

Data analytics for personalized genomics and precision medicine

Feature selection & dimension reduction

Lecturer: Yu LI (李煜) from CSE, Writer: Chan Tung Man 1155142328

Liyu95.com, liyu@cse.cuhk.edu.hk

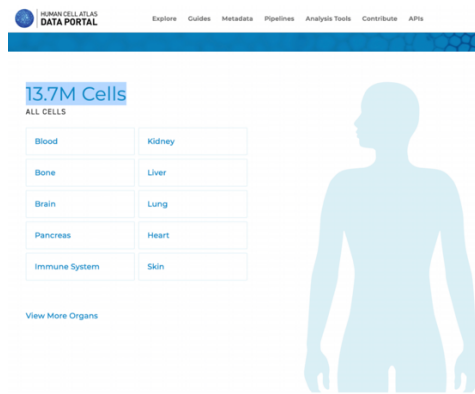
Wednesday, 11 October

The reason of feature selection & dimension reduction

➤ The bio-data always very huge, require lots of storage to store the data.

✓ For example:

❖ Gene expression profile: 25000 gene (features) * 13.7 million cells (cells/ data point) = **1.2 TB!**



❖

✓ For common: normally deal with 10000 cells with 25000 gene.

❖ $25000 * 10000 = 875 \text{ MB}$

➤ We don't need to deal with all the data, some of the data is useless.

✓ For example:

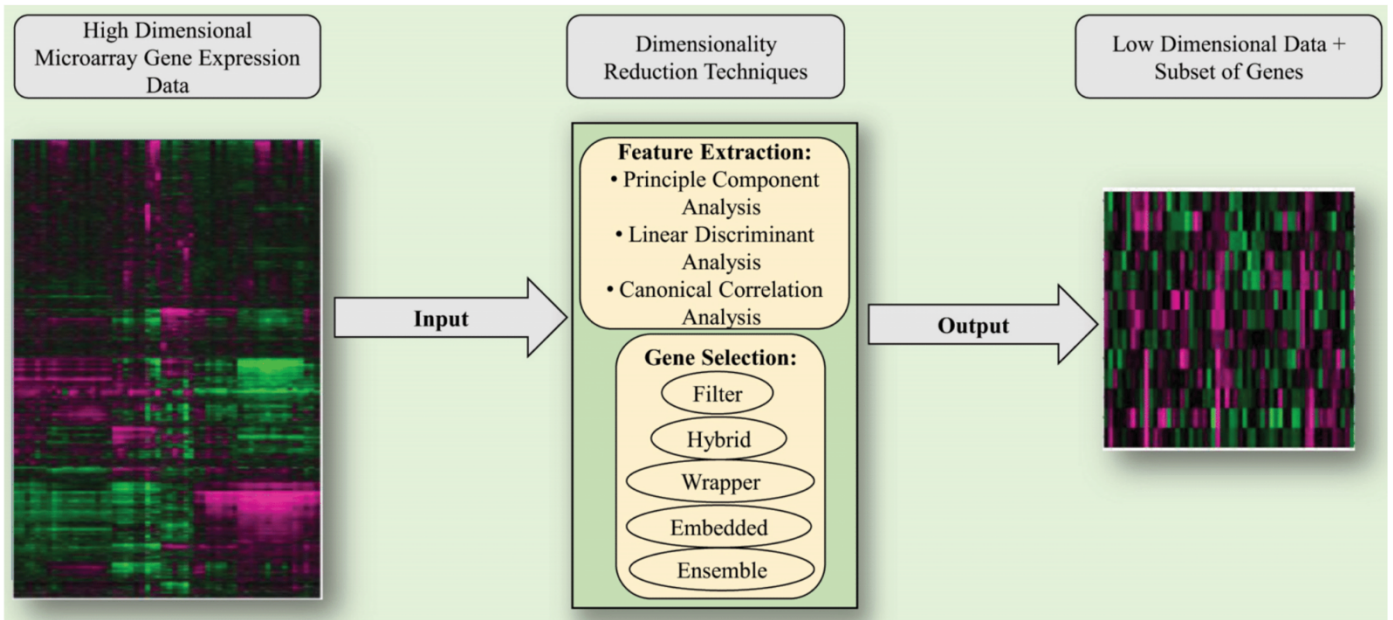
❖ **Irrelevant genes**: we don't have to include them in our analysis.

❖ **Highly correlated genes**: we don't have to include all of them.

❖ **Some genes are complementary**: combine two values into one value may be more useful.

➤ Reduce the data size, more useful and friendly to use the data set.

⇒ Therefore, this procedure is called the **feature (gene) selection and dimension reduction**.



The Left picture can clearly see that the data matrix/ data dimension is huge as compare with the right data matrix.

The color in right/ output matrix is more regularly with smaller dimension represent some useless data points are removed.

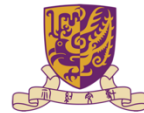
The middle part represents the reduction techniques and algorithm, based on this method for example PCA with different gene selection. The input data matrix can dimension into output matrix.

Benefit

- Data compression
 - ✓ Efficient storage and retrieval
 - ❖ Reduce the size of data matrix and Irrelevant data.
- Improve prediction performance.
 - ✓ Remove unrelated inputs.
 - ❖ Delete useless and irrelevant data.
- Understand the prediction results.
 - ✓ What gene are related to the cancer prediction?
- Facilitate data visualization.
 - ✓ 25000D to 2D, understand the distance between cells visually.

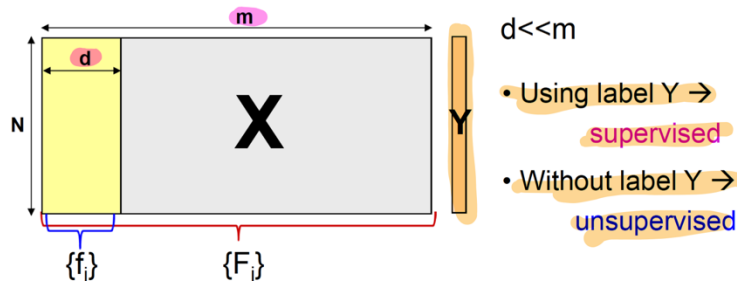
Feature selection

Feature Selection/Extraction



❖ Thousands to millions of **low-level features**: **select/extract** the **most relevant one** to build better, faster, and easier to understand learning machines.

m features reduce to d features



➤ Dimension

Yu Li

Lecture 10-19

1. M numbers of features reduced to d number of features.
 - ❖ $m \gg d$
2. using label $Y \rightarrow$ supervised / without label $Y \rightarrow$ unsupervised
3. N the number of cell/ data pt.

How to reduce dimensionality

➤ Feature selection

Choose the best subset genes from all the genes:

1. Feature ranking
2. Filter/ wrapper

➤ Feature extraction

Extract the feature by linear or non-linear combination.

✓ New gene = gene 1 + gen2

Feature ranking

- Discover the most relevant features and build a better, faster, and easier machines.
- Measurement:
 1. Correlation between feature and class
- ✓ If highly related to class, more useful feature

✓ Example: weight vs gender = 0.714, height vs gender 0.812

❖ Height is more related to class.

2. Mutual information $I(i)$

✓ The higher $I(i)$, the attribute is more related to the class.

3. Fisher score F

✓ The higher F , the attribute is more related to the class.

Problems and issues of individual features ranking.

➤ Relevance vs usefulness:

✓ They are not the same and not directly relation.

➤ Selection of a redundant subset

✓ K best features \neq best k features

❖ Best two features \neq the best combination of the features

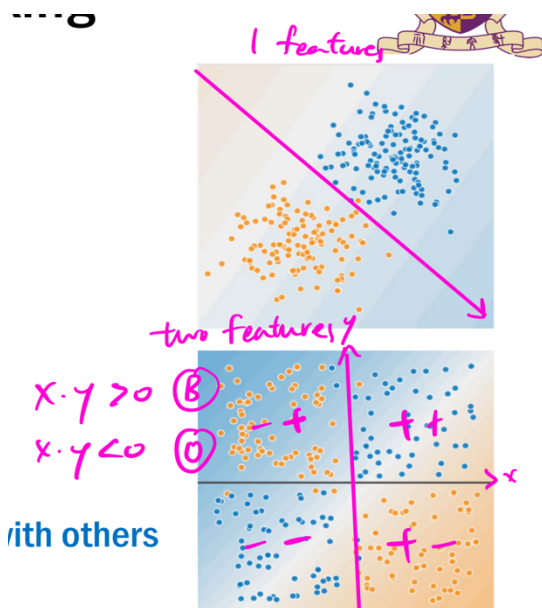
✓ Example: ranking Major (3) \rightarrow weight (2) \rightarrow Height (1)

❖ Height & weight \lll height & major

➤ A variable that is useless by itself can be useful with others.

✓ Salary: occupation + age

❖ Only age is no correlate with salary but combine with occupation, it has the meaning.



Lecture 10-24

➤

➤ two features have two axis to separate in 4 areas.

Filter

- classification performance is not involved in the selection loop.
- variance data required -> age, height, weight.
- Higher variance with more useful information, take a variance threshold to filter our data.

➤ Information gain: Features should be different *G1*

	G1	G2	G3	G4
S1	10		6	8
S2	10		7	8
S3	10		8	6
S4	10		9	5

also + 1
Yu Li

- Pink: delete the gene with same data pt. (information gain)
- Orange: G2 and G3 both increment is 1, choose either one is enough (variance)

Wrapper

Example of wrapper



❖ Wrapper

- Using the classification performance to guide selection
- Computational expensive
- Recursive feature elimination
- Sequential feature selection



	G1	G2	G3	G4	Cancer
S1	10	2	6	8	Yes
S2	10	3	7	8	Yes
S3	10	4	8	6	No
S4	10	5	9	5	No

C₄² times to do CFV
G1 G2 G2 G3 G3 G4
G1 G3 G2 G4
G1 G4
 Yu Li

1. No feature
2. Find the first best feature using cross-fold validation
3. Add the second feature using cross-fold validation
4. ...
5. Until the new feature does not improve the performance

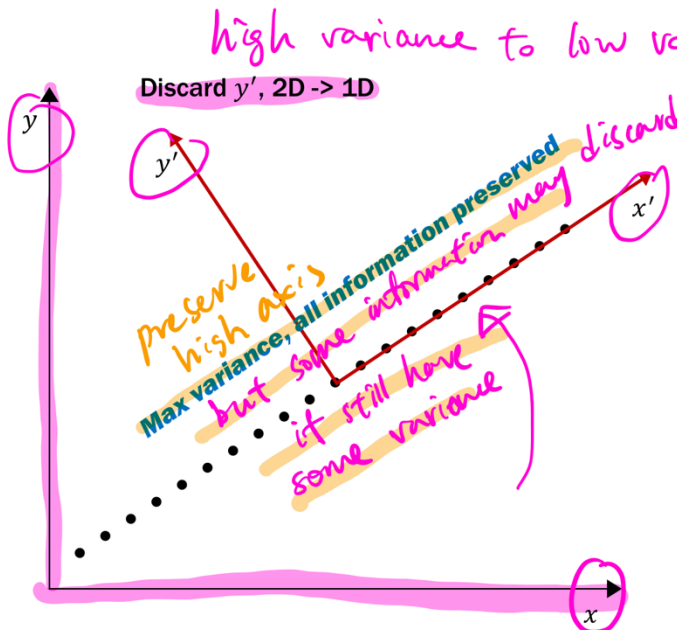
- Use cross-fold validation to update the best feature combination.

Dimension reduction

Principal Components analysis (PCA)



Principal components analysis (PCA)



A two dimensional scatter of points that show a high degree of correlation

We care about **variance** (information) and **distance**

Person	Height (m)	Weight (kg)	Age	Gender
P1	1.79	75	20	M
P2	1.64	54	20	F
P3	1.70	63	20	M
P4	1.88	78	20	M
P5	1.75	70	20	??

Dimension

Yu Li

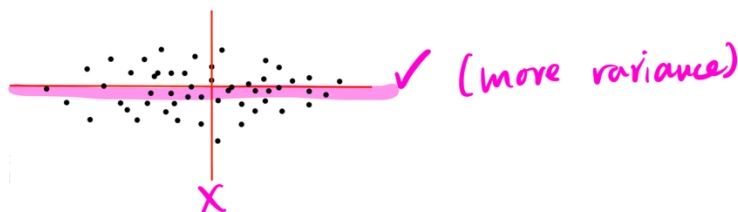
Lecture 10-34

- From 2D xy chart to 1D y' chart
- From high variance to low variance,
- Care about the variance and the distance between point and axis
- The scatter point are directly proportional = high correlation

Principal components analysis (PCA)



- ❖ After vector space transform, we have more “efficient” description
 - 1st dimension captures **max variance**
 - 2nd dimension captures the max amount of **residual variance**, at right angles (orthogonal) to the first
- ❖ The **1st dimension** may capture **so much of the information content in the original data set that we can ignore the remaining axis**



Dimension

Yu Li

Lecture 10-36

- The 1st dimension/ axis is the horizontal one, more data point touch the axis, have more variance.

➤ The 2nd axis perpendicular to the 1st axis, less variance so we can ignore it.

How to do PCA?

How to do PCA?



- ① ❖ Suppose we have a n by d data matrix, X . We first normalize each feature to make the average of each feature 0. Then, we get X'
- ② ❖ Then, we calculate the covariance matrix of X'
 - $\Sigma = \frac{1}{n-1} X'^T X'$, Σ : a d by d matrix
- ③ ❖ Find the eigenvectors and eigenvalues of Σ
 - M eigenvectors with the M largest eigenvalues
 - Principal components
- ④ ❖ Project the data to the M eigenvectors' direction
 - $\hat{X} = X'P$

Variance?
 How important?

direction
 information in those direction

Person	Height (m)	Weight (kg)	Age	Gender
P1	1.79	75	20	M
P2	1.64	54	20	F
P3	1.70	63	20	M
P4	1.88	78	20	M
P5	1.75	70	20	??

1. Normalization
2. Covariance matrix
3. Eigenvectors and eigenvalues
4. Largest eigenvalues
5. Project the data to the eigenvector of its direction

Example:

X	A	B	C
X1	1	1	1
X2	2	2	2
X3	3	3	3

1. Calculate average of each feature and do Normalization:

Average = $1 + 2 + 3 / 3 = 2$

X'	A	B	C
X1	1-2 = -1	1-2 = -1	1-2 = -1
X2	2-2 = 0	2-2 = 0	2-2 = 0

X3	3-2 =1	3-2 =1	3-2 =1
----	--------	--------	--------

2. Covariance matrix

$$\Sigma = \frac{1}{n-1} (X'^T)X'$$

$$X'^T = \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix}$$

$$X' = \begin{bmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix}$$



2 × 2 Matrix Multiplication

$$\begin{bmatrix} a_1 & b_1 \\ c_1 & d_1 \end{bmatrix} \times \begin{bmatrix} a_2 & b_2 \\ c_2 & d_2 \end{bmatrix} = \begin{bmatrix} a_1 a_2 + b_1 c_2 & a_1 b_2 + b_1 d_2 \\ c_1 a_2 + d_1 c_2 & c_1 b_2 + d_1 d_2 \end{bmatrix}$$

3 × 3 Matrix Multiplication

$$\begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} \times \begin{bmatrix} j & k & l \\ m & n & o \\ p & q & r \end{bmatrix} = \begin{bmatrix} (aj + bm + cp) & (ak + bn + cq) & (al + bo + cr) \\ (dj + em + fp) & (dk + en + fq) & (dl + eo + fr) \\ (gj + hm + ip) & (gk + hn + iq) & (gl + ho + ir) \end{bmatrix}$$

$$X'^T X' = \begin{bmatrix} 2 & 2 & 2 \\ 2 & 2 & 2 \\ 2 & 2 & 2 \end{bmatrix}$$

$$\Sigma = \frac{1}{n-1} (X'^T)X' = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

3. Eigenvectors and eigenvalues

Example of PCA



$$X'$$

X1	-1	-1	-1
X2	0	0	0
X3	1	1	1

$$\Sigma = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

- ❖ We first normalize each feature to make the average of each feature 0. Then, we get X'
- ❖ Then, we calculate the covariance matrix of X'
 - $\Sigma = \frac{1}{n-1} X'^T X'$, Σ : a d by d matrix
- ❖ Find the eigenvectors and eigenvalues of Σ
- ❖ M eigenvectors with the M largest eigenvalues
 - Principal components
- ❖ Project the data to the M eigenvectors' direction
 - $\hat{X} = X'P$

$$\Sigma * V = \lambda * V$$

$$|\Sigma - \lambda I| = 0$$

$$\begin{vmatrix} 1-\lambda & 1 & 1 \\ 1 & 1-\lambda & 1 \\ 1 & 1 & 1-\lambda \end{vmatrix} = 0$$

$$\left\{ \begin{array}{l} \lambda_1 = 3 \\ \lambda_2 = 0 \\ \lambda_3 = 0 \end{array} \right.$$

$$\begin{vmatrix} x & x' & x'' \\ y & y' & y'' \\ z & z' & z'' \end{vmatrix} = xy'z'' + x'y''z + x''yz' - xy''z' - x'y'z'' - x''y'z$$

$$(1-\lambda)^3 + 1 + 1 - (1-\lambda) - (1-\lambda) - (1-\lambda) = 0$$

$$\lambda = 3 \text{ or } \lambda = 0$$

eigenvalues

Dimension

Yu Li

Lecture 10-44

$$\Rightarrow I = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

➤ If eigenvalues = 0, represent **No information!**

Example of PCA



$$X'$$

X1	-1	-1	-1
X2	0	0	0
X3	1	1	1

$$\Sigma = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

- ❖ We first normalize each feature to make the average of each feature 0. Then, we get X'
- ❖ Then, we calculate the covariance matrix of X'
 - $\Sigma = \frac{1}{n-1} X'^T X'$, Σ : a d by d matrix
- ❖ Find the eigenvectors and eigenvalues of Σ
- ❖ M eigenvectors with the M largest eigenvalues
 - Principal components
- ❖ Project the data to the M eigenvectors' direction
 - $\hat{X} = X'P$

$$\lambda_1 = 3 \quad V_1 = \begin{bmatrix} \frac{\sqrt{3}}{3} \\ 3 \\ \frac{\sqrt{3}}{3} \\ 3 \\ \frac{\sqrt{3}}{3} \\ 3 \end{bmatrix}$$

preserve

$$\lambda_{2,3} = 0 \quad V_{2,3} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

no information

➤ Dimension

Yu Li

Lecture 10-46

4. Largest eigenvalues

Example of PCA



$$X'$$

X1	-1	-1	-1
X2	0	0	0
X3	1	1	1

$$\Sigma = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

- ❖ We first normalize each feature to make the average of each feature 0. Then, we get X'
- ❖ Then, we calculate the covariance matrix of X'
 - $\Sigma = \frac{1}{n-1} X'^T X'$, Σ : a d by d matrix
- ❖ Find the eigenvectors and eigenvalues of Σ
- ❖ M eigenvectors with the M largest eigenvalues
 - Principal components
- ❖ Project the data to the M eigenvectors' direction
 - $\hat{X} = X'P$

$$\lambda_1 = 3 \quad V_1 = \begin{bmatrix} \frac{\sqrt{3}}{3} \\ \frac{\sqrt{3}}{3} \\ \frac{\sqrt{3}}{3} \end{bmatrix}$$

$$\lambda_{2,3} = 0 \quad V_{2,3} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$2D \quad P = \begin{bmatrix} \frac{\sqrt{3}}{3} & 0 \\ \frac{\sqrt{3}}{3} & 0 \\ \frac{\sqrt{3}}{3} & 0 \end{bmatrix} \quad \hat{X} = X'P = \begin{bmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix} * \begin{bmatrix} \frac{\sqrt{3}}{3} \\ \frac{\sqrt{3}}{3} \\ \frac{\sqrt{3}}{3} \end{bmatrix} = \begin{bmatrix} -\sqrt{3} & 0 \\ 0 & 0 \\ \sqrt{3} & 0 \end{bmatrix}$$

➤ Dimension

Yu Li

Lecture 10-49

5. Project the data to the eigenvector of its direction

Example of PCA



$$X$$

X1	1	1	1
X2	2	2	2
X3	3	3	3

- ❖ We first normalize each feature to make the average of each feature 0. Then, we get X'
- ❖ Then, we calculate the covariance matrix of X'
 - $\Sigma = \frac{1}{n-1} X'^T X'$, Σ : a d by d matrix
- ❖ Find the eigenvectors and eigenvalues of Σ
- ❖ M eigenvectors with the M largest eigenvalues
 - Principal components
- ❖ Project the data to the M eigenvectors' direction
 - $\hat{X} = X'P$

↓
PCA

X1	$-\sqrt{3}$	0
X2	0	0
X3	$\sqrt{3}$	0

\hat{X}

➤ Dimension

Yu Li

Lecture 10-50