

Week 4 - Partial Least Squares Regression

Please indicate at the top of your assignment whether or not you used any AI tools, such as MS Copilot. If you did use one of these tools, please provide a very brief explanation alongside each answer for how you confirmed the correctness of your solution.

- Used MS Copilot for a couple of questions.

We will reimplement and then explore some of the properties of [Cosgrove et al \(http://pubs.rsc.org/en/Content/ArticleLanding/2010/MB/b926287c\)](http://pubs.rsc.org/en/Content/ArticleLanding/2010/MB/b926287c).

In [3]:

```
1 import scipy as sp, numpy as np
2 from sklearn.preprocessing import scale, StandardScaler
3 from sklearn.cross_decomposition import PLSRegression
4 from sklearn.model_selection import LeaveOneGroupOut, LeaveOneOut
5 import matplotlib.pyplot as plt
6 import pandas as pd
7
8 data = sp.io.loadmat('wk4_Cosgrove-data.mat', squeeze_me=True)['s']
9
10 X = data['X'].item() # the untransformed data matrix (66x102)
11 Y = data['Y'].item() # the untransformed LDH release at 48hours. (66x1)
12 phosphoproteins = data['phosphoproteins'].item() # names of phosphoproteins
13 conditions = data['conditions'].item() # cell array of the 66 conditions
14 drugList = data['drugList'].item() # description of the drugs used in each of the 66 conditions
15 drugListToxic = data['drugListToxic'].item() # binary value corresponding to whether drugList[i] is toxic
16 drugs = data['drugs'].item() # binary matrix mapping which measurements correspond to a drug treatment in drugList
17 cytokineList = data['cytokineList'].item() # cell array of cytokine treatments
18 ind4pProtein = data['ind4pProtein'].item() # the column indices corresponding to measurements of the 4 phosphoproteins
```

In [4]:

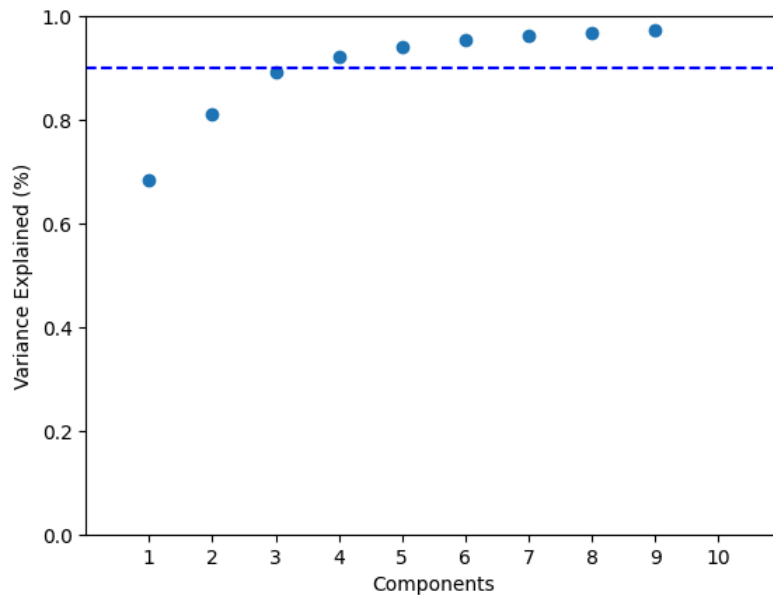
```
1 print('X:', X.shape, 'Y:', Y.shape)
2 print('conditions:', conditions)
3 print('drugList:', drugList)
4 print('drugListToxic:', drugListToxic)
5 print('drugs:', drugs.shape)
6 print('cytokineList:', cytokineList)
7 print('ind4pProtein:', ind4pProtein)

X: (66, 102) Y: (66, 5)
conditions: ['DMSO_NoCyt' 'DMSO_IL-1' 'DMSO_LPS' 'DMSO_TNF' 'DMSO_IL-6' 'DMSO_Mix'
'CIM_NoCyt' 'CIM_IL-1' 'CIM_LPS' 'CIM_TNF' 'CIM_IL-6' 'CIM_Mix'
'RAN_NoCyt' 'RAN_IL-1' 'RAN_LPS' 'RAN_TNF' 'RAN_IL-6' 'RAN_Mix'
'LEV_NoCyt' 'LEV_IL-1' 'LEV_LPS' 'LEV_TNF' 'LEV_IL-6' 'LEV_Mix'
'TRO_NoCyt' 'TRO_IL-1' 'TRO_LPS' 'TRO_TNF' 'TRO_IL-6' 'TRO_Mix'
'BUS_NoCyt' 'BUS_IL-1' 'BUS_LPS' 'BUS_TNF' 'BUS_IL-6' 'BUS_Mix'
'NEF_NoCyt' 'NEF_IL-1' 'NEF_LPS' 'NEF_TNF' 'NEF_IL-6' 'NEF_Mix'
'ASP_NoCyt' 'ASP_IL-1' 'ASP_LPS' 'ASP_TNF' 'ASP_IL-6' 'ASP_Mix'
'NIM_NoCyt' 'NIM_IL-1' 'NIM_LPS' 'NIM_TNF' 'NIM_IL-6' 'NIM_Mix'
'CLA_NoCyt' 'CLA_IL-1' 'CLA_LPS' 'CLA_TNF' 'CLA_IL-6' 'CLA_Mix'
'TEL_NoCyt' 'TEL_IL-1' 'TEL_LPS' 'TEL_TNF' 'TEL_IL-6' 'TEL_Mix']
drugList: ['DMSO' 'CIM' 'RAN' 'LEV' 'TRO' 'BUS' 'NEF' 'ASP' 'NIM' 'CLA' 'TEL']
drugListToxic: [0 0 1 0 1 0 1 0 1 1 1]
drugs: (66, 11)
cytokineList: ['NoCyt' 'IL-1' 'LPS' 'TNF' 'IL-6' 'Mix']
ind4pProtein: [ 7  8  9 10 11 12 31 32 33 34 35 36 37 38 39 40 41 42 85 86 87 88 89 90]
```

(1) Perform PLSR on the matrixes X and Y. Plot the percent variance explained. How many principal components do you need for each to explain 90% of the Y variance? Discuss your findings.

Hint: Be sure you are normalizing each dataset as needed for the analysis.

```
In [6]: 1 # Note: PLSR scales the data by default
2 #X_norm = StandardScaler().fit_transform(X)
3 #Y_norm = StandardScaler().fit_transform(Y)
4
5 scores = []
6
7 for i in range(1, 10):
8     pls = PLSRegression(n_components=i).fit(X, Y[:, -1]) # average of the late time-points used for prediction
9     scores.append(pls.score(X, Y[:, -1]))
10
11
12 plt.scatter(range(1,10), scores)
13 plt.ylabel('Variance Explained (%)')
14 plt.xlabel('Components')
15 plt.plot([0, 11], [0.9, 0.9], '--', color='blue')
16 plt.ylim(0, 1)
17 plt.xlim(0, 11)
18 plt.xticks(np.arange(1, 11, 1))
19 plt.show()
20
```



We would need 4 components to explain 90% of the variance observed under the average of the late time-points 4–48 hrs. We observe that the first component explains almost 68% of variance in Y. The consecutive PCs explain less and less variability eventually approaching 1.

(2) How would you expect the percent of X variance explained to compare between PLSR and PCA? Why?

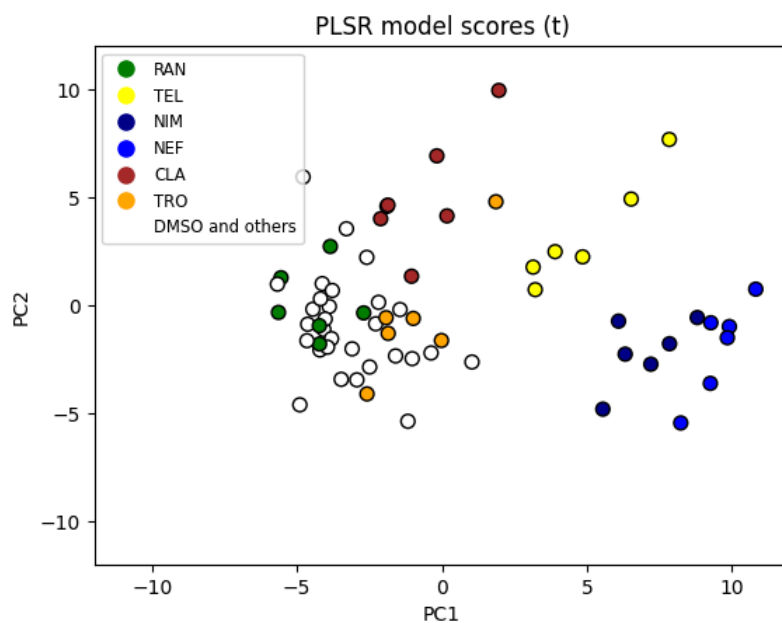
The percent of X variance explained in PCR would be higher compared to PLSR. PLSR focuses on maximizing covariance between X and Y, whereas PCA focuses on maximizing variance in X. Therefore 4 components are more than enough to explain 90% of the covariance in the dataset. PCA would effectively explain the maximum variability observed in X, whereas PLSR would not.

(3) Recreate the Figure S2A plot from Cosgrove et al. supplementary information. This is the PLSR scores plot (for PC1 and PC2), with toxic drugs colored according to the drug type and all other drugs are not colored. Use the drugList, drugListToxic to identify these categories.

```

In [7]: 1 plsr = PLSRegression(n_components=4).fit(X, Y[:, -1])
2 scores = plsr.x_scores_
3 pcs = scores[:, :2]
4
5 drug_colors = {
6     'RAN': 'green',
7     'TEL': 'yellow',
8     'NIM': 'darkblue',
9     'NEF': 'blue',
10    'CLA': 'brown',
11    'TRO': 'orange',
12    'DMSO and others': 'white'
13 }
14
15 colors = [
16     drug_colors[drug] if (np.any(drugs[i] == 1) and (drug := drugList[np.argwhere(drugs[i] == 1)][0][0]) in drug_colors)
17     else 'white'
18     for i in range(len(scores))
19 ]
20
21 # Scatter plot with legend
22 plt.scatter(x=pcs[:, 0], y=pcs[:, 1], c=colors, s=50, edgecolor='k')
23 plt.xlim(-12, 12)
24 plt.ylim(-12, 12)
25 plt.legend(handles=[plt.Line2D([0], [0], marker='o', color='w', markerfacecolor=color, markersize=10, label=drug) for drug in drug_colors],
26            loc='upper left',
27            fontsize='small',
28            bbox_to_anchor=(0, 1))
29
30 # Labels and title
31 plt.xlabel('PC1')
32 plt.ylabel('PC2')
33 plt.title('PLSR model scores (t)')
34 plt.show()
35
36

```



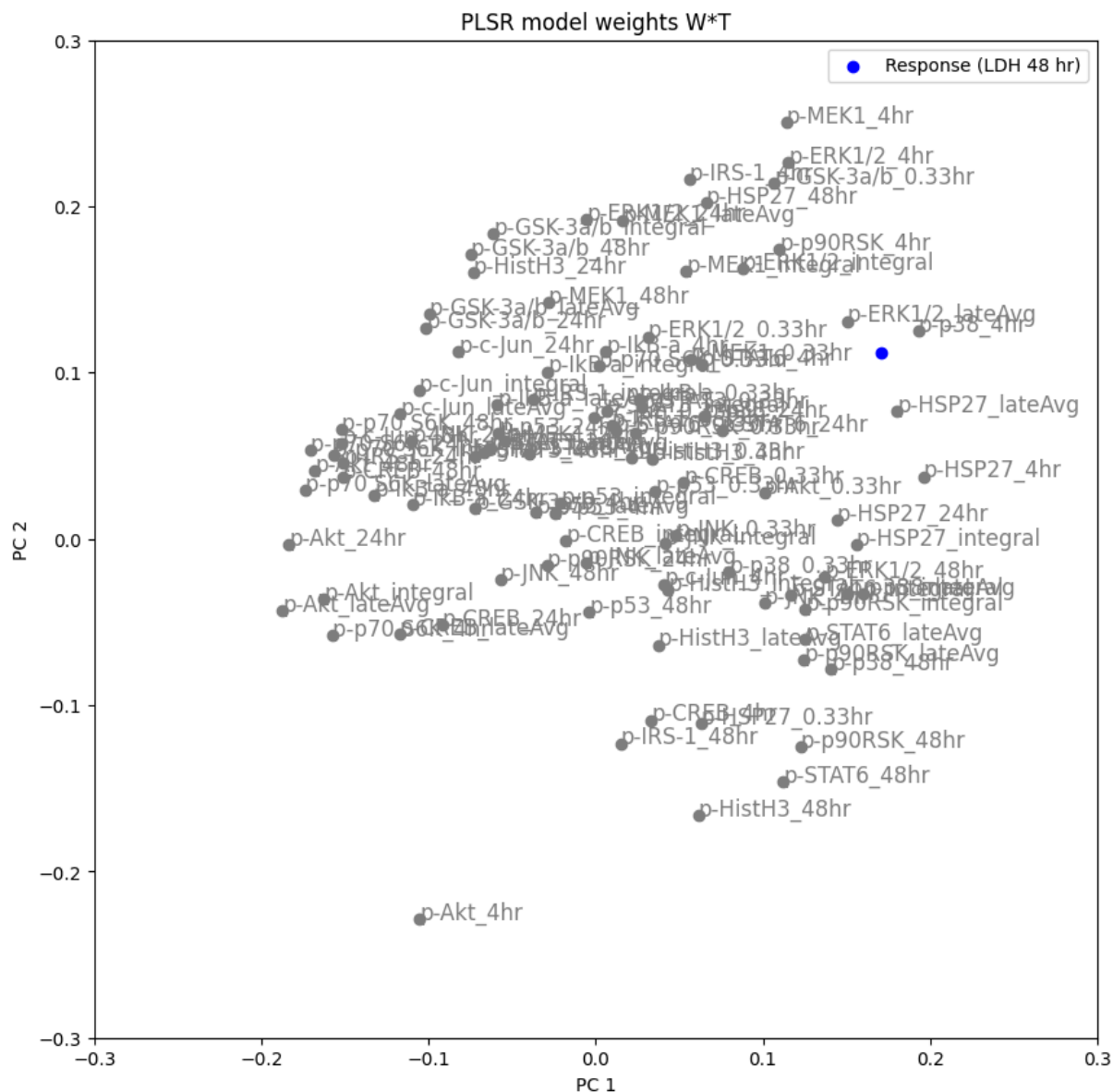
(4) Create the loadings plot corresponding to (3). Interpret the results shown on the plot.

```

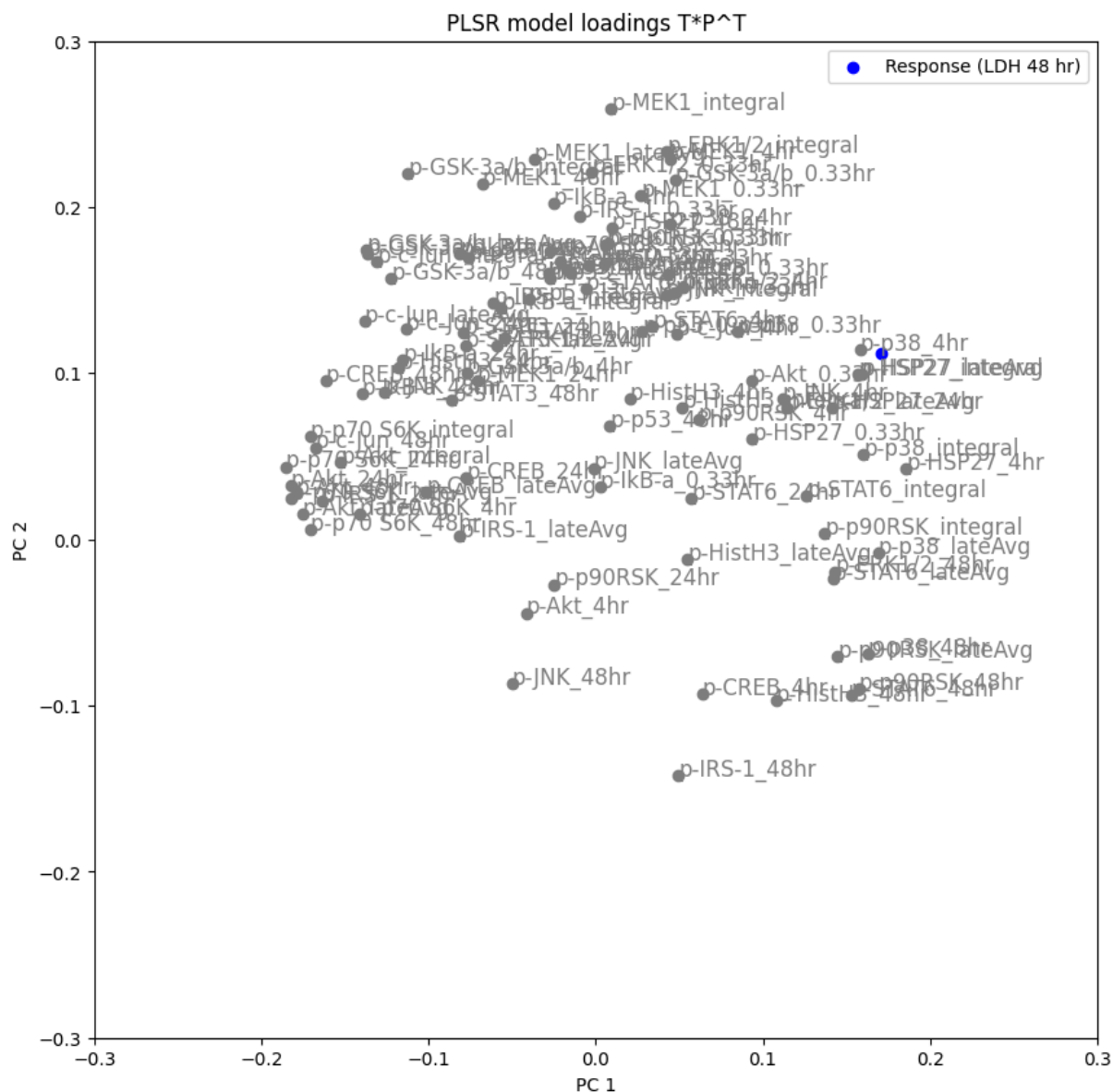
1 weights = plsrx_weights[:, :2]
2 response = plsry_weights_
3
4 print(weights.shape, response.shape)
5
6 plt.figure(figsize=(10,10))
7 plt.scatter(weights[:, 0], weights[:, 1], color='grey')
8 plt.scatter(response[:, 0], response[:, 1], color='blue', label= 'Response (LDH 48 hr)')
9
10 for i, feature in enumerate(phosphoproteins):
11     plt.text(weights[i, 0], weights[i, 1], feature, fontsize=12, color='grey')
12
13 # The six phosphoproteins (p-MEK1, p-ERK1/2, p-Akt, p-70 S6K, p-p38, p-HSP27) should surround the response variable
14
15 plt.legend()
16 plt.xlim(-0.3, 0.3)
17 plt.ylim(-0.3, 0.3)
18 plt.xlabel('PC 1')
19 plt.ylabel('PC 2')
20 plt.title('PLSR model weights W*T')
21 plt.show()
22

```

(102, 2) (1, 4)



```
In [9]: 1 loadings = plsrx.loadings[:, :2]
2 response = plsry.loadings_
3
4 print(loadings.shape, response.shape)
5
6 plt.figure(figsize=(10,10))
7 plt.scatter(loadings[:, 0], loadings[:, 1], color='grey')
8 plt.scatter(response[:, 0], response[:, 1], color='blue', label= 'Response (LDH 48 hr)')
9
10 for i, feature in enumerate(phosphoproteins):
11     plt.text(loadings[i, 0], loadings[i, 1], feature, fontsize=12, color='grey')
12
13 plt.legend()
14 plt.xlim(-0.3, 0.3)
15 plt.ylim(-0.3, 0.3)
16 plt.xlabel('PC 1')
17 plt.ylabel('PC 2')
18 plt.title('PLSR model loadings T*P*T')
19 plt.show()
```



We use loadings to show how each variable (phosphoproteins and LDH release) contributes to each component in the PCA space. They can be thought of as correlations between the original data and the components.

The weights (W) define the directions of the latent space that best predict Y . Weights are used to calculate scores.

- $T = X * W$

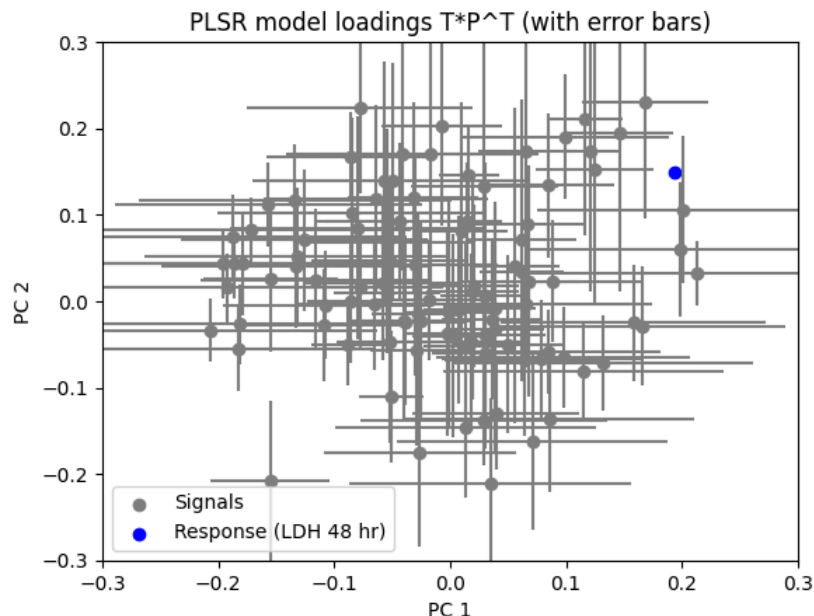
The loadings (P) allow reconstructing X from the scores, explaining how much variance of X is captured by each component. Loadings show how variables relate to components.

- $X \approx T * \text{transpose}(P)$

(5) Add the variance of the loadings to your loadings plot (this can be shown as error bars). How does the variance of component one compare to that of component two? Would you expect a trend in the general variance versus component number?

```
In [13]: 1 from sklearn.utils import resample
2 import pandas as pd
3
4 loadings = pls.x_weights[:, :2]
5 response = pls.y_weights_
6
7 print(loadings.shape, response.shape)
8
9
10 plt.scatter(loadings[:, 0], loadings[:, 1], color='grey', label = "Signals")
11 plt.scatter(response[:, 0], response[:, 1], color='blue', label= 'Response (LDH 48 hr)')
12 plt.legend()
13 plt.xlim(-0.3, 0.3)
14 plt.ylim(-0.3, 0.3)
15 plt.xlabel('PC 1')
16 plt.ylabel('PC 2')
17 plt.title('PLSR model loadings T*P^T (with error bars)')
18
19
20 output_data = []
21 n_bootstraps = 1000
22
23 for i in range(n_bootstraps):
24     x_boot, y_boot = resample(X, Y[:, -1], replace=True)
25     pls = PLSRegression(n_components=2).fit(x_boot, y_boot)
26     output_data.append({
27         'X_predict' : pls.predict(X),
28         'X_loadings' : pls.x_loadings[:, :2],
29         'Y_loadings' : pls.y_loadings_,
30         'Scores' : pls.x_scores_
31     })
32 )
33
34 output_df = pd.DataFrame(output_data)
35
36 vars = np.std(np.array(output_df['X_loadings']), axis=0)
37 plt.errorbar(loadings[:, 0], loadings[:, 1], xerr=vars[:, 0], yerr=vars[:, 1], linestyle='', color='grey')
38 plt.show()
```

(102, 2) (1, 2)



The variance of PC1 is generally higher compared to PC2. We should expect to see lower variance with increasing component numbers because the consecutive PCs are built with the residuals from the previous component. According to the paper, the positive model loadings imply signaling metrics with pro-death contributions, and negative model loadings imply pro-survival contributions. There is no direct correlation with variance, because we optimize for co-variance

(6) Recreate a 4-principal component model using PLSR with just the final 4 phosphoprotein model. Plot predicted v. observed LDH for this model. Report the model fitness (R^2). Define here how you are calculating R^2 .

```

In [86]: 1 from sklearn.metrics import r2_score
2
3 '''
4 Metrics with variable importance score (VIP) score >1 have significant importance in the model whereas metrics with
5 significantly lack unique information in the model. Inspection of model loadings and VIP scores identified
6 four signaling pathways (MEK-ERK, Akt, p70 S6K, and p38-HSP27) with phosphoproteins having informative model
7 contributions from four or more signaling metrics.
8
9 5(A) An OPLSR model of LDH release was trained on the time-dependent signaling metrics from
10 six phosphoproteins (p-MEK1, p-ERK1/2, p-Akt, p-70 S6K, p-p38, p-HSP27) using data from donor #1.
11 This trained model was used to predict LDH release responses based on the same signaling metrics measured in hepato
12 A correlation plot relating the observed and predicted LDH release responses for the both training data (donor #1)
13 which contained six treatments present and eight treatments not present in the training data. Test conditions (CHL
14 in the training data but nonetheless predicted accurately are noted. The model-predicted responses are presented as
15 mean prediction  $\pm$  cross-validation standard error, calculated by jack-knifing. The test data are presented as mean
16 three biological replicates. Experimental and prediction uncertainties are not shown for the training data.
17
18 '''
19 print(ind4pProtein)
20 new_X = X[:, ind4pProtein]
21
22 plsr = PLSRegression(n_components=4)
23 plsr.fit(new_X, Y[:, -1])
24 predicted = plsr.predict(new_X)
25 score = plsr.score(new_X, Y[:, -1])
26
27 plt.figure(figsize=(5, 5))
28 plt.ylabel('Predicted LDH release')
29 plt.xlabel('Observed LDH release')
30 plt.title('6-phosphoproteins PLSR model predictions')
31
32 plt.plot([0, 10], [0, 10], '--r')
33 plt.scatter(Y[:, -1], predicted, color='k')
34 plt.show()
35
36 print(f'Model Fiting (R-squared): {score}')

```

```
[ 7  8  9 10 11 12 31 32 33 34 35 36 37 38 39 40 41 42 85 86 87 88 89 90]
```

/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "orientation" which is no longer supported as of 3.3 and will become an error two minor releases later

```
fig.canvas.print_figure(bytes_io, **kw)
```

/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "dpi" which is no longer supported as of 3.3 and will become an error two minor releases later

```
fig.canvas.print_figure(bytes_io, **kw)
```

/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "facecolor" which is no longer supported as of 3.3 and will become an error two minor releases later

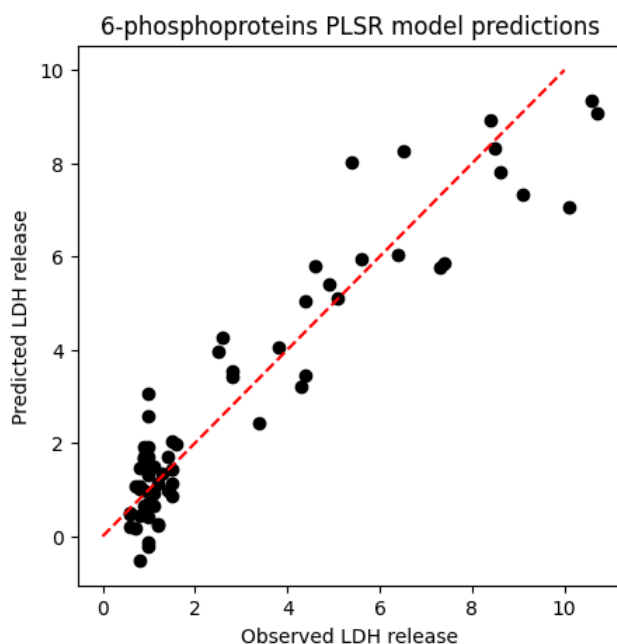
```
fig.canvas.print_figure(bytes_io, **kw)
```

/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "edgecolor" which is no longer supported as of 3.3 and will become an error two minor releases later

```
fig.canvas.print_figure(bytes_io, **kw)
```

/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: savefig() got unexpected keyword argument "bbox_inches_restore" which is no longer supported as of 3.3 and will become an error two minor releases later

```
fig.canvas.print_figure(bytes_io, **kw)
```



Model Fiting (R-squared): 0.8821898904528919

R-squared was calculated by fitting our PLSR model to X and Y matrices. Mathematically, it is 1 minus the sum of squares residuals divided by the total sum of squares. The model predicts Y from X. We calculate R² value with sklearn's r2_score function which uses the predicted and measured outcomes to inform about the goodness of fit.

(7) Cosgrove *et al* discusses their method for model validation using leave-one-out cross-validation. Calculate all LDH predictions for leave-one-out cross-validation and calculate the R² value for the resulting yfit values. What is the R² value? Why do you think it's important to perform cross-validation?

```
In [90]: 1 from sklearn.model_selection import cross_val_predict
2 predicted = cross_val_predict(plsr, X, Y[:, -1], cv=LeaveOneOut())
3 print(f'Model Fiting (R-squared): {r2_score(Y[:, -1], predicted)}')
```

Model Fiting (R-squared): 0.847813385573984

The R² value tells us about the goodness of fit. Here, our R² value has decreased from 0.88 to 0.85 after leave-one-out-CV. It is important to cross-validate data in order to see how well our model is able to generalize. In our previous approach, we included the testing set in the training process leading to an overestimation by the model. In CV, we avoid this by splitting our data into test and training sets which enables an unbiased training process.

(8) Now, instead of performing LOOCV, let's perform leave-one-cytokine-out cross-validation. That is, one cytokine at a time, leave out all the data for the NoCyt, IL-1, LPS, TNF, IL-6, or Mix conditions.

Hint: Look at sklearn.model_selection.LeaveOneGroupOut .

How does this affect your cross-validation? How do the two approaches here differ? When might each be most appropriate?

```
In [92]: 1 from sklearn.model_selection import LeaveOneGroupOut
2
3 print(cytokineList)
4 print(conditions)
5
6 groups = [next(i for i, cytokine in enumerate(cytokineList) if cytokine in cond) for cond in conditions]
7 predicted = cross_val_predict(plsr, X, Y[:, -1], cv=LeaveOneGroupOut().split(X, Y[:, -1], groups=groups))
8 print(f'Model Fiting (R-squared): {r2_score(Y[:, -1], predicted)}')
```

```
['NoCyt' 'IL-1' 'LPS' 'TNF' 'IL-6' 'Mix']
['DMSO_NoCyt' 'DMSO_IL-1' 'DMSO_LPS' 'DMSO_TNF' 'DMSO_IL-6' 'DMSO_Mix'
 'CIM_NoCyt' 'CIM_IL-1' 'CIM_LPS' 'CIM_TNF' 'CIM_IL-6' 'CIM_Mix'
 'RAN_NoCyt' 'RAN_IL-1' 'RAN_LPS' 'RAN_TNF' 'RAN_IL-6' 'RAN_Mix'
 'LEV_NoCyt' 'LEV_IL-1' 'LEV_LPS' 'LEV_TNF' 'LEV_IL-6' 'LEV_Mix'
 'TRO_NoCyt' 'TRO_IL-1' 'TRO_LPS' 'TRO_TNF' 'TRO_IL-6' 'TRO_Mix'
 'BUS_NoCyt' 'BUS_IL-1' 'BUS_LPS' 'BUS_TNF' 'BUS_IL-6' 'BUS_Mix'
 'NEF_NoCyt' 'NEF_IL-1' 'NEF_LPS' 'NEF_TNF' 'NEF_IL-6' 'NEF_Mix'
 'ASP_NoCyt' 'ASP_IL-1' 'ASP_LPS' 'ASP_TNF' 'ASP_IL-6' 'ASP_Mix'
 'NIM_NoCyt' 'NIM_IL-1' 'NIM_LPS' 'NIM_TNF' 'NIM_IL-6' 'NIM_Mix'
 'CLA_NoCyt' 'CLA_IL-1' 'CLA_LPS' 'CLA_TNF' 'CLA_IL-6' 'CLA_Mix'
 'TEL_NoCyt' 'TEL_IL-1' 'TEL_LPS' 'TEL_TNF' 'TEL_IL-6' 'TEL_Mix']
Model Fiting (R-squared): 0.8624270850931157
```

The results show that our R² has slightly improved with leave-one-cytokine-out cross validation compared to the previous iteration with LOOCV. This shows that the models predictive power had increased slightly. LOGO CV here has helped us determine whether coupling a cytokine with a drug really helps model learn unique features compared to there is no cytokine coupled or vice versa. LOOCV might be useful when testing a models overall ability to generalize.