the module of the imaginary number  $z = x + yi$  be  $|z| = \sqrt{x^2 + y^2}$ , and its principal argument  $\arg(z) = \tan^{-1}(\frac{y}{x})$ .

In our regression with complex values  $Y_i = \beta_0 + \beta_1 X_{i,1} + \beta_2 X_{i,2} + \beta_3 X_{i,3} + \beta_4 X_{i,4} + \beta_5 X_{i,5}$  multiplication by  $\beta_p$  will rescale the independent variable  $X_{i,p}$  by the modulus  $|\beta_p|$  and rotate it around the origin by  $\arg(\beta_p)$ .

The LS-estimator  $\hat{\beta}$  that we want to find can be calculated as in the simple-number case with  $\hat{\beta} = (X^T X)^{-1} X^T y$ , only that the transpose operation  $X^T$  is now replaced with the conjugate transpose  $X^* = \bar{X}^T$ , where conjugating means inverting the sign of the imaginary part of the complex number, i.e:  $\bar{z} = x - yi$ . Thus,  $\hat{\beta} = (X^* X)^{-1} X^* y$

#### 4.1.2 Data preprocessing

After dealing with outliers and making necessary transformations to cast into complex format, the feature matrix looks like this:

	0	1	2	3	4
0	11.728000+0.000000j	-0.648360+0.000000j	-5.000000+0.000000j	2.302600+3.141600j	-23.000000+0.000000j
1	11.771000+0.000000j	-0.514550+0.000000j	-4.950000+0.000000j	2.292500+3.141600j	-22.700000+0.000000j
2	11.817000+0.000000j	-0.393880+0.000000j	-4.900000+0.000000j	2.282400+3.141600j	-22.400000+0.000000j
3	11.864000+0.000000j	-0.282390+0.000000j	-4.850000+0.000000j	2.272100+3.141600j	-22.100000+0.000000j
4	11.913000+0.000000j	-0.177040+0.000000j	-4.800000+0.000000j	2.261800+3.141600j	-21.800000+0.000000j
...	...	...	...	...	...
196	12.087000+0.000000j	0.177040+0.000000j	4.800000+0.000000j	2.261800+0.000000j	35.800000+0.000000j
197	12.136000+0.000000j	0.282390+0.000000j	4.850000+0.000000j	2.272100+0.000000j	36.100000+0.000000j
198	12.183000+0.000000j	0.393880+0.000000j	4.900000+0.000000j	2.282400+0.000000j	36.400000+0.000000j
199	12.229000+0.000000j	0.514550+0.000000j	4.950000+0.000000j	2.292500+0.000000j	36.700000+0.000000j
200	12.272000+0.000000j	0.648360+0.000000j	5.000000+0.000000j	2.302600+0.000000j	37.000000+0.000000j

201 rows × 5 columns

With conjugate transpose looking like:

	0	1	2	3	4	5	6
0	11.728000-0.000000j	11.771000-0.000000j	11.817000-0.000000j	11.864000-0.000000j	11.913000-0.000000j	11.962000-0.000000j	12.012000-0.000000j
1	-0.648360-0.000000j	-0.514550-0.000000j	-0.393880-0.000000j	-0.282390-0.000000j	-0.177040-0.000000j	-0.075364-0.000000j	0.024783-0.000000j
2	-5.000000-0.000000j	-4.950000-0.000000j	-4.900000-0.000000j	-4.850000-0.000000j	-4.800000-0.000000j	-4.750000-0.000000j	-4.700000-0.000000j
3	2.302600-3.141600j	2.292500-3.141600j	2.282400-3.141600j	2.272100-3.141600j	2.261800-3.141600j	2.251300-3.141600j	2.240700-3.141600j
4	-23.000000-0.000000j	-22.700000-0.000000j	-22.400000-0.000000j	-22.100000-0.000000j	-21.800000-0.000000j	-21.500000-0.000000j	-21.200000-0.000000j

5 rows × 201 columns

#### 4.1.3 Least Squares Estimator

As explained above, the function to obtain  $\beta$  now looks like:

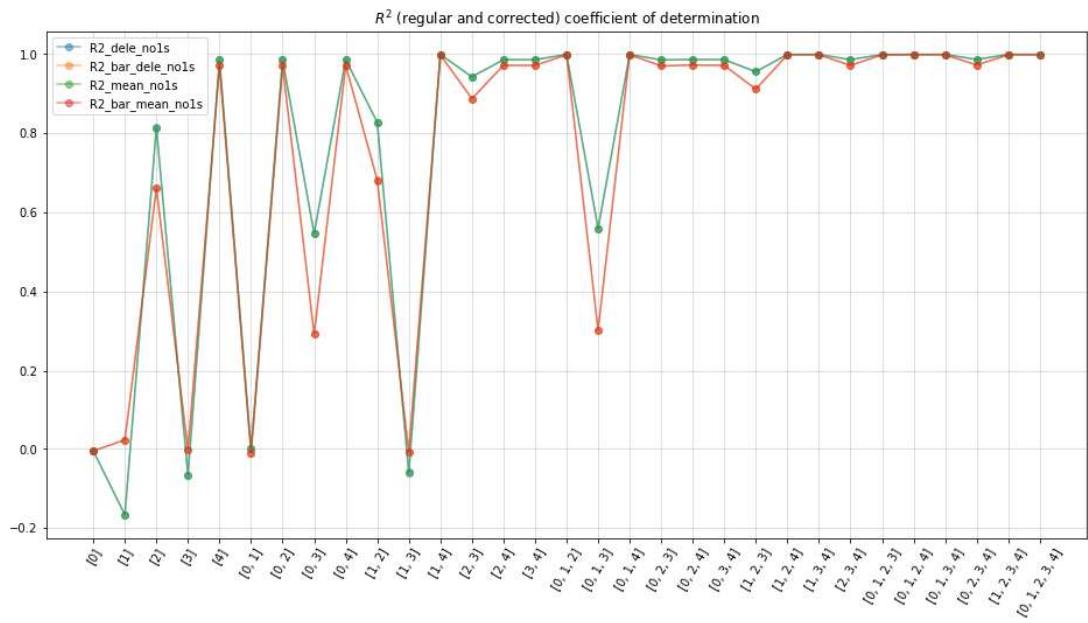
```
[ ]: def least_squares_complex(X,y):
    # np.linalg.inv outputs the inverse of a matrix
    # .H outputs the conjugate transpose as explained above = .conj().T
    return np.linalg.inv(X.H @ X) @ X.H @ y
```

#### 4.1.4 Results

The feature combinations and the coefficients of determination  $R^2$  and its corrected counterpart  $\bar{R}^2$  can be calculated just the same as in the simple-matrix method.

Let “deletion” and “mean” denote the technique used to deal with the outlier in row 100 from the feature matrix. The coefficients of determination shed the following results:

	R2_bar_deletion	R2_bar_mean	R2_deletion	R2_mean
[0]	-0.005025	-0.004999	-0.005058	-0.005059
[1]	0.022643	0.022990	-0.165994	-0.166959
[2]	0.661964	0.660387	0.814655	0.813686
[3]	-0.000708	-0.000487	-0.065730	-0.067200
[4]	0.971270	0.971255	0.985603	0.985596
[0, 1]	-0.010152	-0.010101	0.000372	0.000370
[0, 2]	0.970253	0.970247	0.985166	0.985162
[0, 3]	0.291900	0.291933	0.546824	0.546822
[0, 4]	0.971495	0.971489	0.985790	0.985786
[1, 2]	0.682449	0.680846	0.828034	0.827066
[1, 3]	-0.006591	-0.006361	-0.059372	-0.060851
[1, 4]	0.998049	0.998034	0.999034	0.999026
[2, 3]	0.887344	0.887185	0.942590	0.942504
[2, 4]	0.971532	0.971527	0.985808	0.985805
[3, 4]	0.971486	0.971479	0.985785	0.985781
[0, 1, 2]	0.997985	0.997978	0.999007	0.999003
[0, 1, 3]	0.302013	0.302064	0.559049	0.559046
[0, 1, 4]	0.998348	0.998341	0.999186	0.999183
[0, 2, 3]	0.970109	0.970104	0.985170	0.985166
[0, 2, 4]	0.972322	0.972317	0.986275	0.986272
[0, 3, 4]	0.971368	0.971363	0.985799	0.985795
[1, 2, 3]	0.911435	0.911273	0.955390	0.955303
[1, 2, 4]	0.998351	0.998344	0.999188	0.999184
[1, 3, 4]	0.998282	0.998274	0.999154	0.999149
[2, 3, 4]	0.971394	0.971389	0.985812	0.985808
[0, 1, 2, 3]	0.997979	0.997972	0.999009	0.999006
[0, 1, 2, 4]	0.998352	0.998345	0.999192	0.999189
[0, 1, 3, 4]	0.998340	0.998333	0.999186	0.999183
[0, 2, 3, 4]	0.972224	0.972220	0.986297	0.986294
[1, 2, 3, 4]	0.998343	0.998336	0.999188	0.999184
[0, 1, 2, 3, 4]	0.998343	0.998336	0.999192	0.999189



#### 4.1.5 Conclusions

The above results and plot can be used to infer the most meaningful combinations of features, which seem to be aligned to what has been concluded in the simple-matrix method. We can see, for instance, how the 5th feature alone (index 4) or the combination of the 2nd and 5th (index 1 and 5), reach coefficients as good as some the 4 and 5-feature combinations.

[ ]: