

Singular Value Decomposition: many
other applications

Using SVD to solve square system of linear equations

If \mathbf{A} is a $n \times n$ square matrix and we want to solve $\mathbf{A} \mathbf{x} = \mathbf{b}$, we can use the SVD for \mathbf{A} such that

$$\begin{aligned} \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T \mathbf{x} &= \mathbf{b} \\ \mathbf{\Sigma} \mathbf{V}^T \mathbf{x} &= \mathbf{U}^T \mathbf{b} \end{aligned}$$

Solve: $\mathbf{\Sigma} \mathbf{y} = \mathbf{U}^T \mathbf{b}$ (diagonal matrix, easy to solve!)

Evaluate: $\mathbf{x} = \mathbf{V} \mathbf{y}$

Cost of solve: $O(n^2)$

Cost of decomposition $O(n^3)$ (recall that SVD and LU have the same cost asymptotic behavior, however the number of operations - constant factor before n^3 - for the SVD is larger than LU)

Matrix norms and condition number

- **The Euclidean norm of an orthogonal matrix is equal to 1**

$$\|U\|_2 = \max_{\|x\|_2=1} \|Ux\|_2 = \max_{\|x\|_2=1} \sqrt{(Ux)^T(Ux)} = \max_{\|x\|_2=1} \sqrt{x^T x} = \max_{\|x\|_2=1} \|x\|_2 = 1$$

- **The Euclidean norm of a matrix is given by the largest singular value**

$$\|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2 = \max_{\|x\|_2=1} \|U \Sigma V^T x\|_2$$

Since $\|U\|_2 = 1$, $\|V\|_2 = 1$ and Σ is diagonal then

$$\|A\|_2 = \max(\sigma_i) = \sigma_1$$

Matrix norms and condition number

- **The Euclidean norm of the inverse of a square-matrix is given by:**

$$\|A^{-1}\|_2 = \max_{\|x\|_2=1} \|(U \Sigma V^T)^{-1}x\|_2 = \max_{\|x\|_2=1} \|V \Sigma^{-1}U^T x\|_2$$

Since $\|U\|_2 = 1$, $\|V\|_2 = 1$ and Σ is diagonal then

$$\|A^{-1}\|_2 = \frac{1}{\sigma_{min}} \text{ (where } A \text{ is full rank, so that } A^{-1} \text{ exists)}$$

- **The norm of the pseudo-inverse of a $m \times n$ matrix is:**

$$\|A^+\|_2 = \frac{1}{\sigma_{min}} \text{ where } \sigma_{min} \text{ is the smallest non-zero singular value (for any matrix, regardless if it is full rank or rank deficient)}$$

- **The condition number of a matrix is given by**

$$cond_2(A) = \frac{\sigma_{max}}{\sigma_{min}} \text{ if } A \text{ is full rank or } cond_2(A) = \infty \text{ otherwise}$$

Matrix norms and condition number

- **The Euclidean norm of the inverse of a matrix is given by:**

$$\|A^+\|_2 = \max_{\|x\|_2=1} \|(U \Sigma V^T)^{-1}x\|_2 = \max_{\|x\|_2=1} \|V \Sigma^+ U^T x\|_2$$

Since $\|U\|_2 = 1$, $\|V\|_2 = 1$ and Σ is diagonal then

$$\text{rank}(A) = \min(m, n) \rightarrow \|A^+\|_2 = \frac{1}{\sigma_{\min}}$$

$$\text{rank}(A) < \min(m, n) \rightarrow \|A^+\|_2 = \infty$$

- **The condition number of a matrix is given by**

$$\text{rank}(A) = \min(m, n) \rightarrow \text{cond}_2(A) = \frac{\sigma_{\max}}{\sigma_{\min}}$$

$$\text{rank}(A) < \min(m, n) \rightarrow \text{cond}_2(A) = \infty$$

Matrix norms and SVD (more detailed derivation)

$$\|U\|_2 = \max_{\|x\|_2=1} \|Ux\|_2 = \max_{\|x\|_2=1} \sqrt{(Ux)^T(Ux)} \quad A = U \Sigma V^T$$

$$= \max_{\|x\|_2=1} \sqrt{x^T x} = \max_{\|x\|_2=1} \|x\|_2 = 1$$

$$\|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2 = \max_{\|x\|_2=1} \|U \Sigma V^T x\|_2 = \max_{\|x\|_2=1} \|\Sigma V^T x\|_2 =$$

$\|U\|_2 = 1$

$$= \max_{\|V^T x\|_2=1} \|\Sigma V^T x\|_2 = \max_{\|y\|_2=1} \|\Sigma y\|_2 = \sigma_1 = \max(\sigma_i)$$

$$\|V\|_2 = 1$$

Σ is diagonal

Condition number and SVD (more detailed derivation)

If matrix \mathbf{A} is square and non-singular:

$$\text{cond}_2(\mathbf{A}) = \|\mathbf{A}\|_2 \|\mathbf{A}^{-1}\|_2$$

$$\begin{aligned} \|\mathbf{A}^{-1}\|_2 &= \max_{\|\mathbf{x}\|_2=1} \|(\mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T)^{-1} \mathbf{x}\|_2 = \max_{\|\mathbf{x}\|_2=1} \|\mathbf{V}^{-T} \boldsymbol{\Sigma}^{-1} \mathbf{U}^{-1} \mathbf{x}\|_2 = \\ &= \max_{\|\mathbf{x}\|_2=1} \|\mathbf{V} \boldsymbol{\Sigma}^{-1} \mathbf{U}^T \mathbf{x}\|_2 = \max_{\|\mathbf{y}\|_2=1} \|\boldsymbol{\Sigma}^{-1} \mathbf{y}\|_2 = \frac{1}{\sigma_{min}} \end{aligned}$$

$$\text{cond}_2(\mathbf{A}) = \frac{\sigma_{max}}{\sigma_{min}}$$

What happens when you don't satisfy the conditions above?

Condition number and SVD (cont.)

$$\mathit{cond}_2(\mathbf{A}) = \|\mathbf{A}\|_2 \|\mathbf{A}^+\|_2$$

$$\|\mathbf{A}^+\|_2 = \max_{\|\mathbf{x}\|_2=1} \|\mathbf{V} \boldsymbol{\Sigma}^+ \mathbf{U}^T \mathbf{x}\|_2 = \max_{\|\mathbf{y}\|_2=1} \|\boldsymbol{\Sigma}^+ \mathbf{y}\|_2$$

If $\mathit{rank}(\mathbf{A}) = \min(m, n)$

$$\mathit{cond}_2(\mathbf{A}) = \frac{\sigma_{\max}}{\sigma_{\min}}$$

otherwise

$$\mathit{cond}_2(\mathbf{A}) = \infty$$

Low-Rank Approximation

Another way to write the SVD (assuming for now $m > n$ for simplicity)

$$\begin{aligned} A &= \begin{pmatrix} \vdots & \dots & \vdots \\ \mathbf{u}_1 & \dots & \mathbf{u}_m \\ \vdots & \dots & \vdots \end{pmatrix} \begin{pmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_n & \\ & & & 0 \\ & & & \vdots \\ & & & 0 \end{pmatrix} \begin{pmatrix} \dots & \mathbf{v}_1^T & \dots \\ \vdots & \vdots & \vdots \\ \dots & \mathbf{v}_n^T & \dots \end{pmatrix} \\ &= \begin{pmatrix} \vdots & \dots & \vdots \\ \mathbf{u}_1 & \dots & \mathbf{u}_n \\ \vdots & \dots & \vdots \end{pmatrix} \begin{pmatrix} \dots & \sigma_1 \mathbf{v}_1^T & \dots \\ \vdots & \vdots & \vdots \\ \dots & \sigma_n \mathbf{v}_n^T & \dots \end{pmatrix} \\ &= \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \dots + \sigma_n \mathbf{u}_n \mathbf{v}_n^T \end{aligned}$$

The SVD writes the matrix A as a sum of outer products (of left and right singular vectors).

Low-Rank Approximation

$$\mathbf{A} = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \cdots + \sigma_n \mathbf{u}_n \mathbf{v}_n^T \quad \sigma_1 \geq \sigma_2 \geq \sigma_3 \dots \geq 0$$

What is the rank of $\sigma_1 \mathbf{u}_1 \mathbf{v}_1^T$?

One (1 linearly independent column!).

What is the rank of $\sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T$?

Two (2 linearly independent columns!).

Low-Rank Approximation

The best **rank- k** approximation for a $m \times n$ matrix \mathbf{A} is given by:

$$\mathbf{A}_k = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \cdots + \sigma_k \mathbf{u}_k \mathbf{v}_k^T$$

$$\sigma_1 \geq \sigma_2 \geq \sigma_3 \dots \geq 0$$

Note that:

- \mathbf{A} is full rank if $\text{rank}(\mathbf{A}) = n$
- Best **rank- k** approximation \mathbf{A}_k has $\text{rank}(\mathbf{A}_k) = k$
- The error of this approximation can be expressed as:

$$\|\mathbf{A} - \mathbf{A}_k\|_2 = \|\sigma_{k+1} \mathbf{u}_{k+1} \mathbf{v}_{k+1}^T + \sigma_{k+2} \mathbf{u}_{k+2} \mathbf{v}_{k+2}^T + \cdots + \sigma_n \mathbf{u}_n \mathbf{v}_n^T\|_2 = \sigma_{k+1}$$

Low-Rank Approximation

The best **rank- k** approximation for a $m \times n$ matrix \mathbf{A} , (where $k \leq \min(m, n)$) is the one that minimizes the following problem:

$$\begin{aligned} \min_{A_k} \|\mathbf{A} - A_k\| \\ \text{such that } \text{rank}(A_k) \leq k. \end{aligned}$$

When using the induced 2-norm, the best **rank- k** approximation is given by:

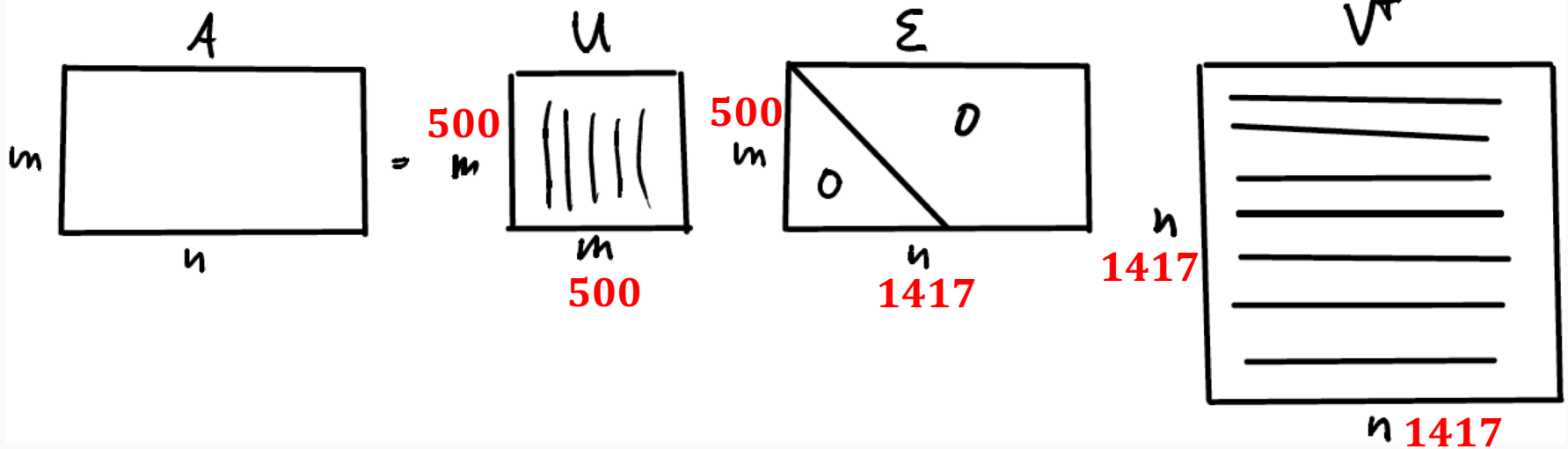
$$\begin{aligned} \mathbf{A}_k &= \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \cdots + \sigma_k \mathbf{u}_k \mathbf{v}_k^T \\ \sigma_1 &\geq \sigma_2 \geq \sigma_3 \dots \geq 0 \end{aligned}$$

Note that $\text{rank}(\mathbf{A}) = n$ and $\text{rank}(\mathbf{A}_k) = k$ and the norm of the difference between the matrix and its approximation is

$$\|\mathbf{A} - \mathbf{A}_k\|_2 = \left\| \sum_{i=k+1}^n \sigma_i \mathbf{u}_i \mathbf{v}_i^T \right\|_2 = \sigma_{k+1}$$

Image compression demo

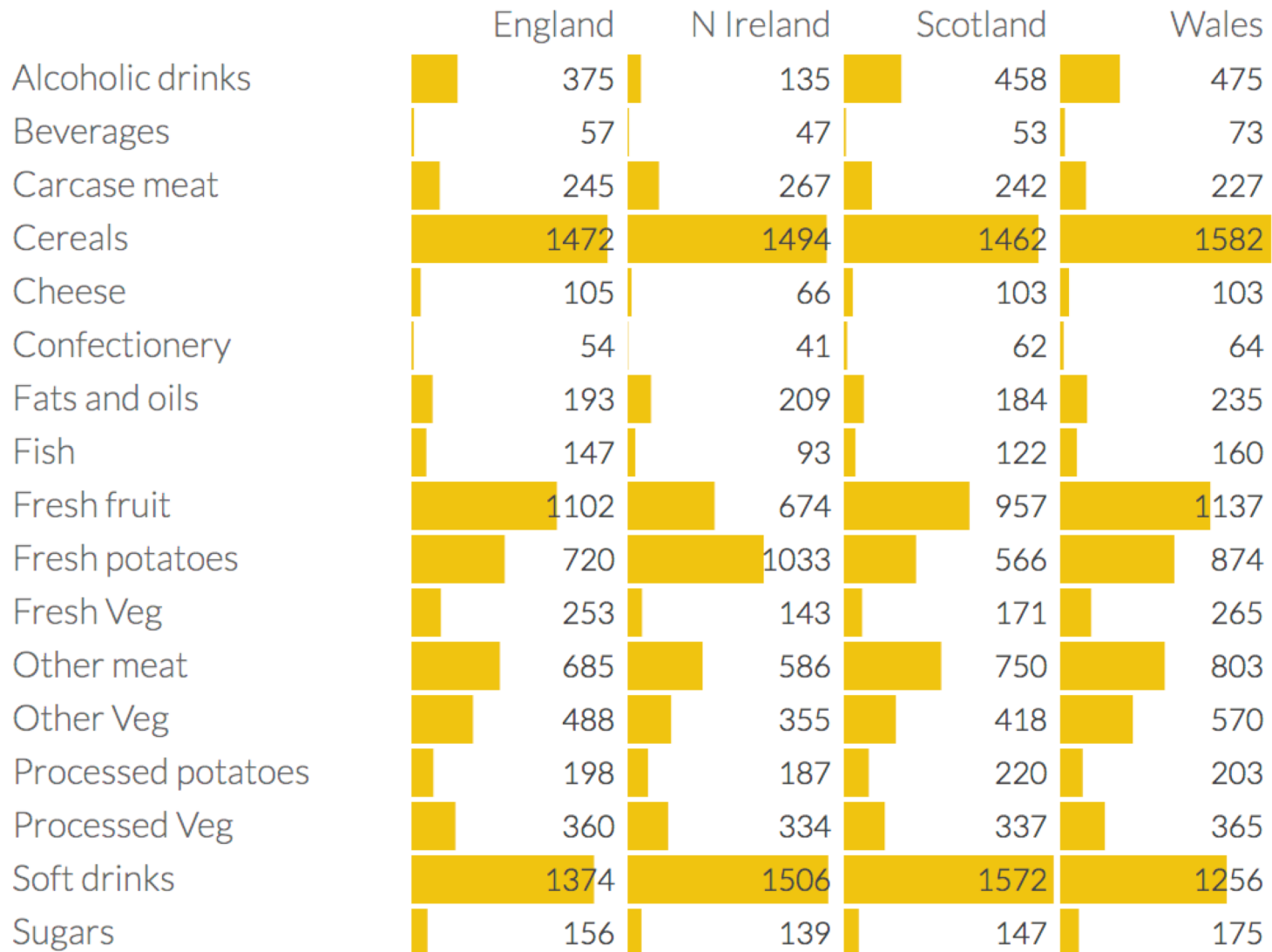
Demo "Image Compression"



Principal Component Analysis

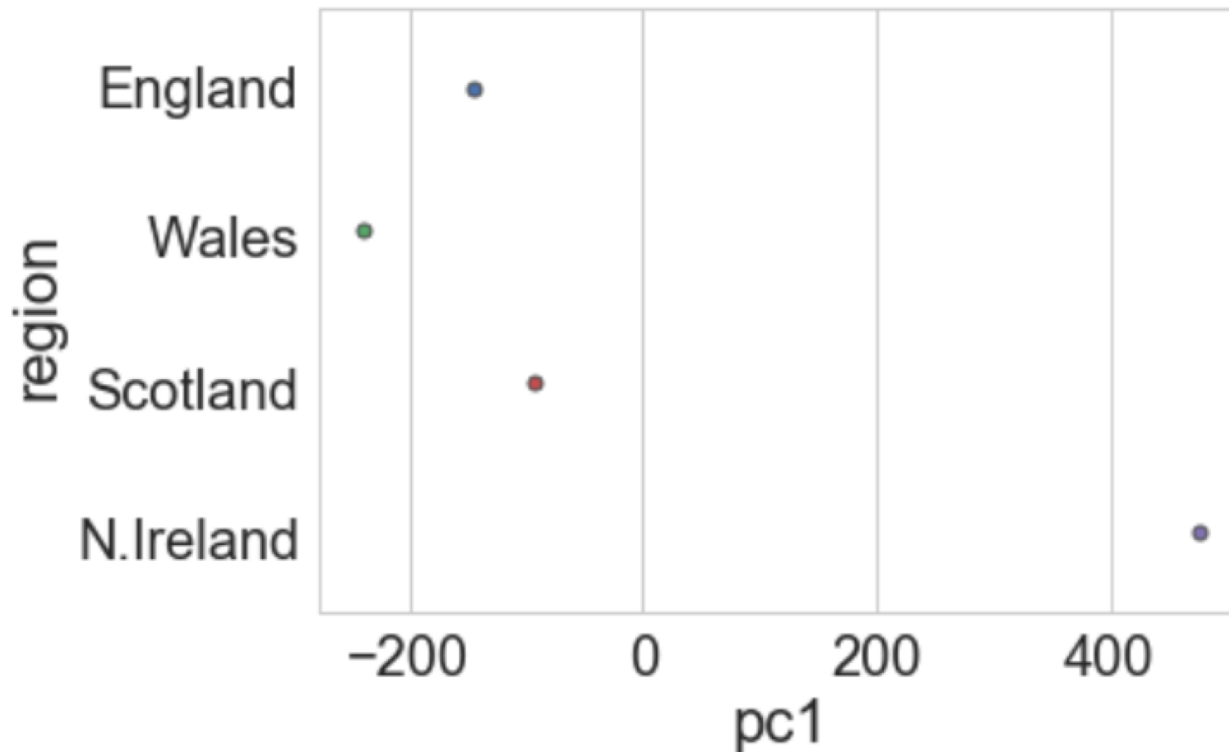
Food consumption in the UK

<http://setosa.io/ev/principal-component-analysis/>



How can we focus in just a few of the variables?

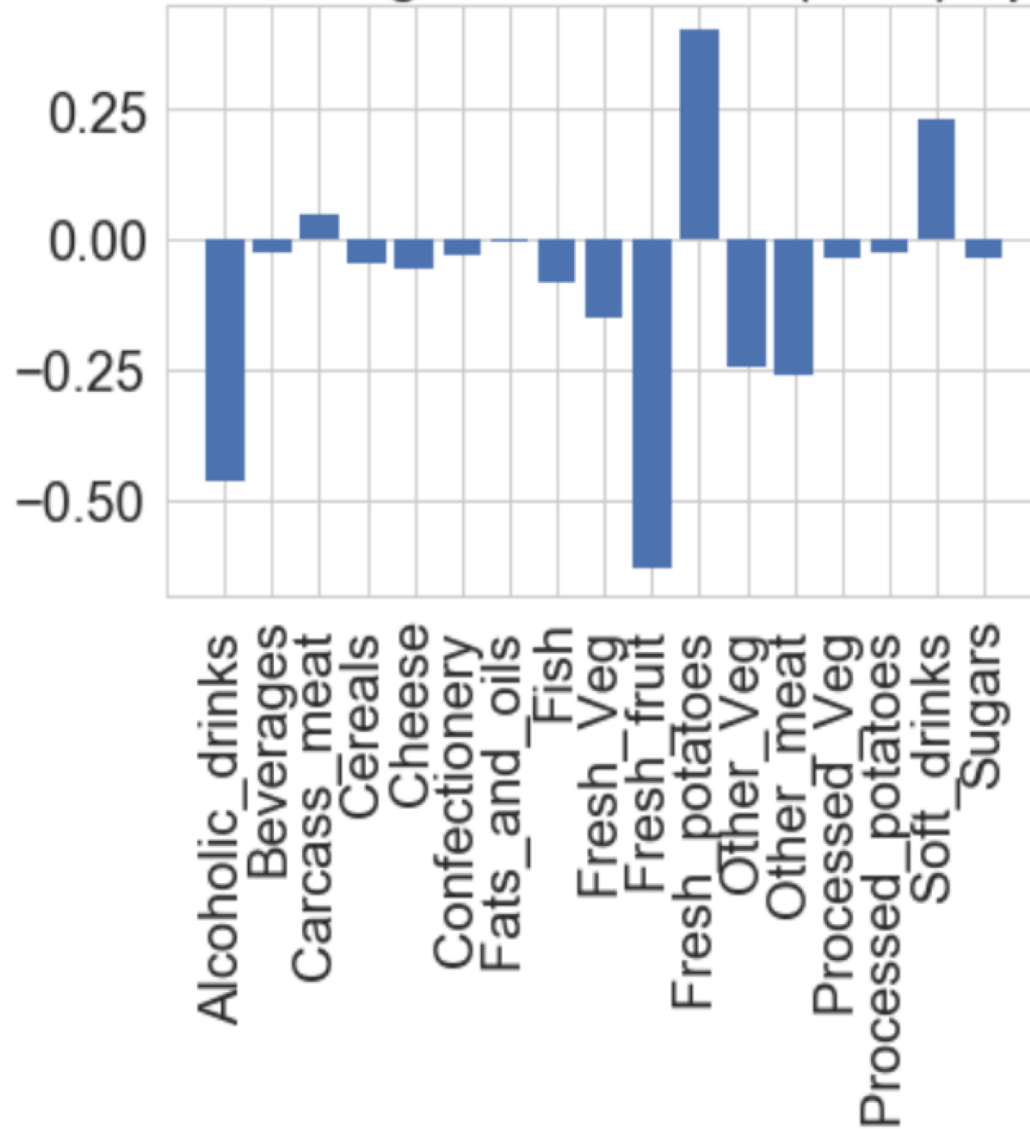
**We want to reduce the dimension of the feature space,
Let's try to reduce to one dimension:**



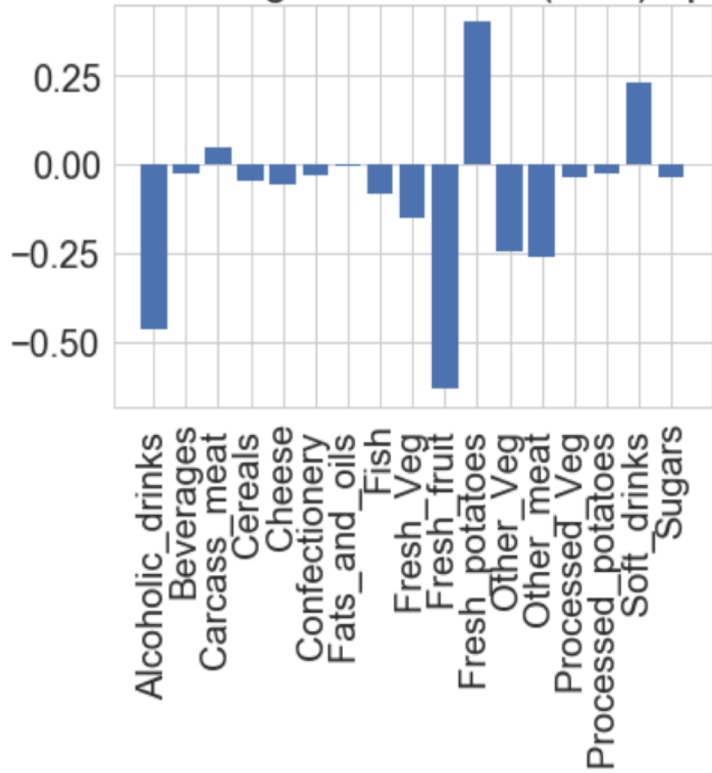
pc1: Principal component 1 - linear combination of the other 17 variables

$$pc1 = x1 \text{ Alcoholic Drinks} + x2 \text{ Beverages} + x3 \text{ Carcass meat} + \dots + x17 \text{ Sugars}$$

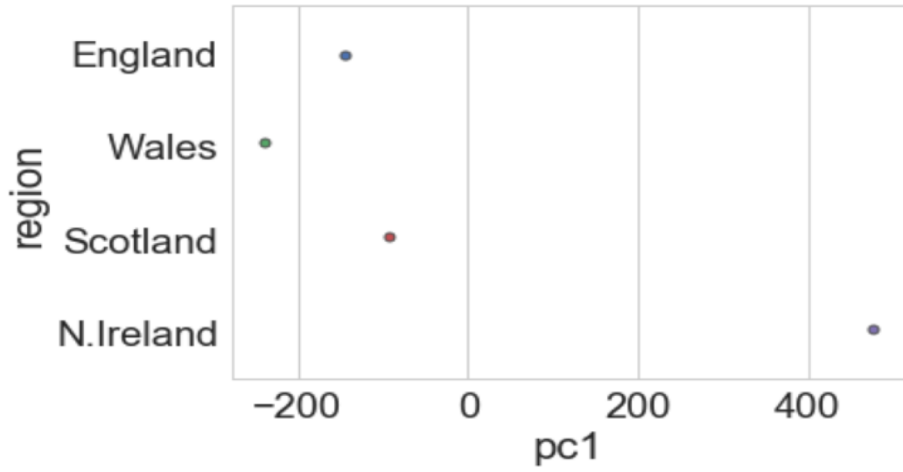
influence of original variables(food) upon pc1



influence of original variables(food) upon pc1

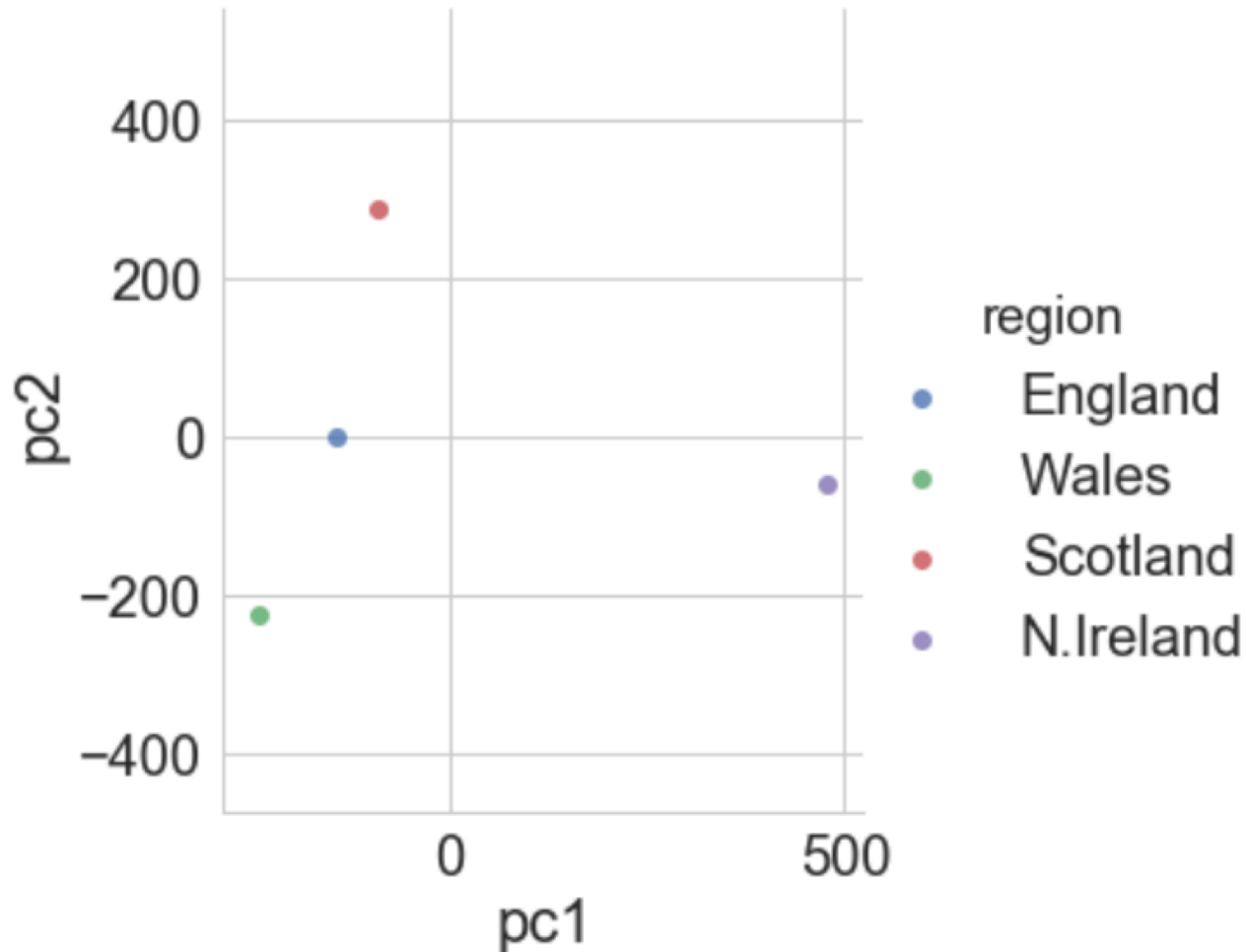


| | England | N Ireland | Scotland | Wales |
|--------------------|---------|-----------|----------|-------|
| Alcoholic drinks | 375 | 135 | 458 | 475 |
| Beverages | 57 | 47 | 53 | 73 |
| Carcass meat | 245 | 267 | 242 | 227 |
| Cereals | 1472 | 1494 | 1462 | 1582 |
| Cheese | 105 | 66 | 103 | 103 |
| Confectionery | 54 | 41 | 62 | 64 |
| Fats and oils | 193 | 209 | 184 | 235 |
| Fish | 147 | 93 | 122 | 160 |
| Fresh fruit | 1102 | 674 | 957 | 1137 |
| Fresh potatoes | 720 | 1033 | 566 | 874 |
| Fresh Veg | 253 | 143 | 171 | 265 |
| Other meat | 685 | 586 | 750 | 803 |
| Other Veg | 488 | 355 | 418 | 570 |
| Processed potatoes | 198 | 187 | 220 | 203 |
| Processed Veg | 360 | 334 | 337 | 365 |
| Soft drinks | 1374 | 1506 | 1572 | 1256 |
| Sugars | 156 | 139 | 147 | 175 |

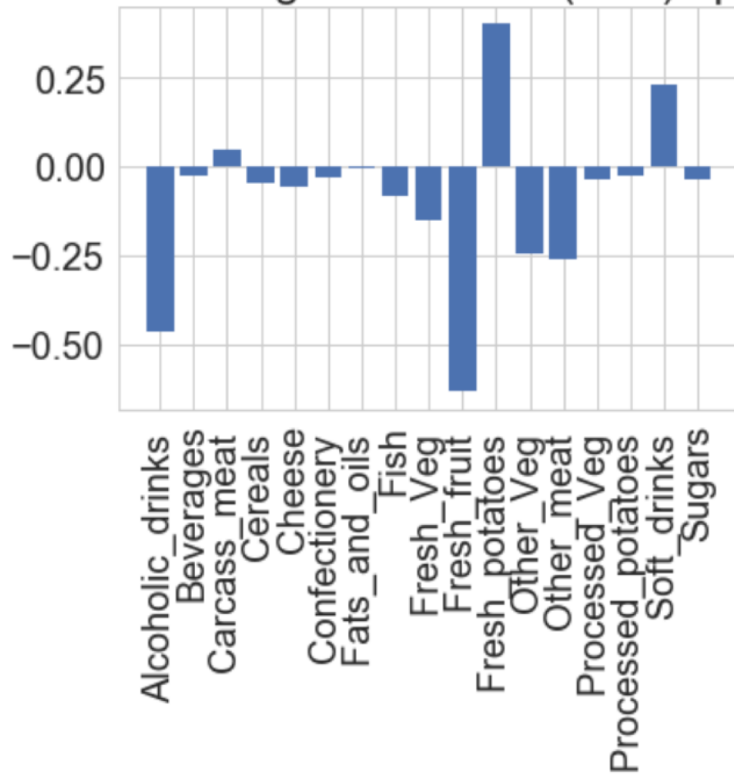


How can we focus in just a few of the variables?

What about reducing to two dimensions?

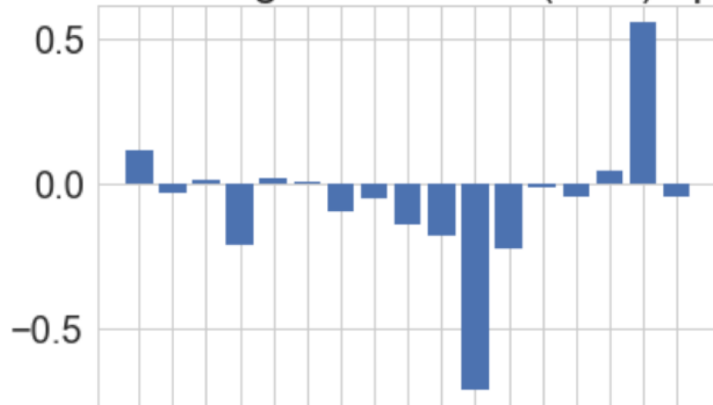


influence of original variables(food) upon pc1



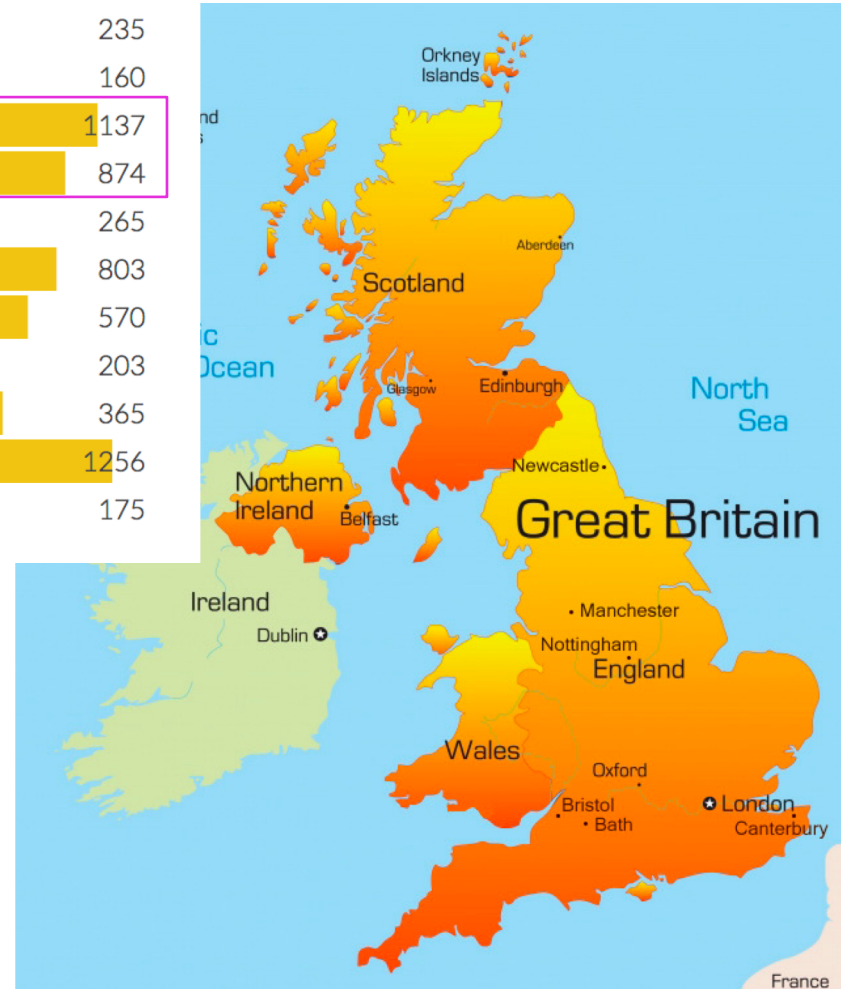
| | England | N Ireland | Scotland | Wales |
|--------------------|---------|-----------|----------|-------|
| Alcoholic drinks | 375 | 135 | 458 | 475 |
| Beverages | 57 | 47 | 53 | 73 |
| Carcass meat | 245 | 267 | 242 | 227 |
| Cereals | 1472 | 1494 | 1462 | 1582 |
| Cheese | 105 | 66 | 103 | 103 |
| Confectionery | 54 | 41 | 62 | 64 |
| Fats and oils | 193 | 209 | 184 | 235 |
| Fish | 147 | 93 | 122 | 160 |
| Fresh fruit | 1102 | 674 | 957 | 1137 |
| Fresh potatoes | 720 | 1033 | 566 | 874 |
| Fresh Veg | 253 | 143 | 171 | 265 |
| Other meat | 685 | 586 | 750 | 803 |
| Other Veg | 488 | 355 | 418 | 570 |
| Processed potatoes | 198 | 187 | 220 | 203 |
| Processed Veg | 360 | 334 | 337 | 365 |
| Soft drinks | 1374 | 1506 | 1572 | 1256 |
| Sugars | 156 | 139 | 147 | 175 |

influence of original variables(food) upon pc2



The three variables, Fresh potatoes, Alcoholic drinks and Fresh fruit, there is a noticeable difference between the values for England, Wales and Scotland, which are roughly similar, and Northern Ireland, which is usually significantly higher or lower.

| | England | N Ireland | Scotland | Wales |
|--------------------|---------|-----------|----------|-------|
| Alcoholic drinks | 375 | 135 | 458 | 475 |
| Beverages | 57 | 47 | 53 | 73 |
| Carcase meat | 245 | 267 | 242 | 227 |
| Cereals | 1472 | 1494 | 1462 | 1582 |
| Cheese | 105 | 66 | 103 | 103 |
| Confectionery | 54 | 41 | 62 | 64 |
| Fats and oils | 193 | 209 | 184 | 235 |
| Fish | 147 | 93 | 122 | 160 |
| Fresh fruit | 1102 | 674 | 957 | 1137 |
| Fresh potatoes | 720 | 1033 | 566 | 874 |
| Fresh Veg | 253 | 143 | 171 | 265 |
| Other meat | 685 | 586 | 750 | 803 |
| Other Veg | 488 | 355 | 418 | 570 |
| Processed potatoes | 198 | 187 | 220 | 203 |
| Processed Veg | 360 | 334 | 337 | 365 |
| Soft drinks | 1374 | 1506 | 1572 | 1256 |
| Sugars | 156 | 139 | 147 | 175 |



Predicting breast cancer

Goal (MP): Use data about tumor cell features to create a model to predict if a breast tumor is malign or benign.

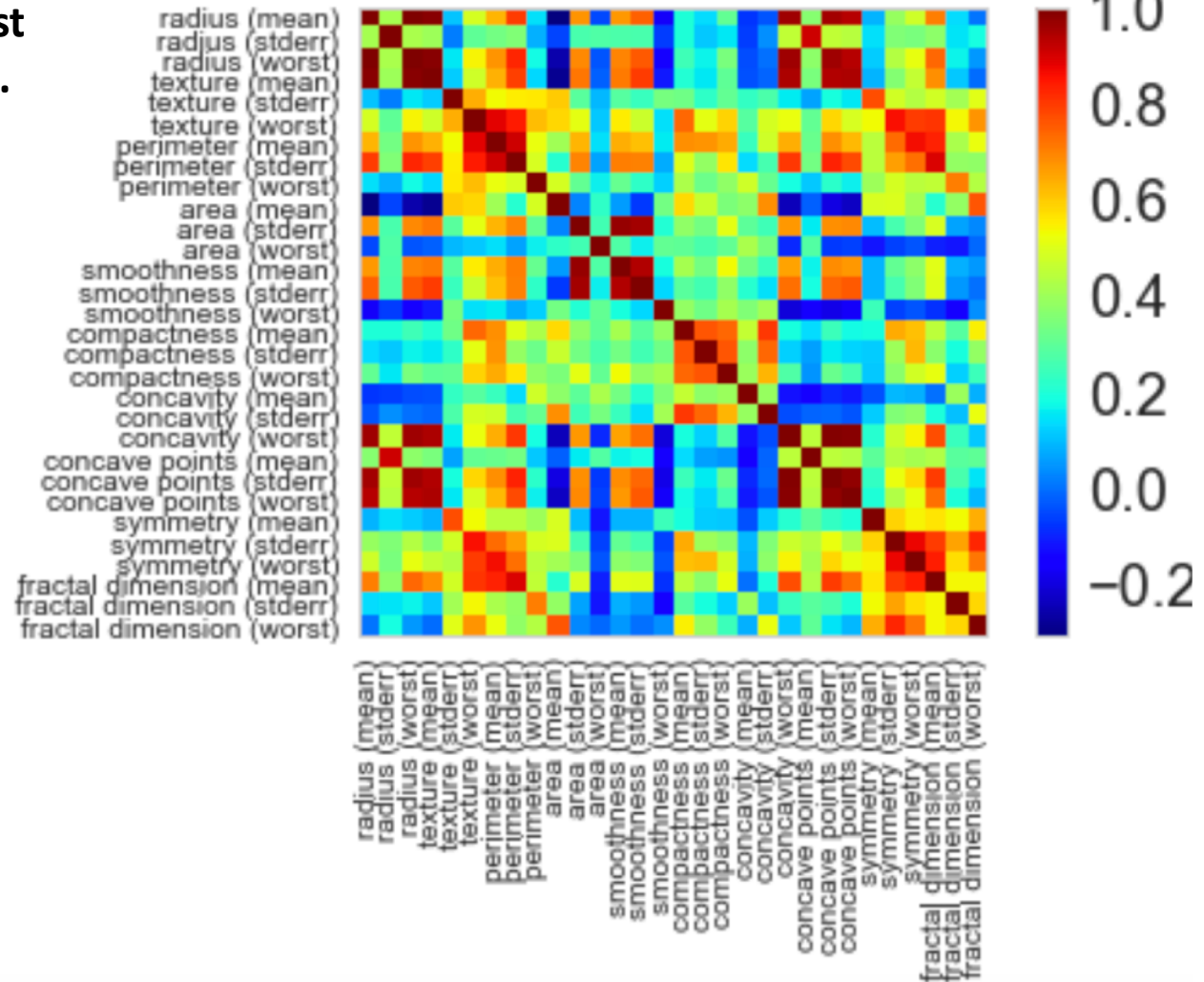
The data includes 30 different cell features.

There are many variables that are highly correlated with each other.

Reduce the feature space:

Approach 1: remove some of the feature variables.

Tumor features correlation matrix



Example: Reduce the feature space by including only the features regarding the mean

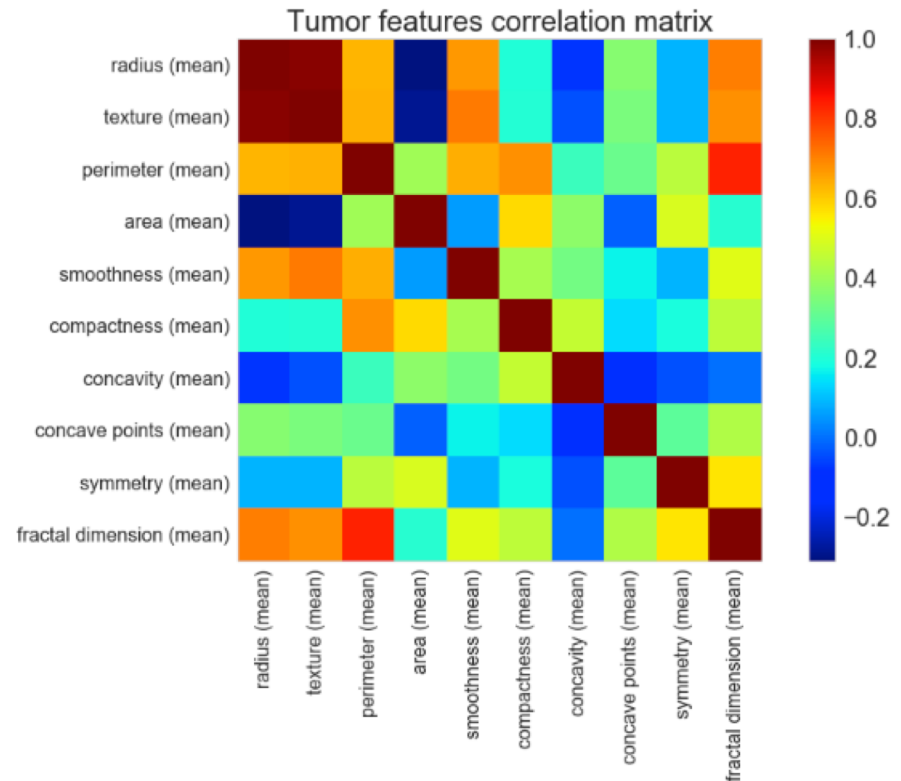
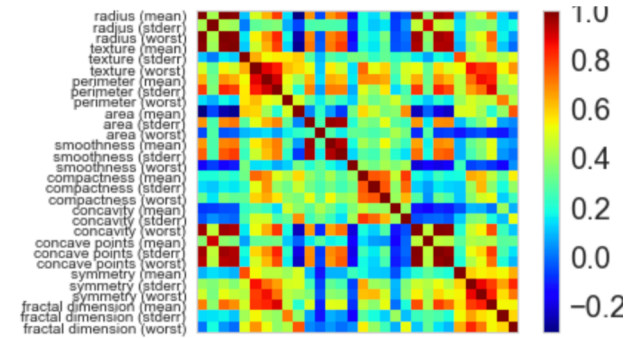
$$A = \begin{bmatrix} \vdots & \vdots & \vdots \\ F_1 & \dots & F_{30} \\ \vdots & \vdots & \vdots \end{bmatrix}$$



$$A^* = \begin{bmatrix} \vdots & \vdots & \vdots \\ F_1 & \dots & F_{10} \\ \vdots & \vdots & \vdots \end{bmatrix}$$

PROS: simple and maintain interpretation of the feature variables

CONS: lose information from the variables that were dropped



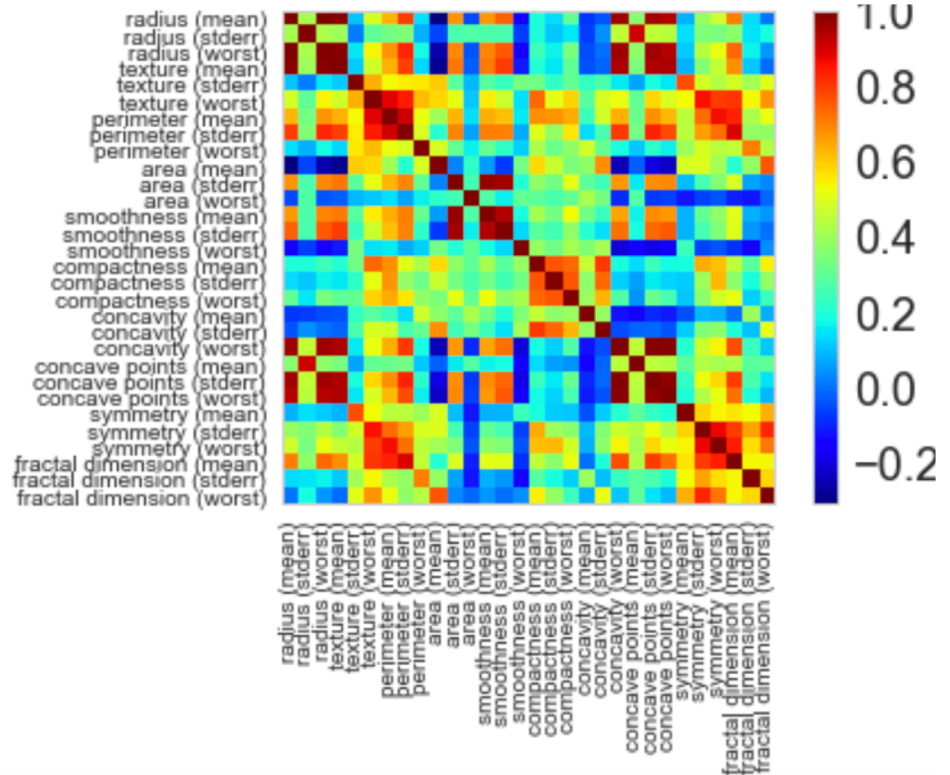
Get a new data set, resulting from a linear combination of the original dataset

$$A = \begin{bmatrix} \vdots & \vdots & \vdots \\ F_1 & \dots & F_{30} \\ \vdots & \vdots & \vdots \end{bmatrix}$$



$$A^* = \begin{bmatrix} \vdots & \vdots & \vdots \\ F_1^* & F_2^* & F_3^* \\ \vdots & \vdots & \vdots \end{bmatrix}$$

$$F_1^* = \sum_{i=1}^n a_i F_i$$



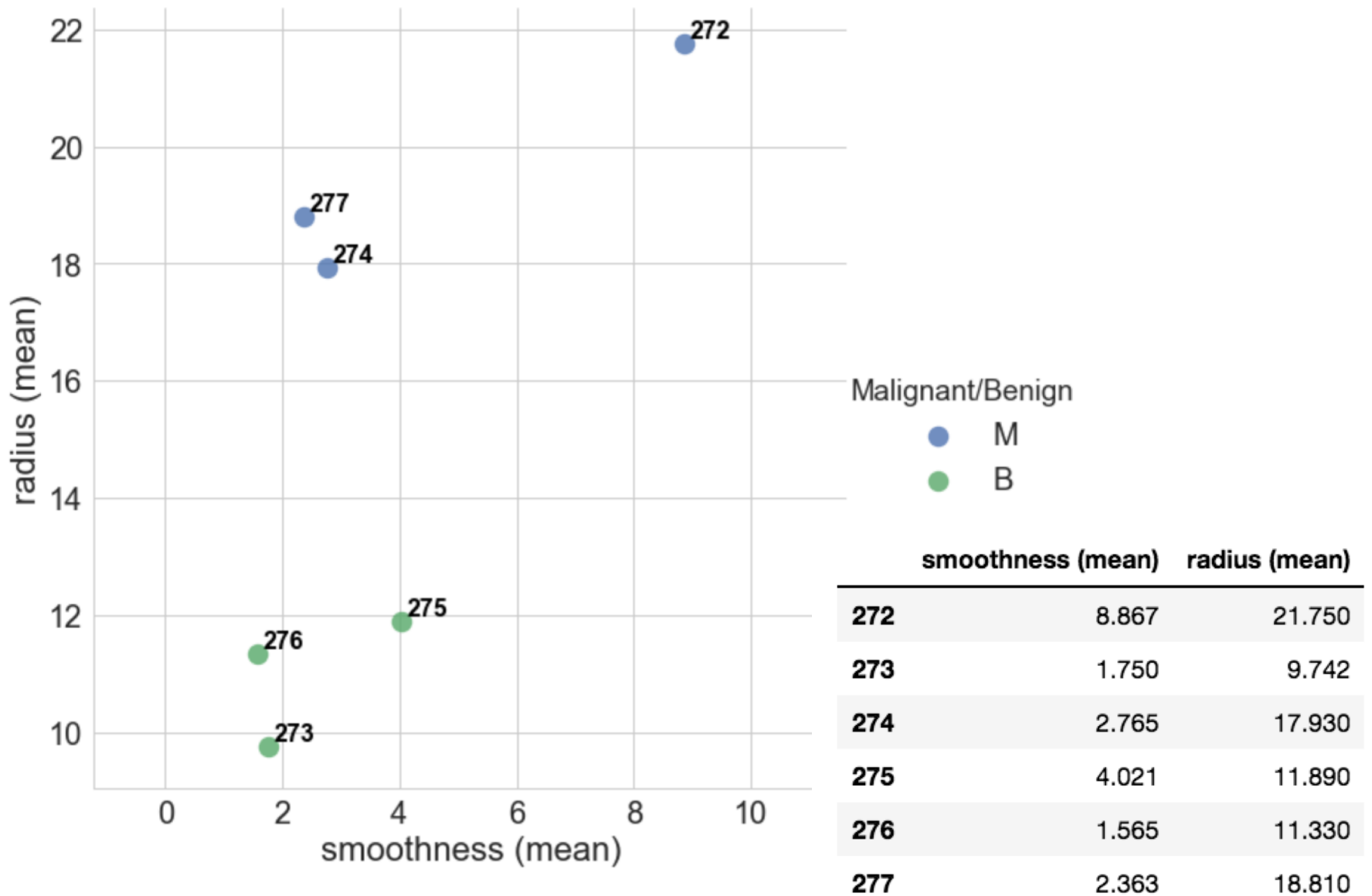
PROS: less variables containing information of all features

CONS: the new features no longer have a “meaningful” interpretation (here a characteristic of a tumor cell)

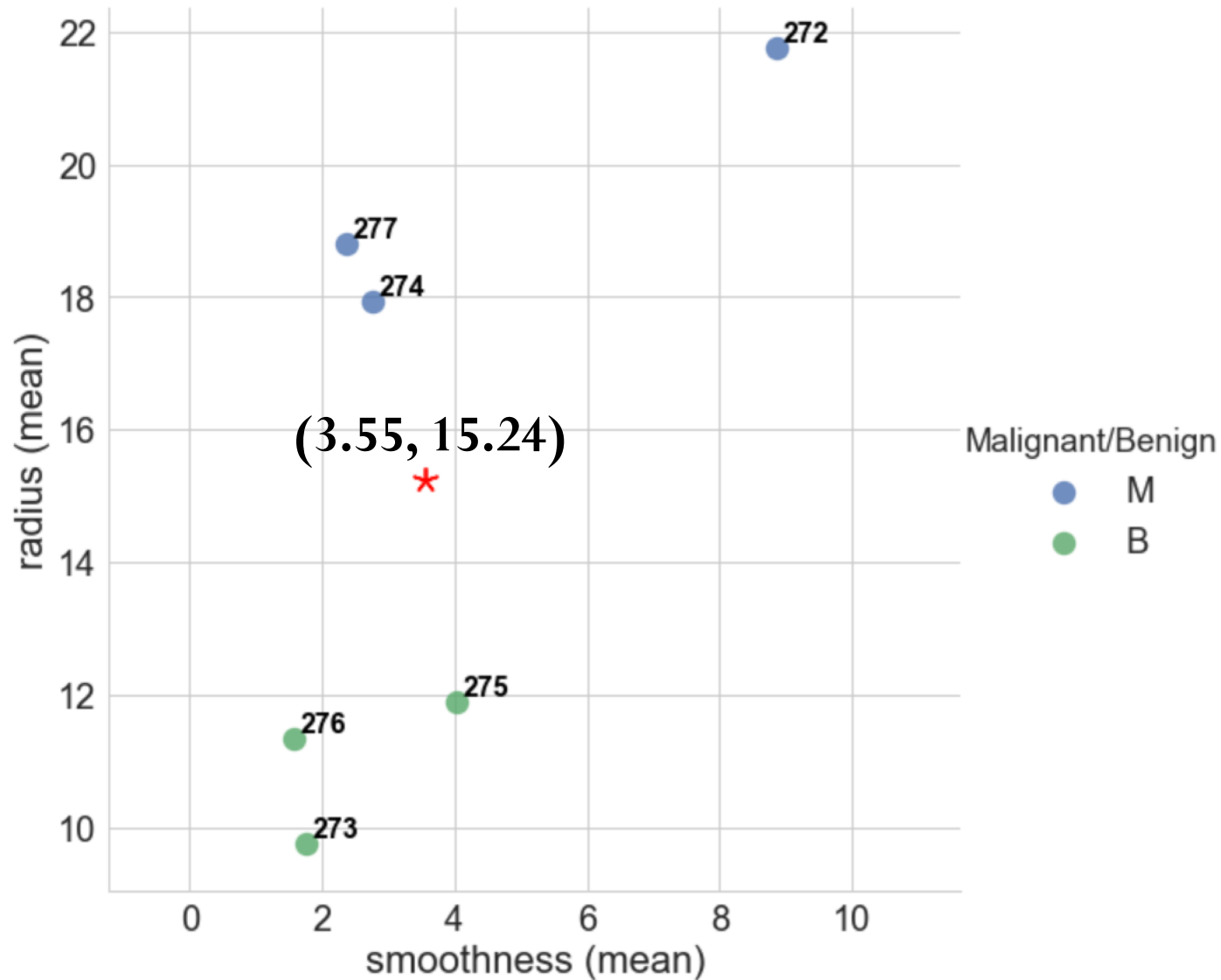
Principal component analysis

- PCA will combine the feature variables in a specific way, creating “new variables”.
- We can now drop the “least important” new variables while still retaining the most valuable parts of all of the feature variables!
- As an added benefit, each of the “new variables” after PCA are all independent of one another (important requirement for linear models).
- Cons: the new variables don't have the same meaning as the feature variables (loss of interpretability)

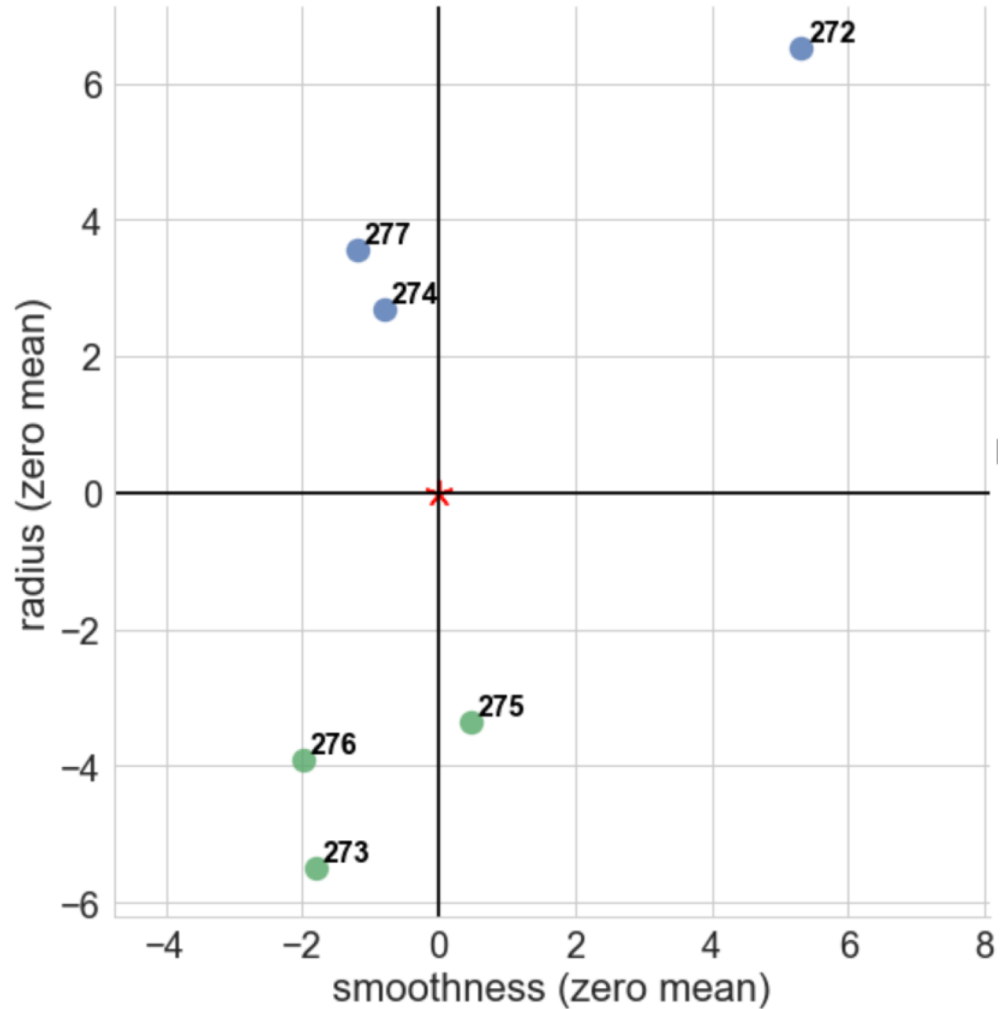
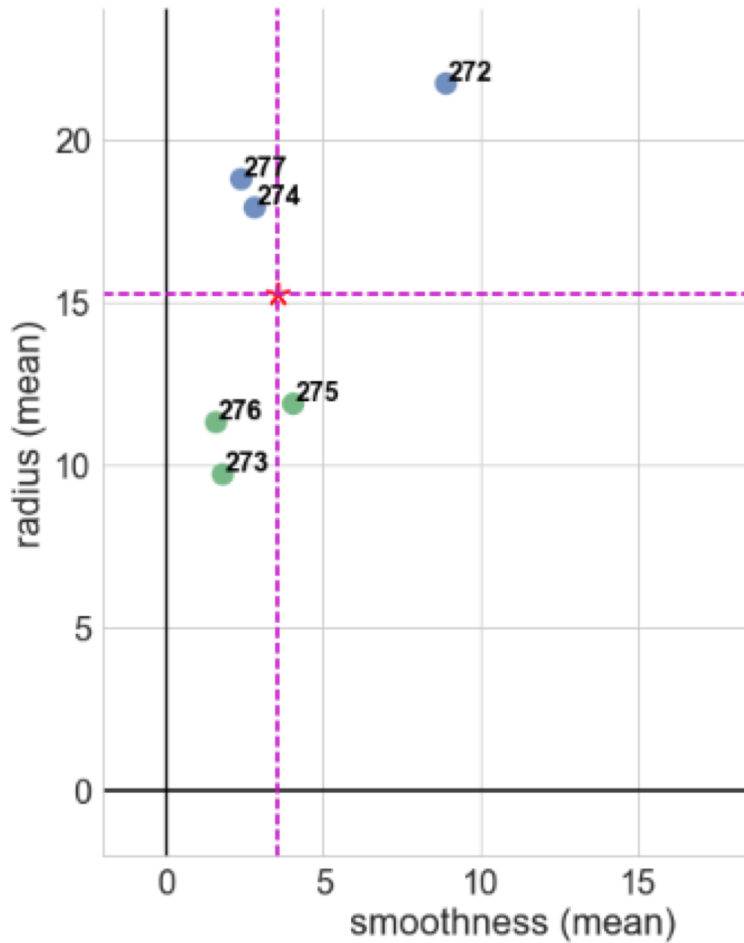
Let's start with a subset of 6 patients, and take a look at only two of the features: smoothness and radius



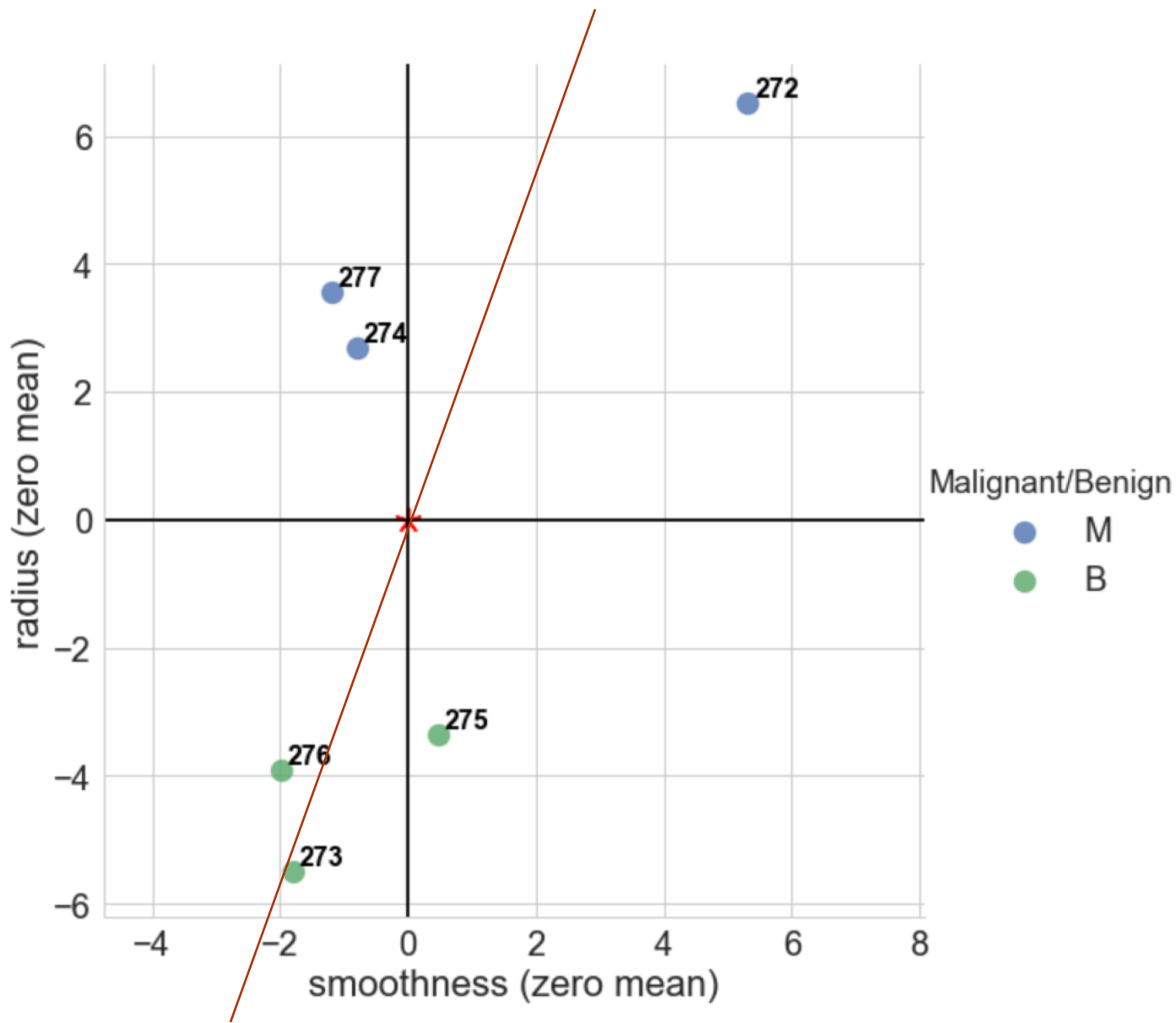
Determine the “center” of the dataset – the mean value of each feature



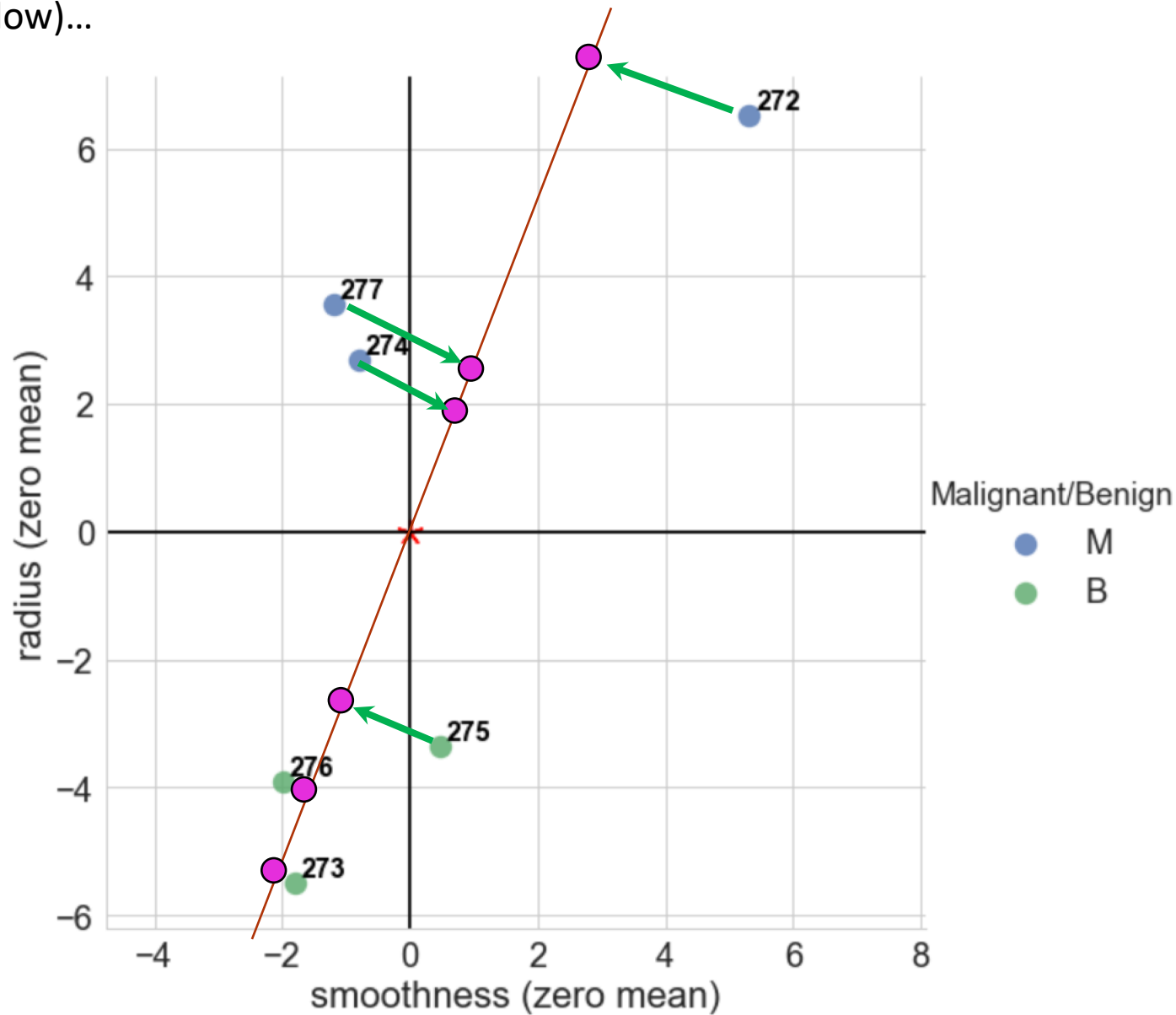
We will shift the dataset such that the “center” of the dataset (mean value) is at the origin (0,0) – the new dataset has zero mean value.



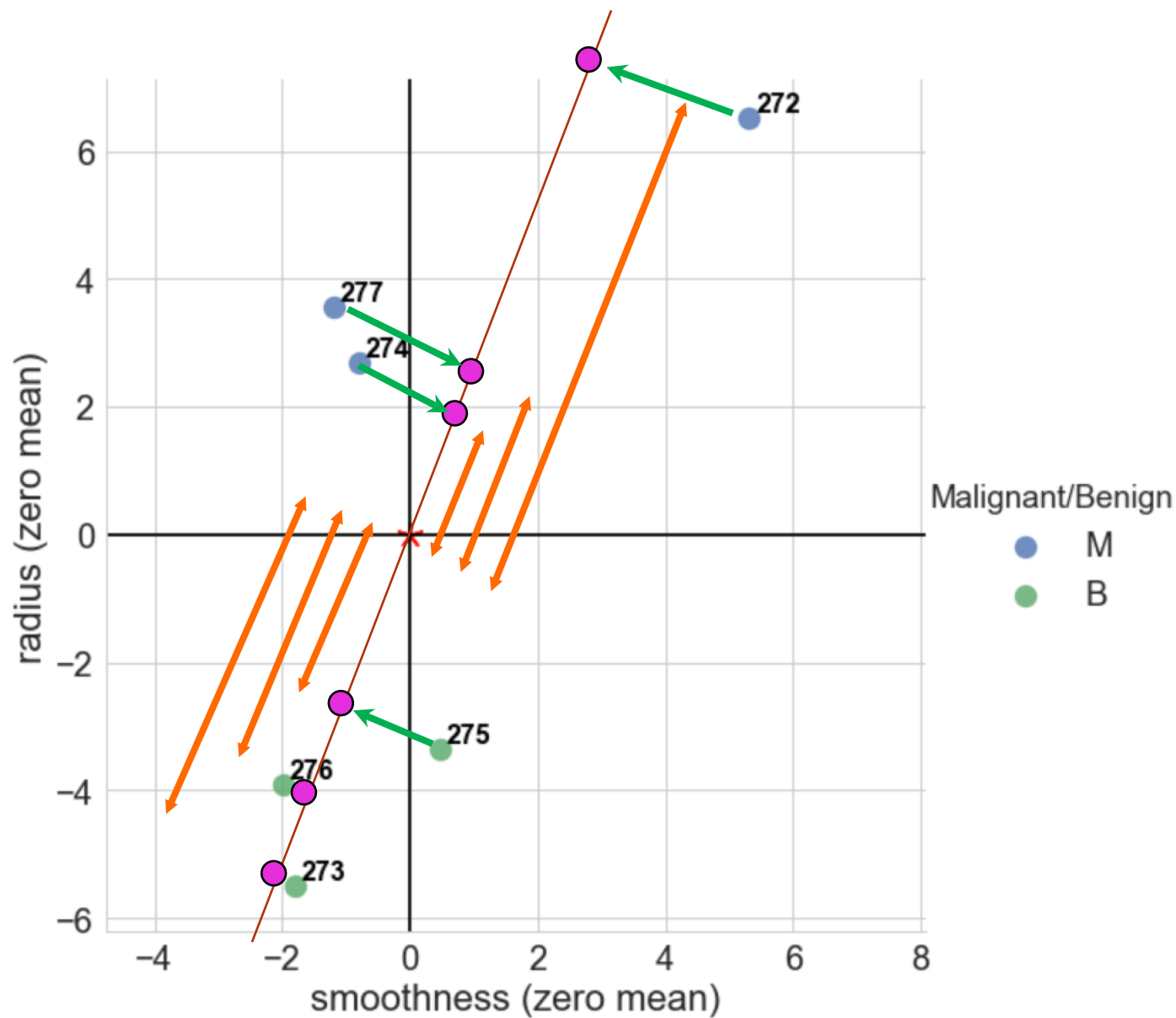
We want to find a straight line that fits the dataset.



Let's propose the red line below. To quantify how good the fit is, PCA projects the data onto the line. The best fit minimizes the distances from the points to the line (indicated in green below)...

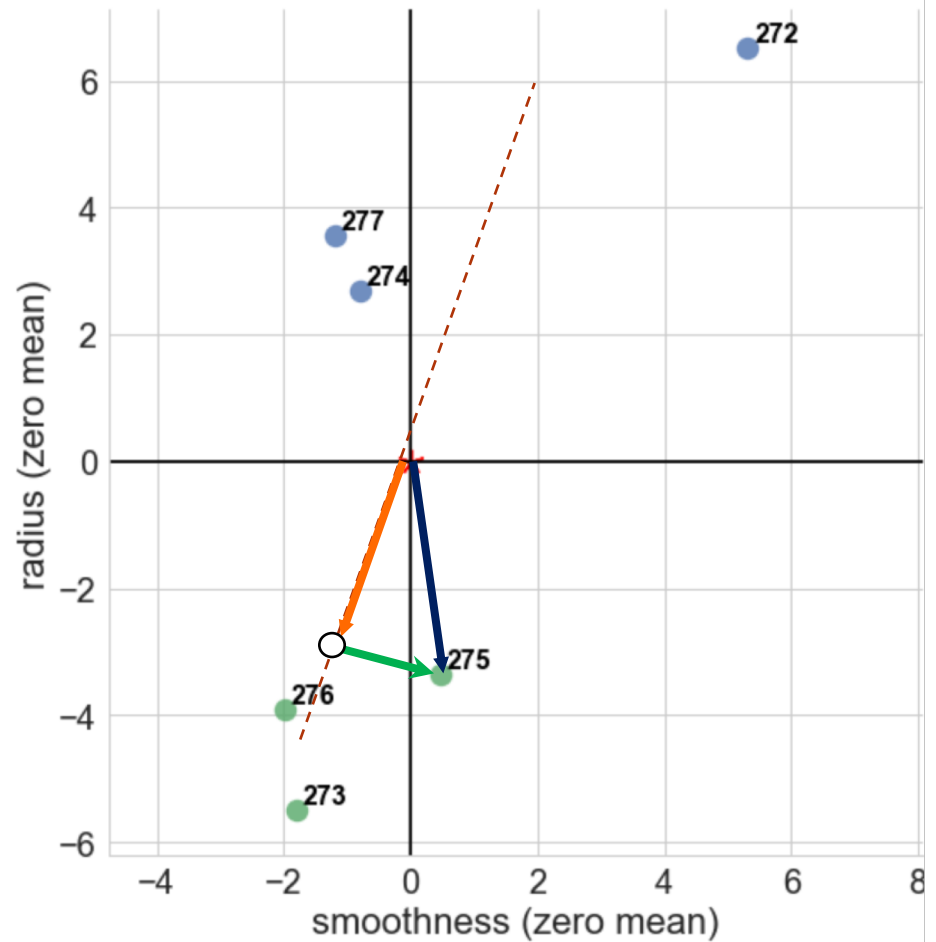
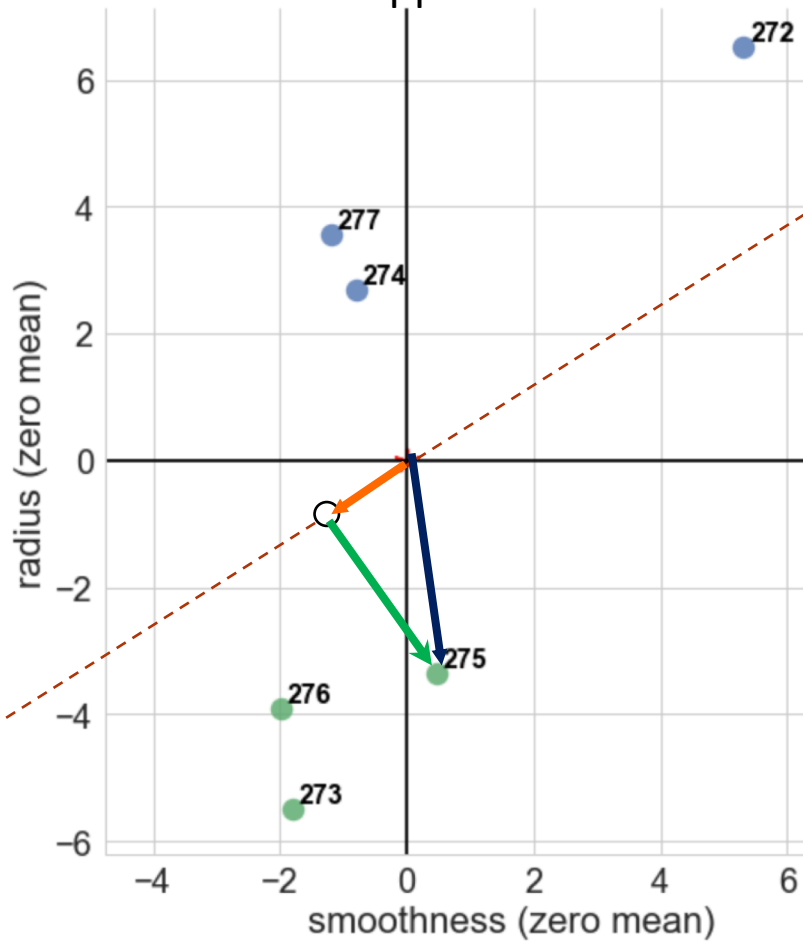


Or maximizes the distances from the projected points to the origin (indicated in orange)

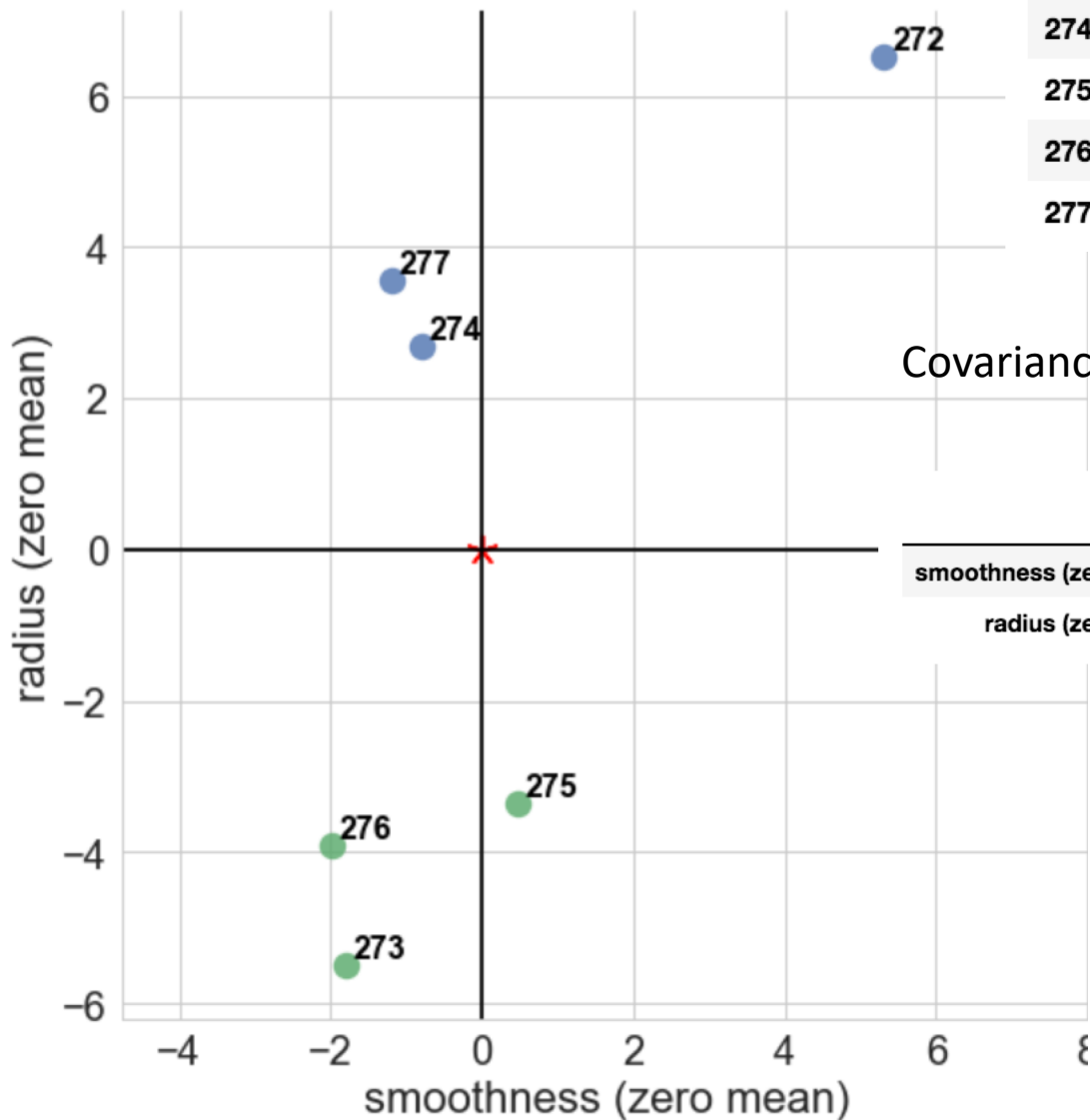


Why are they the same?

Take a look at what happens to the vectors below when we change the fit curve.



Let's talk about the variance of the dataset



| | smoothness (zero mean) | radius (zero mean) |
|-----|------------------------|--------------------|
| 272 | 5.311833 | 6.508 |
| 273 | -1.805167 | -5.500 |
| 274 | -0.790167 | 2.688 |
| 275 | 0.465833 | -3.352 |
| 276 | -1.990167 | -3.912 |
| 277 | -1.192167 | 3.568 |

$A =$

Covariance matrix: $\frac{1}{(n-1)} A^T A$

| | smoothness (zero mean) | radius (zero mean) |
|------------------------|------------------------|--------------------|
| smoothness (zero mean) | 7.539518 | 8.868854 |
| radius (zero mean) | 8.868854 | 23.819936 |

Covariance matrix: $\frac{1}{(n-1)} A^T A$

| | smoothness (zero mean) | radius (zero mean) |
|-----|------------------------|--------------------|
| 272 | 5.311833 | 6.508 |
| 273 | -1.805167 | -5.500 |
| 274 | -0.790167 | 2.688 |
| 275 | 0.465833 | -3.352 |
| 276 | -1.990167 | -3.912 |
| 277 | -1.192167 | 3.568 |

$A =$

| | smoothness (zero mean) | radius (zero mean) |
|------------------------|------------------------|--------------------|
| smoothness (zero mean) | 7.539518 | 8.868854 |
| radius (zero mean) | 8.868854 | 23.819936 |

Diagonalization of covariance matrix:

$$A^T A = XDX^T$$

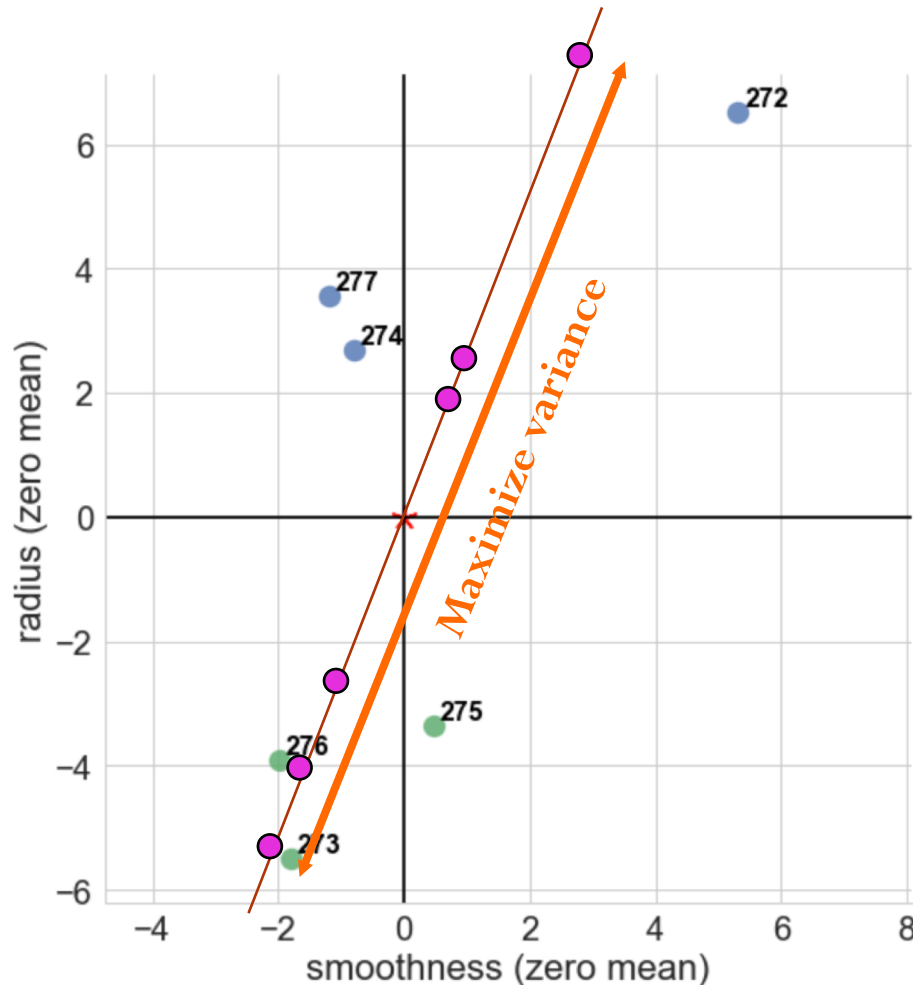
X : eigenvectors of $A^T A$

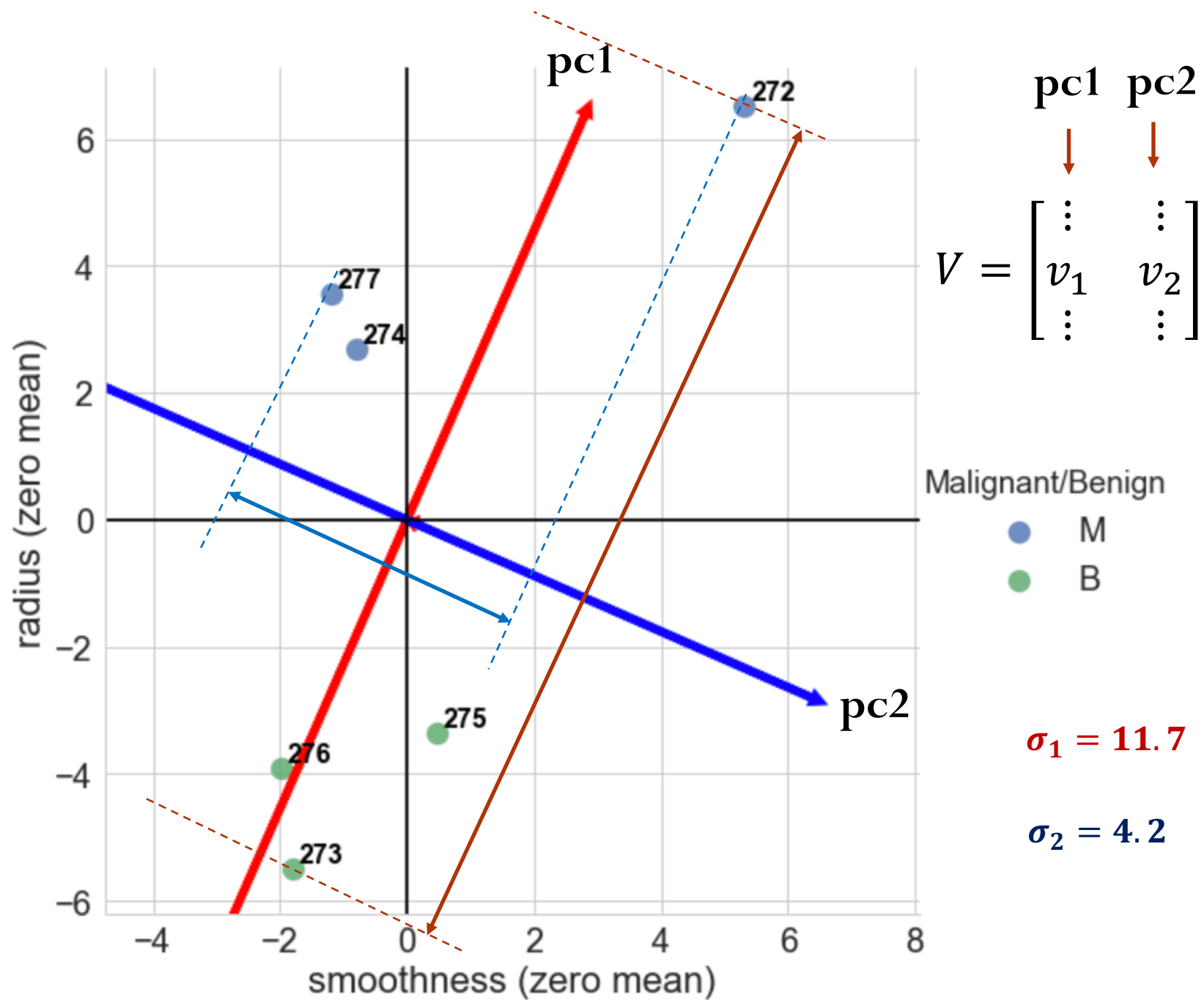
D : eigenvalues of $A^T A$

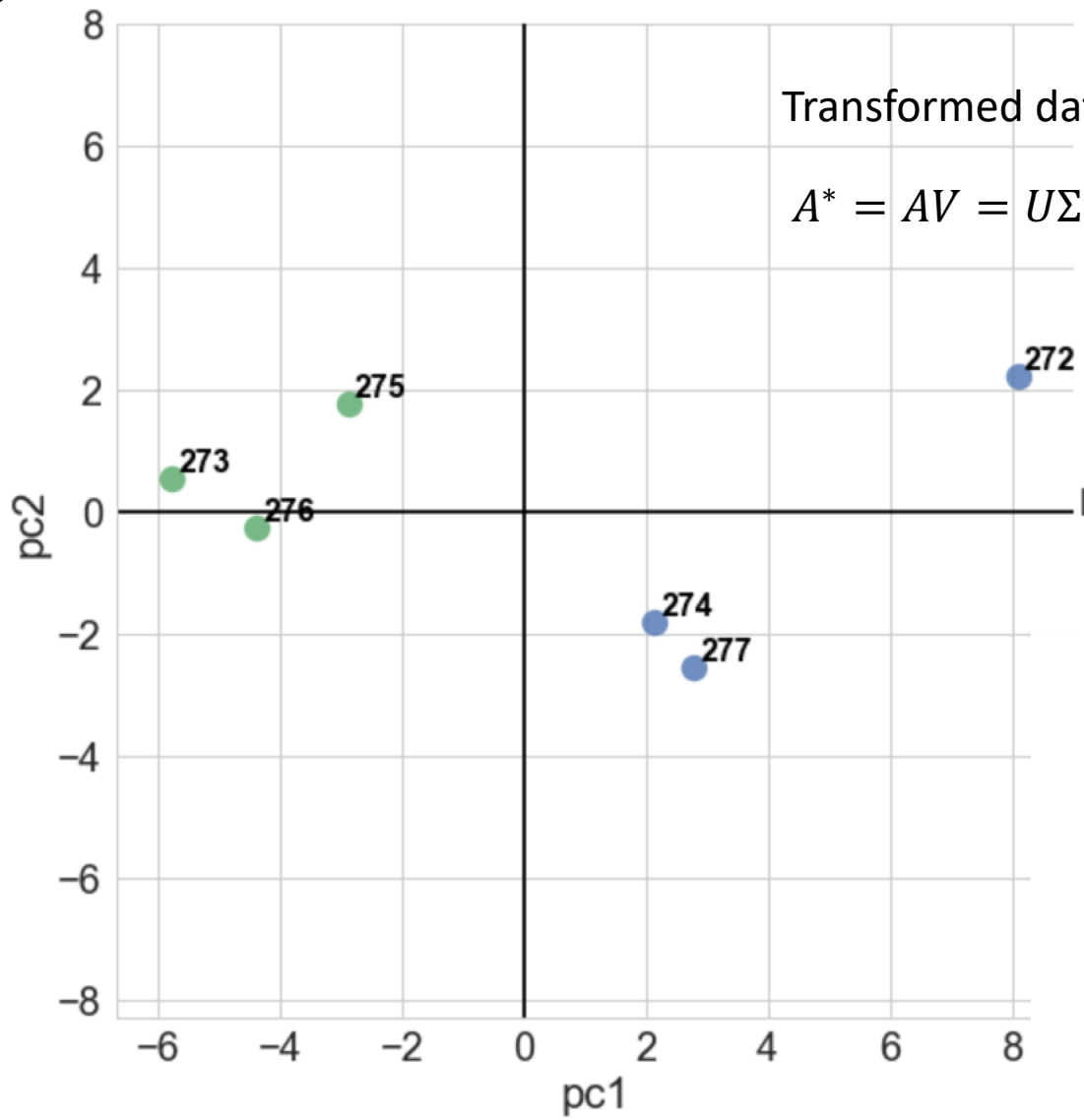
From SVD: $A = U\Sigma V^T$

Maximum variance:
largest singular value of Σ

Direction of maximum variance:
Corresponding column of V

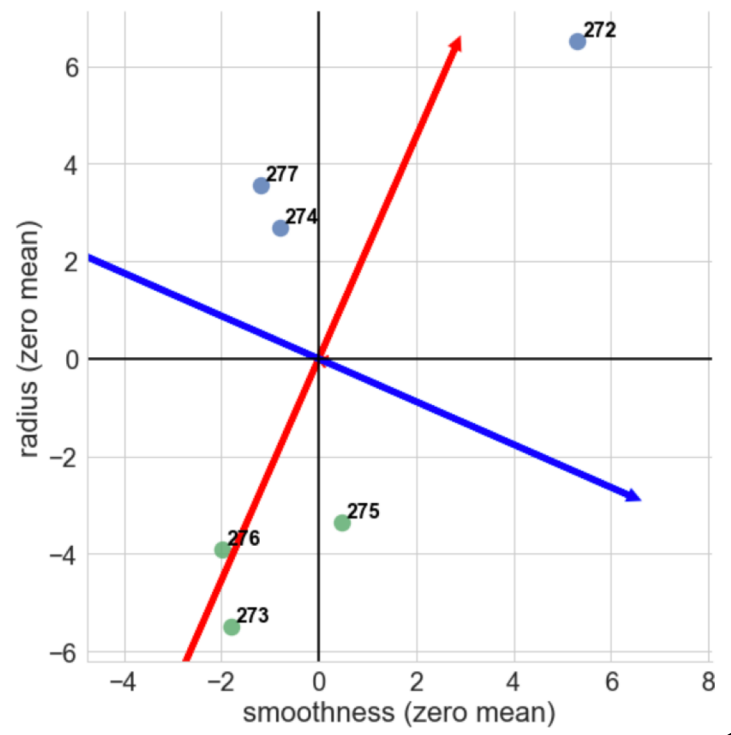


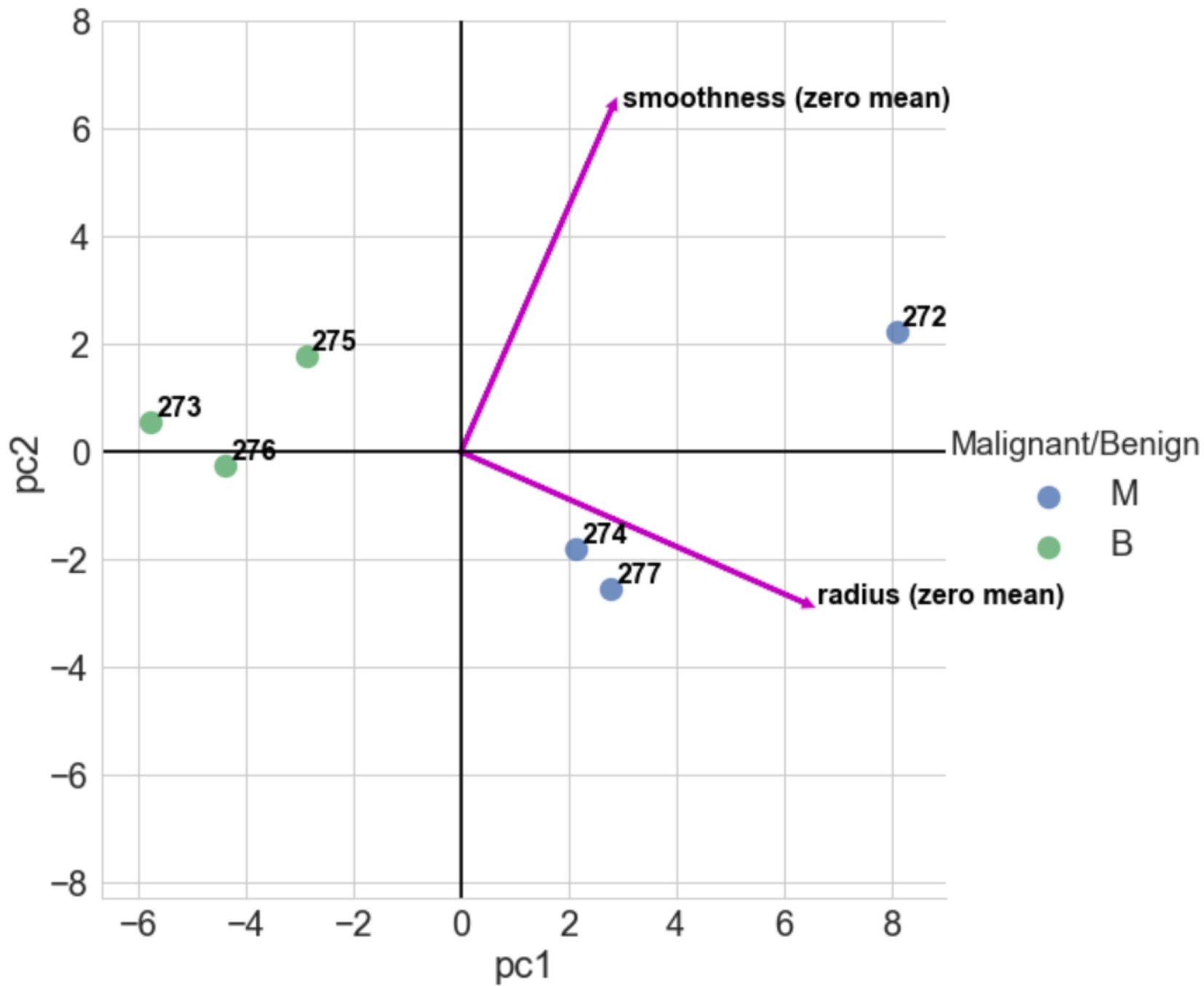




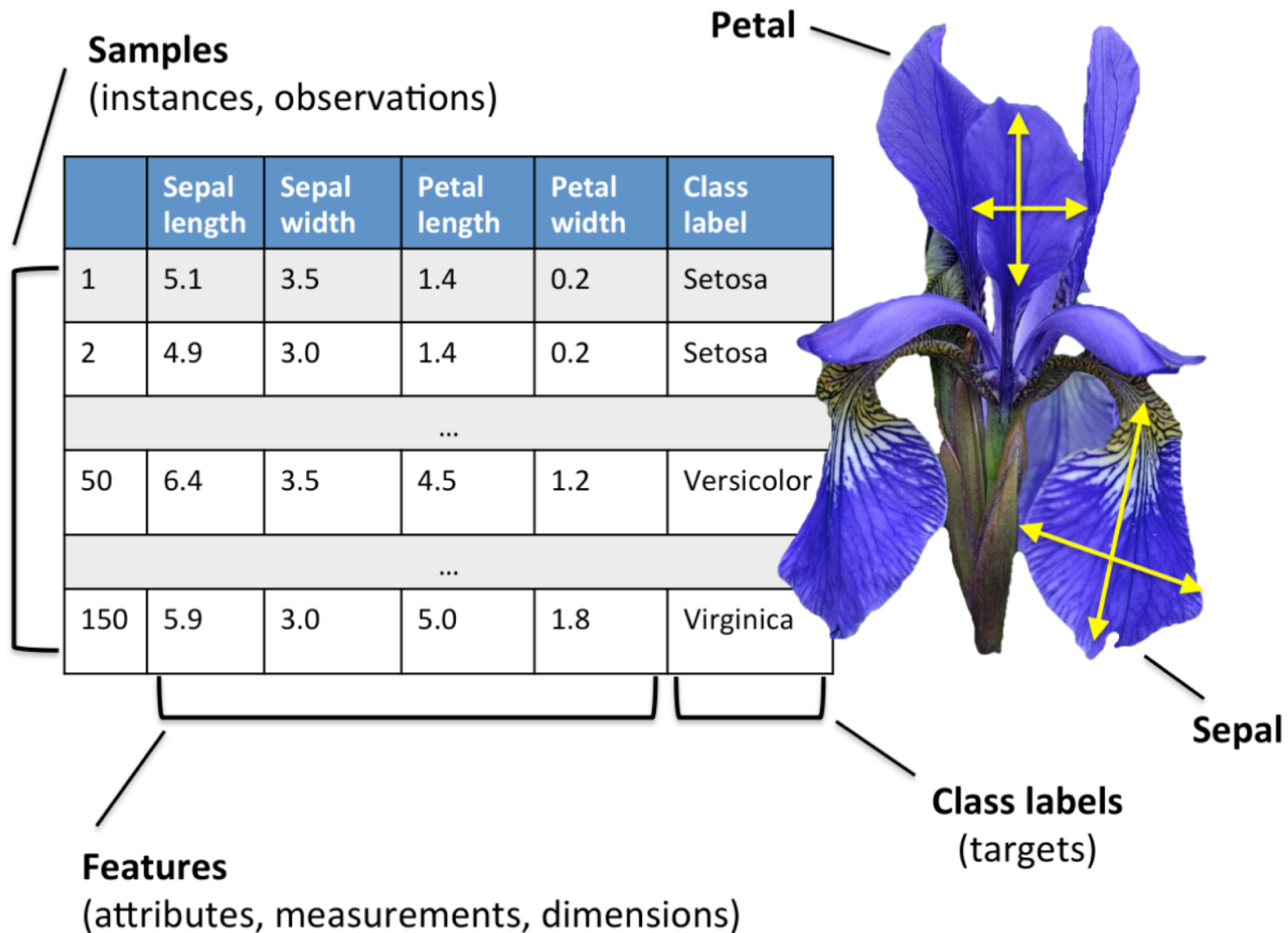
Malignant/Benign

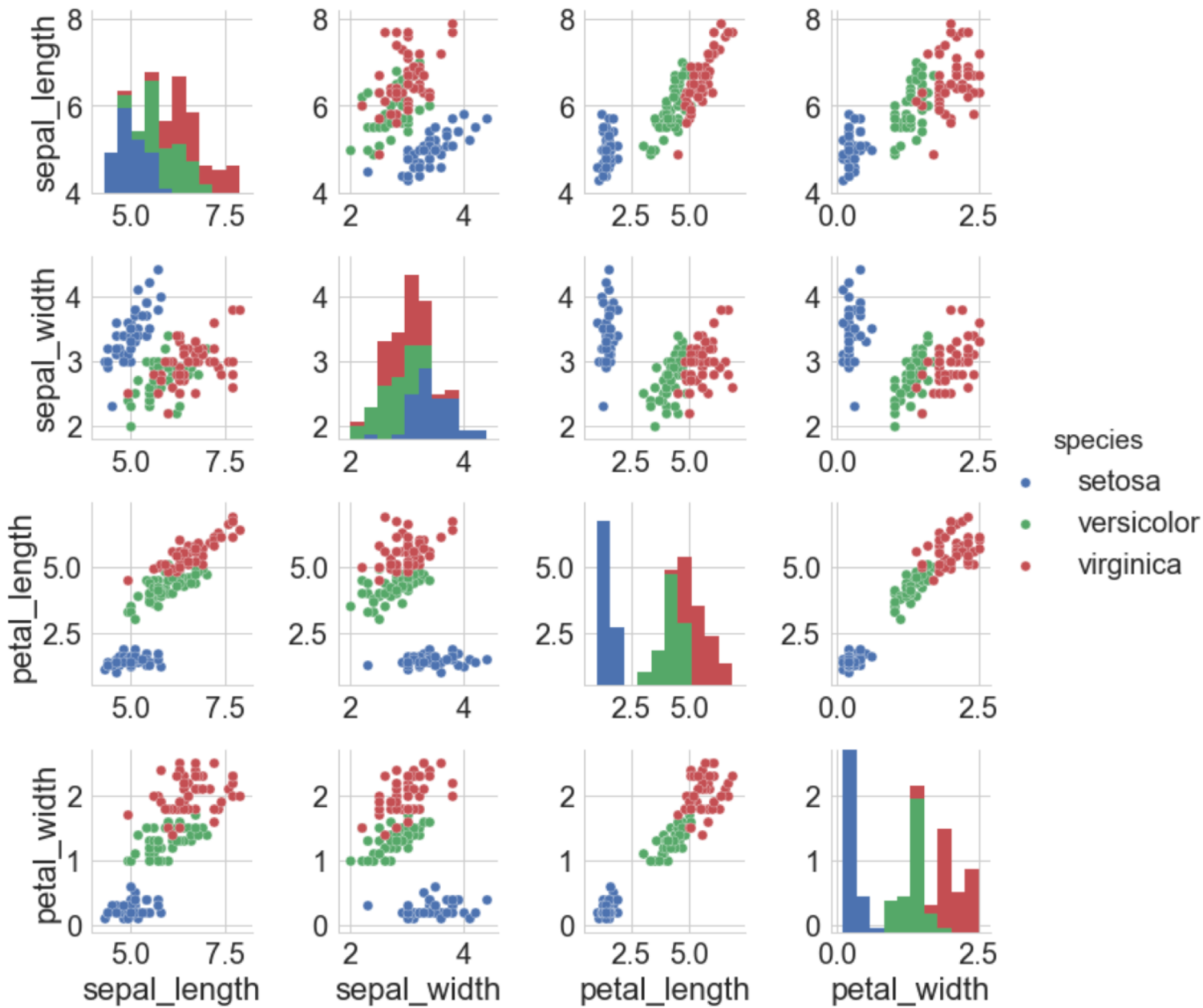
- M
- B





Let's add more features! Flower classification





Principal component analysis

How can we reduce the dimension of a dataset without missing important information?

Detect correlation between variables, if a strong correlation exists, then reducing the dimension of the dataset makes sense.

Overall idea: Find the directions of maximum variance in high-dimensional dataset (n dimension) and project it onto a subspace with smaller dimension (k dimension, with $k < n$), while retaining most of the information.

What is the adequate value for k ?

- 1) Shift the dataset to zero mean: $A = A - A.mean()$
- 2) Compute SVD: $A = U\Sigma V^T$
- 3) Principal components: variances = singular values squared
- 4) Principal directions: columns of V
- 5) New dataset: $A^* = A V$

Note how the variances of the new dataset correspond to the singular values squared of the original dataset:

$$(A^*)^T A = V^T A^T A V = V^T (U\Sigma V^T)^T U\Sigma V^T V = \Sigma^T \Sigma$$

6) In general: $A^* = A V$

$m \times n$ $m \times n$ $n \times n$

7) But since we want to reduce the dimension of the dataset, we only use the first k columns of V

$A^* = A V$

$m \times k$ $m \times n$ $n \times k$

Iris dataset

1) Shift the dataset to zero mean:

```
1 X = iris.iloc[:, :4] - iris.iloc[:, :4].mean()
```

Optional (modeling choice!): decide whether or not to standardize. If you want to standardize, divide each observation in a column by that column's standard deviation.

```
1 Z = X / iris.iloc[:, :4].std()
```

In this new dataset Z each feature has mean zero and standard deviation 1.

This decision depends on the problem you are solving. If some variables have a large variance and some small, since PCA maximizes the variance, it will weight more the features with large variance. If you want your PCA to be independent of the variance, standardizing the features will do that.

```
1 iris.head()
```

```
:
```

| | sepal_length | sepal_width | petal_length | petal_width | species |
|---|--------------|-------------|--------------|-------------|---------|
| 0 | 5.1 | 3.5 | 1.4 | 0.2 | setosa |
| 1 | 4.9 | 3.0 | 1.4 | 0.2 | setosa |
| 2 | 4.7 | 3.2 | 1.3 | 0.2 | setosa |
| 3 | 4.6 | 3.1 | 1.5 | 0.2 | setosa |
| 4 | 5.0 | 3.6 | 1.4 | 0.2 | setosa |

Explained variance

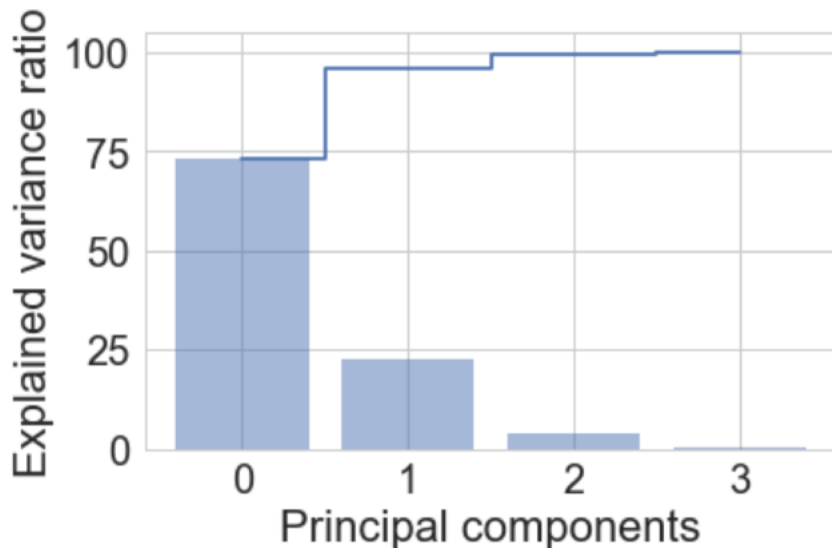
2) Compute SVD: $A = U\Sigma V^T$

3) Principal components: variances = singular values squared

```
1 U, S, Vt = np.linalg.svd(Z, full_matrices=False)
2 variances = S**2
3 print(variances)
```

```
[ 434.85617466  136.19054025  21.86677446   3.08651063]
```

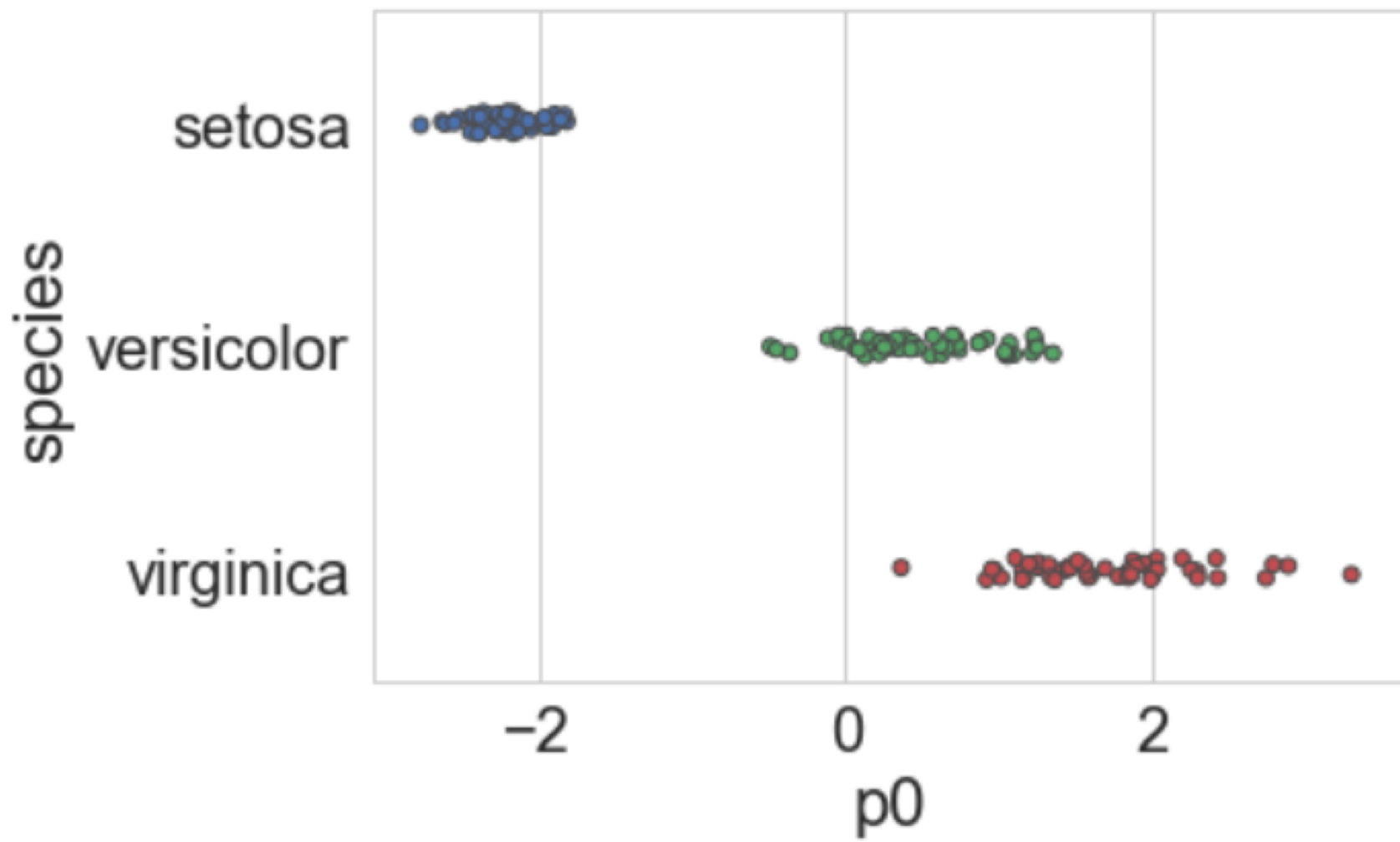
Explained variance: $\text{expvar}_i = \frac{\text{variance}_i}{\text{sum}(\text{variance})}$



— cumulative explained variance
■ individual explained variance

What is the adequate value for k ?

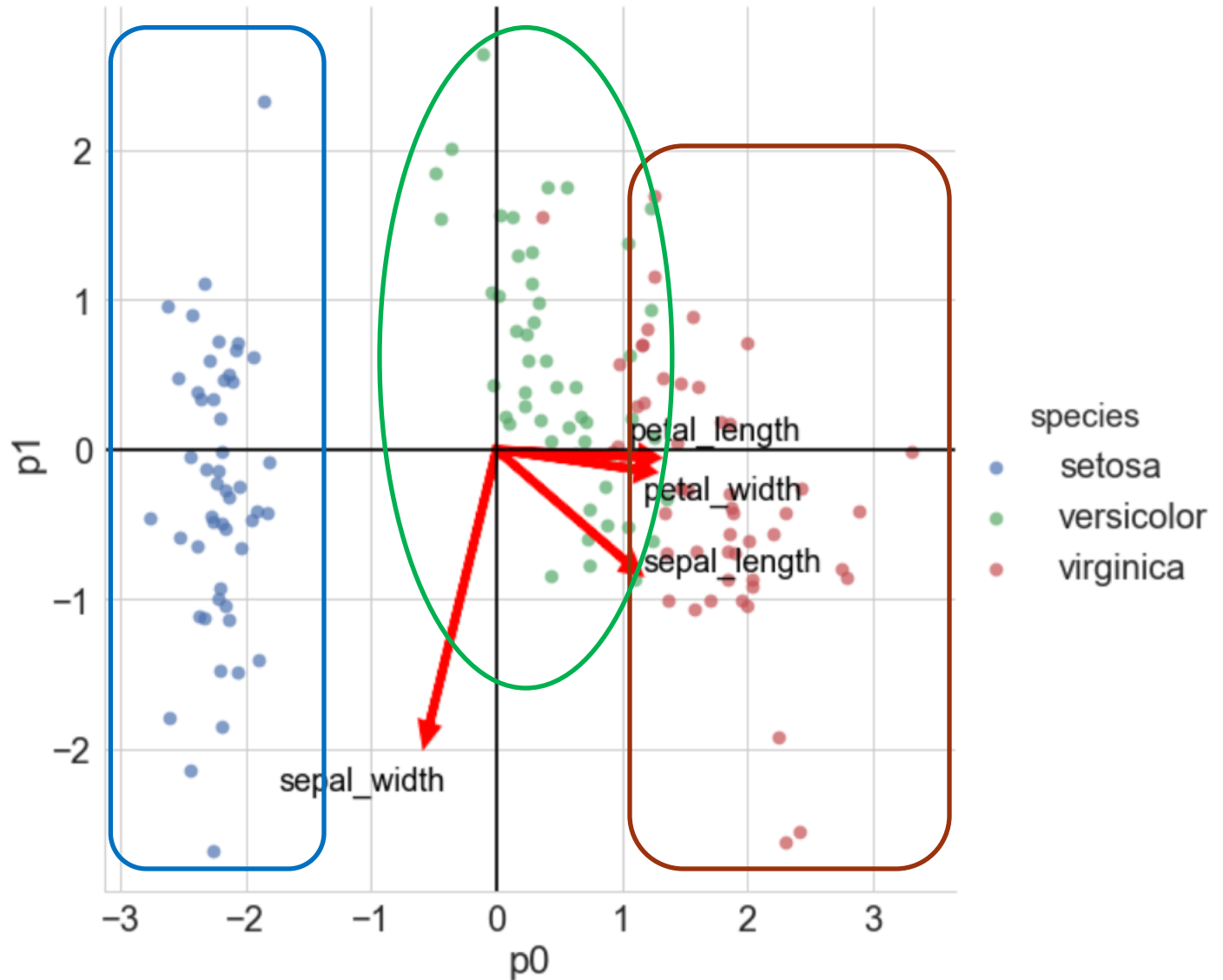
Note that the first two principal components account for about 96% of the variance. It makes sense here to make $k = 2$



5) New REDUCED dataset:

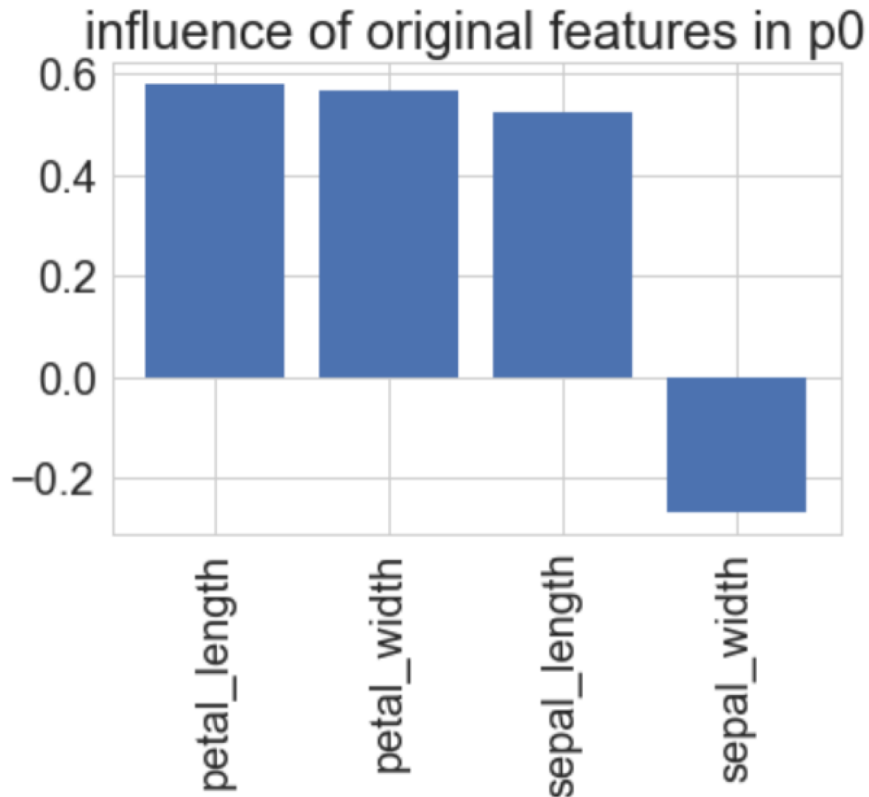
```
1 zstar = z@V[:, :2]
```

$$Z^* = \begin{bmatrix} \vdots & \vdots \\ p0 & p1 \\ \vdots & \vdots \end{bmatrix}$$

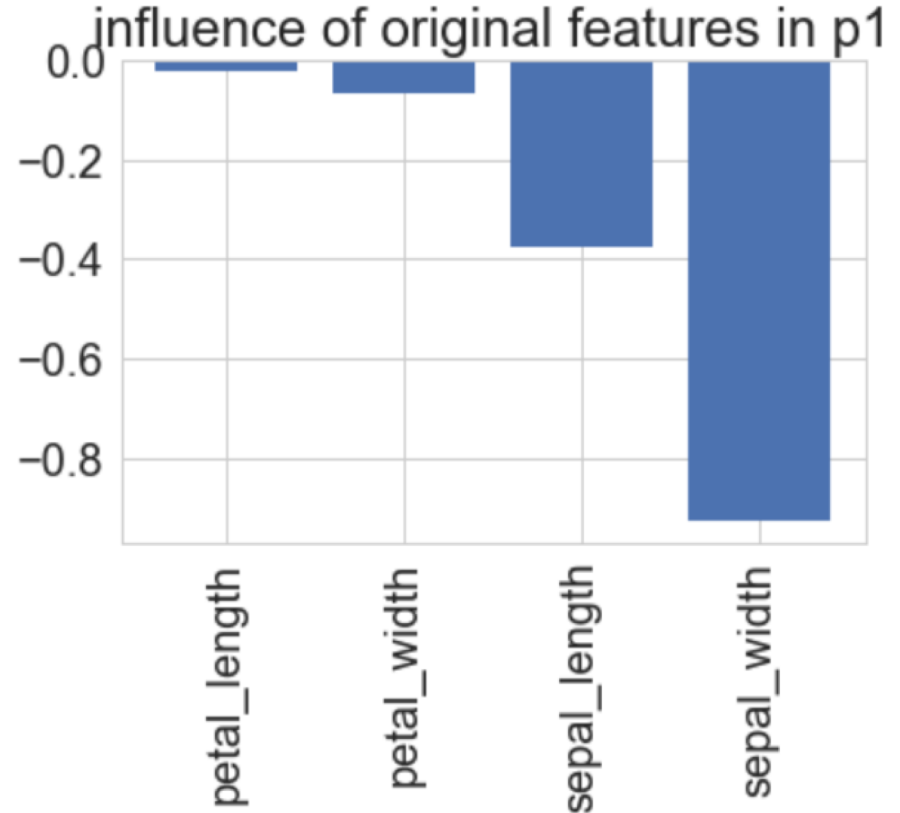


Weight (importance) of each feature in the principal components

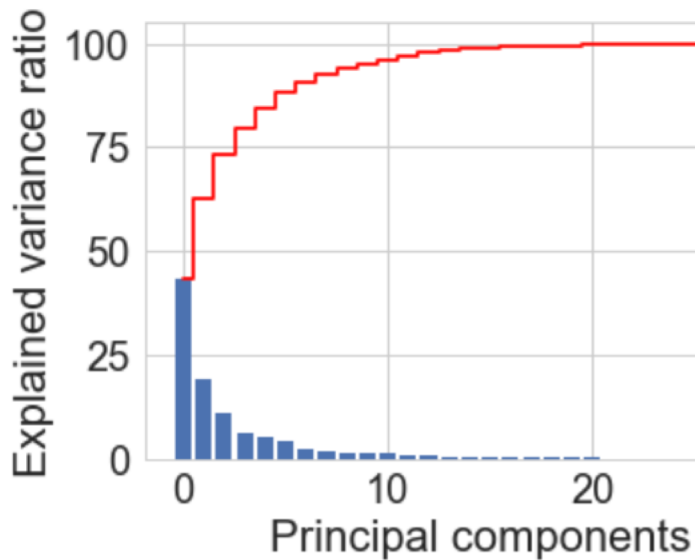
```
1 plt.figure()
2 plt.bar(headers[:4],V[:,0])
3 plt.xticks(rotation=90)
4 plt.title('influence of original features in p0')
```



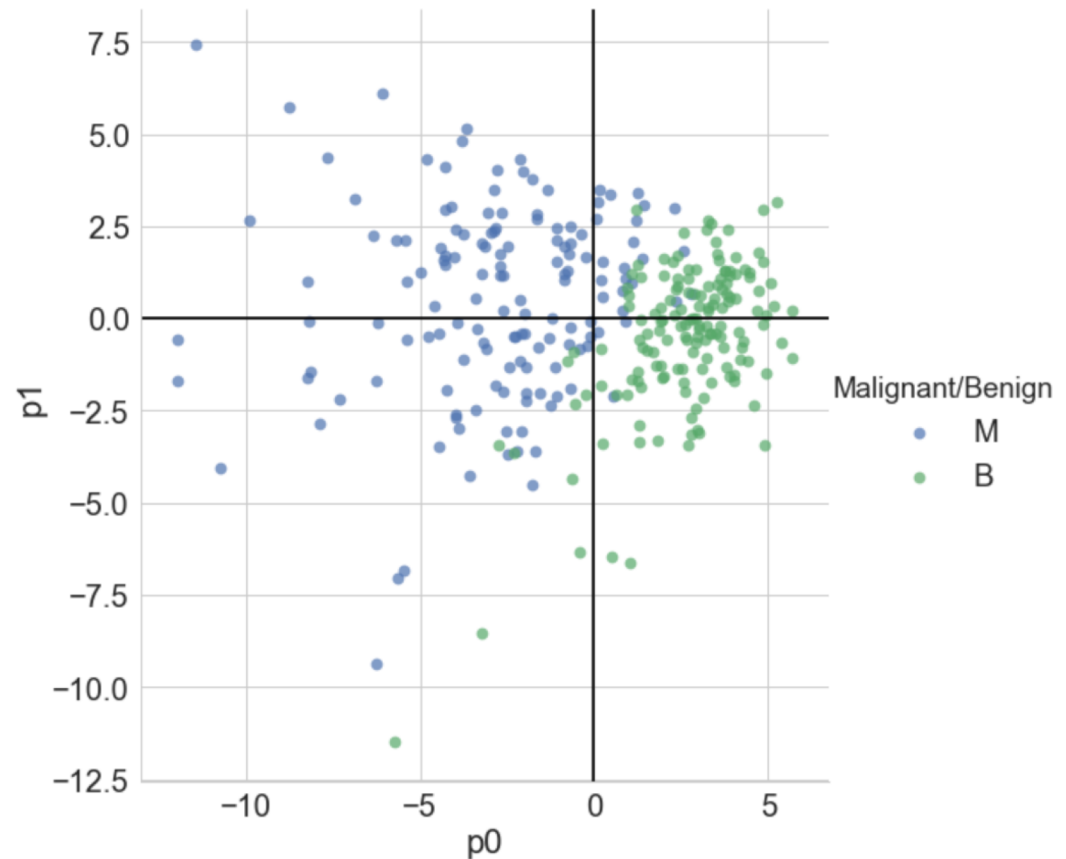
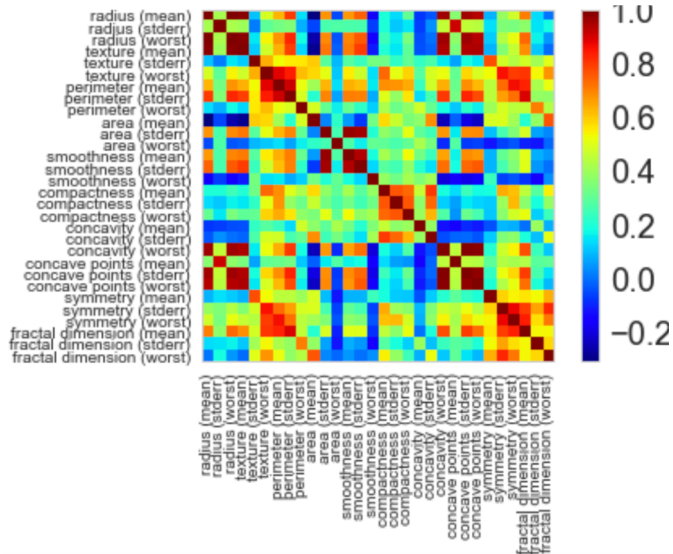
```
1 plt.figure()
2 plt.bar(headers[:4],V[:,1])
3 plt.xticks(rotation=90)
4 plt.title('influence of original features in p1')
```

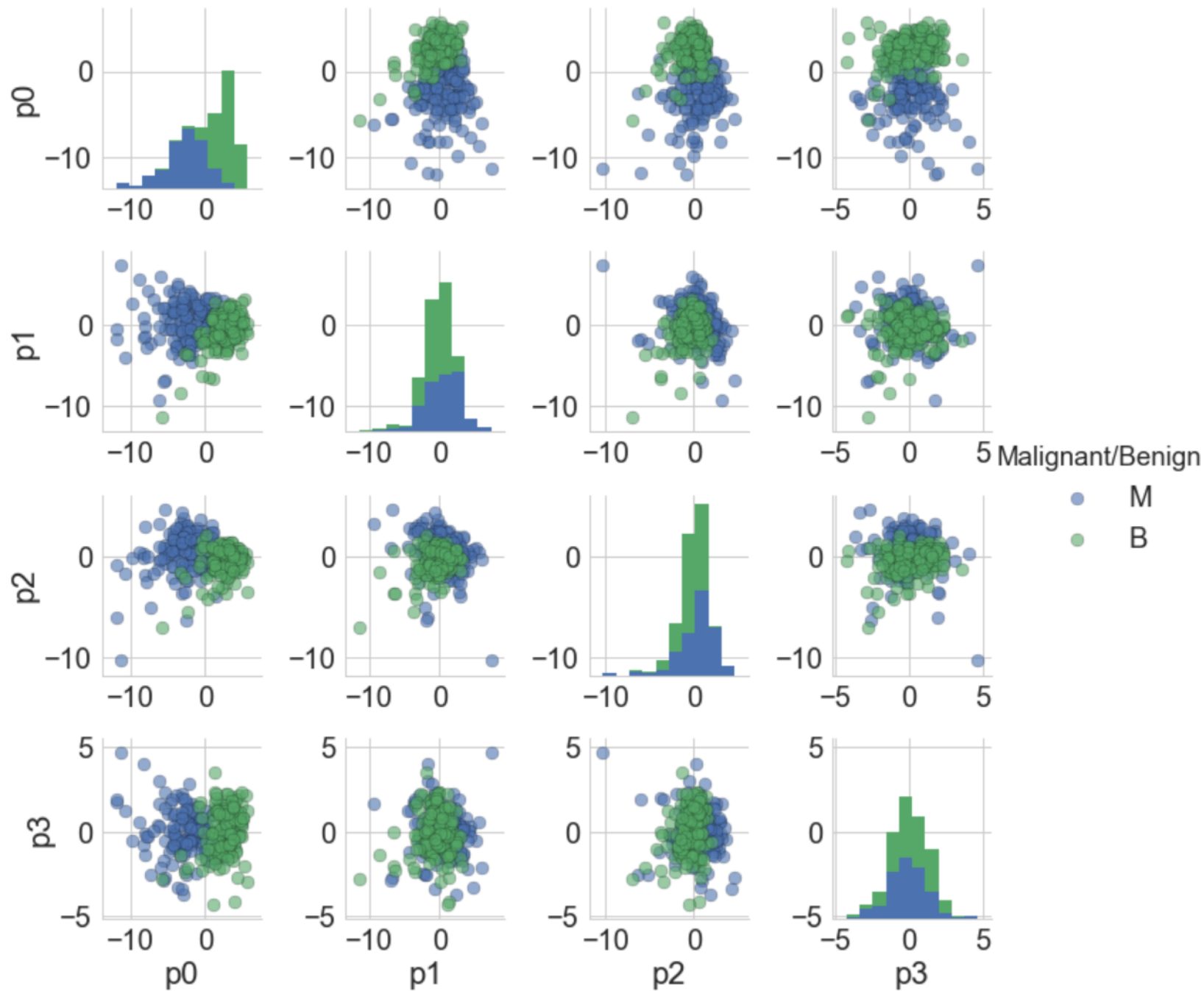


Let's go back to a dataset with many features!



— cumulative explained variance
 ■ individual explained variance





influence of original features in p0

