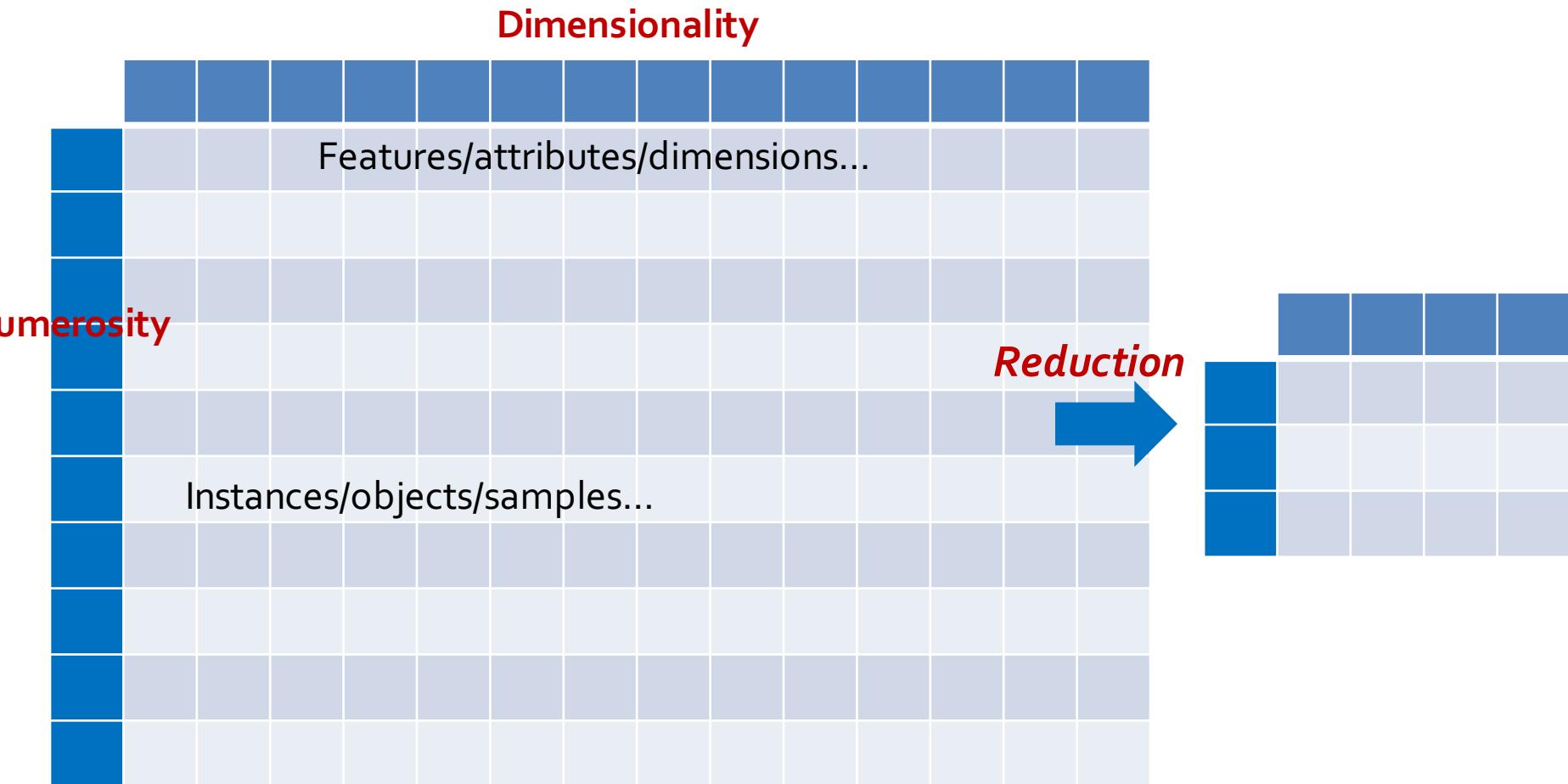


Chapter 3. Data Preprocessing: Data Reduction

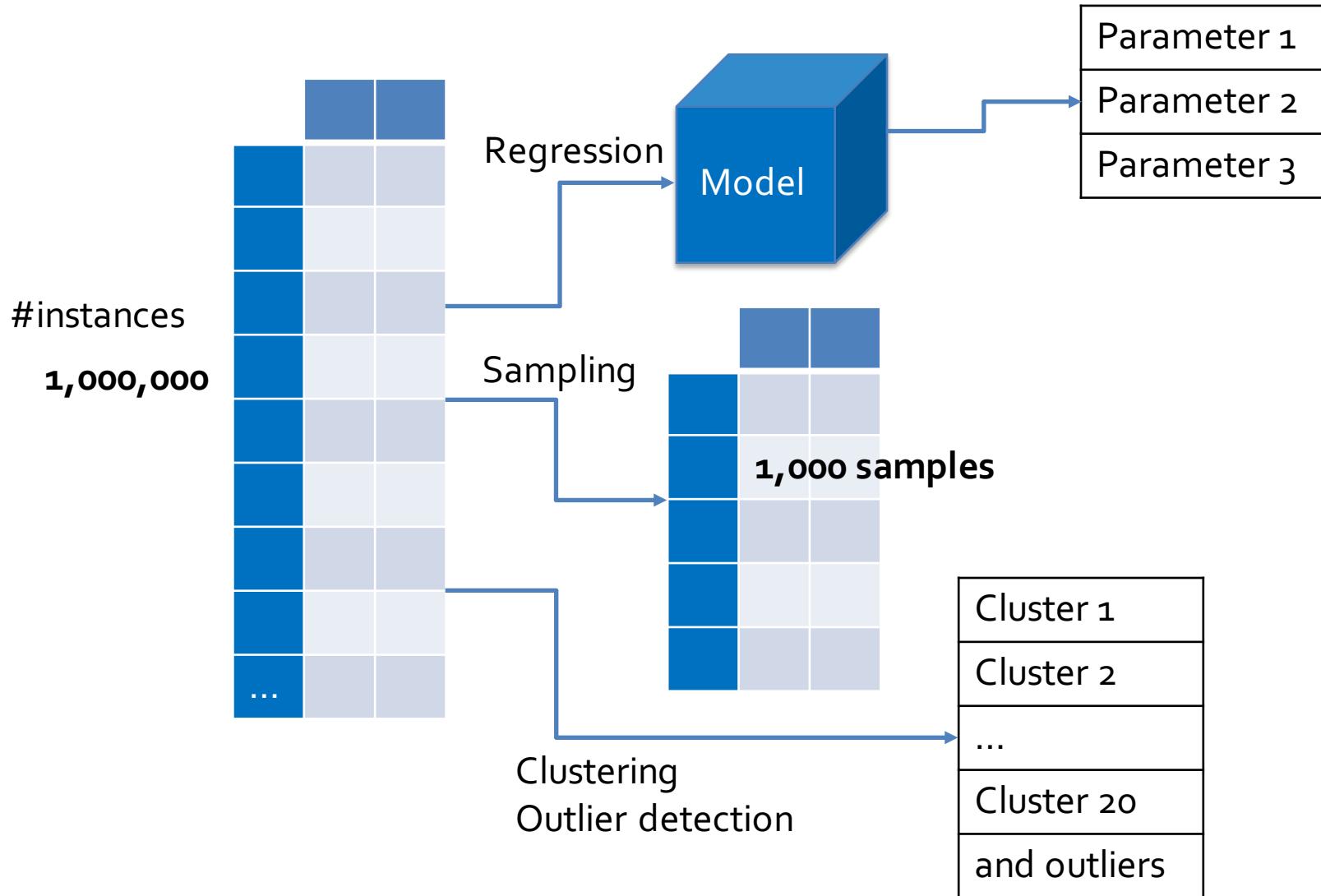
Meng Jiang
Data Science

error23

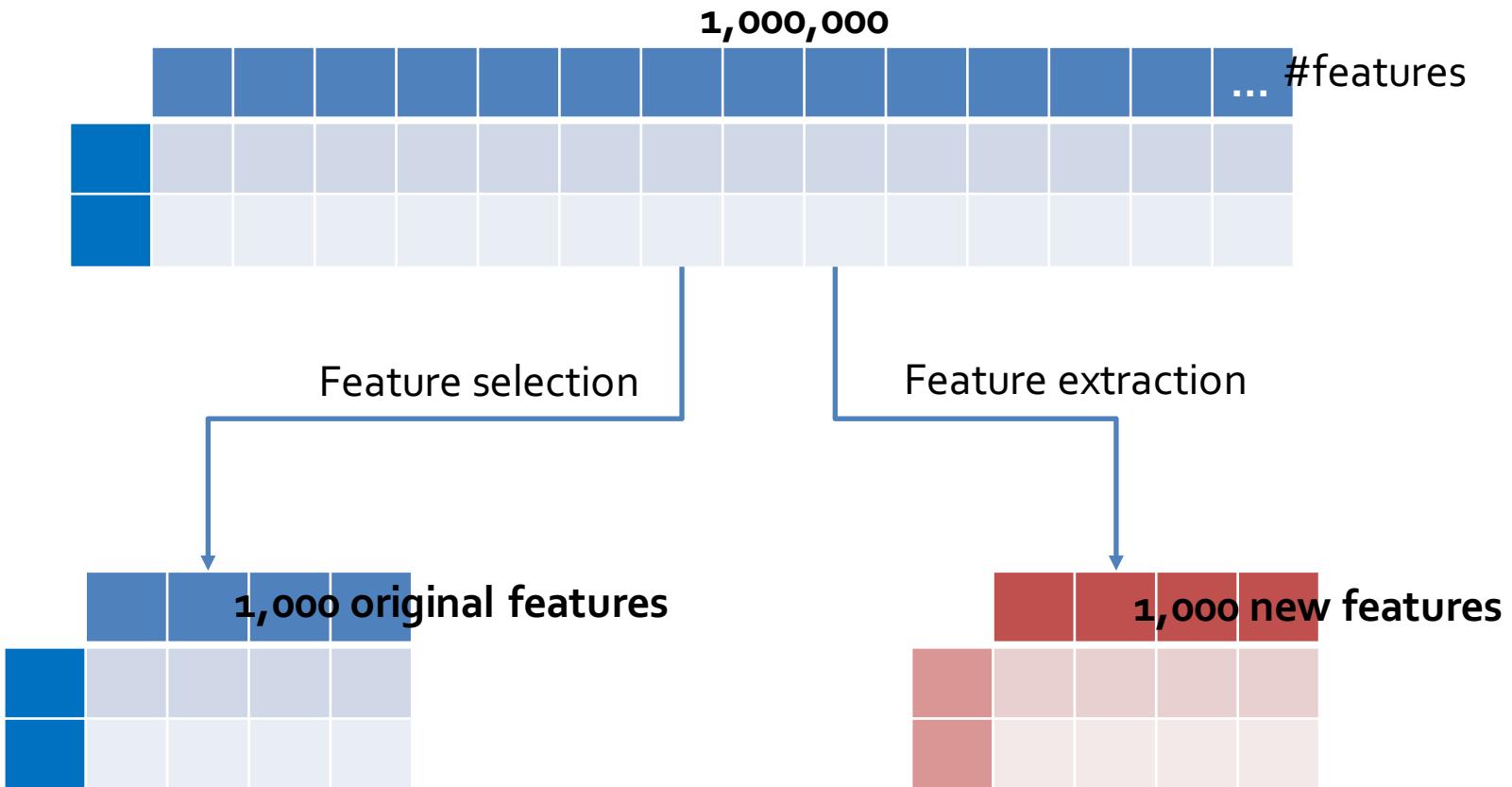
Numerosity and Dimensionality



Numerosity Reduction



Dimensionality Reduction

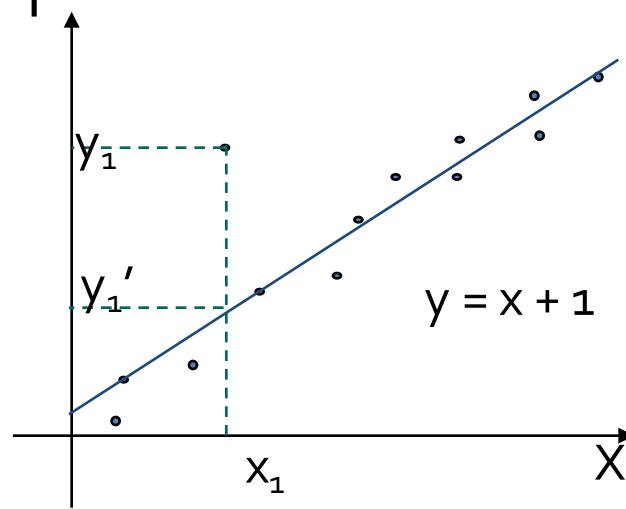


Today: Data Reduction

- Describe **numerosity reduction** (reducing #instances)
 - Parametric methods: Fit some model and estimate model parameters
 - Regression: Describe linear/non-linear regression models
 - Nonparametric methods
 - Histograms
 - Clustering
 - Sampling: Describe stratified sampling
- Describe **dimensionality reduction** (reducing #features)
 - Feature selection
 - Heuristic search
 - Feature extraction
 - Principal component analysis (PCA)
 - Singular value decomposition (SVD)

Regression Analysis

- Regression analysis: A collective name for techniques for the modeling and analysis of numerical data consisting of values
 - of a ***dependent variable*** (also called ***response variable*** or ***measurement***): Y
 - and of one or more ***independent variables*** (also known as ***explanatory variables*** or ***predictors***): X, or $X_1, X_2, \dots X_n$
- Parameters are estimated to give a “**best fit**” of the data
 - Data: (x_1, y_1)
 - Fit of the data: (x_1, y_1')
 - Ex. $y_1' = x_1 + 1$

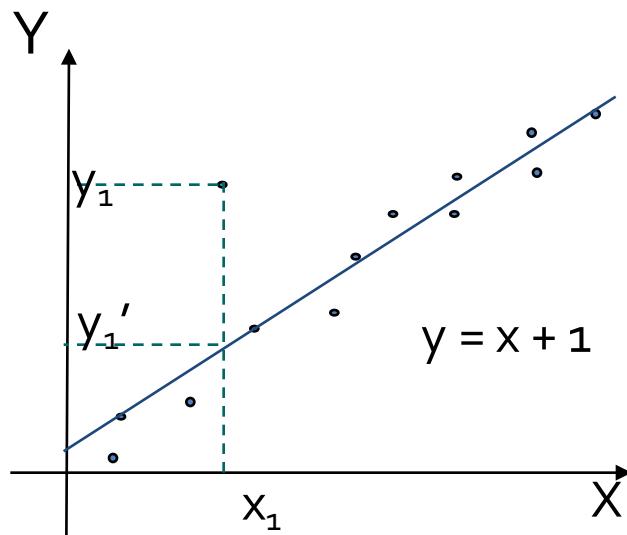


Least Square Method for Regression

- Most commonly the best fit is evaluated by using the ***least square method***, but other criteria have also been used

$$\min g = \sum_{i=1}^n (y_i - y'_i)^2, \text{ where } y'_i = f(x_i, \beta)$$

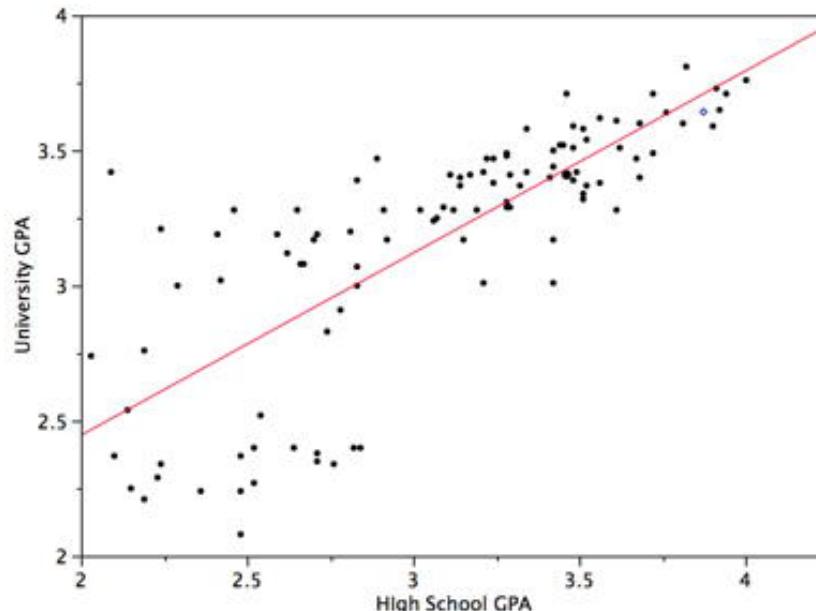
- Used for **prediction** (including forecasting of time-series data), **inference**, and **hypothesis testing**.



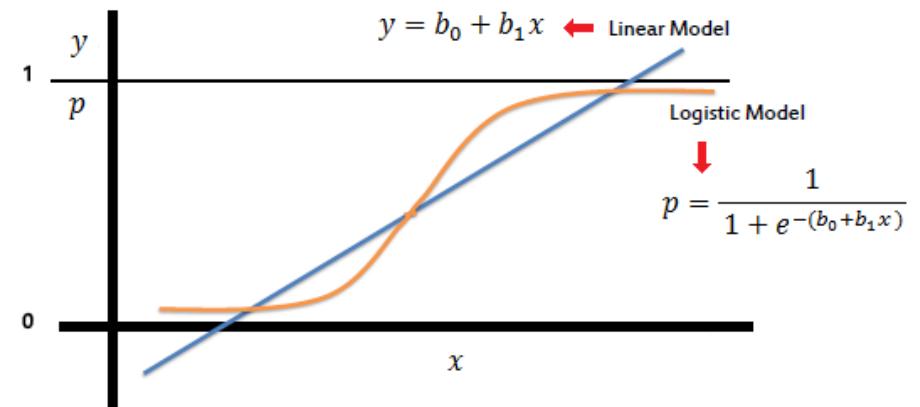
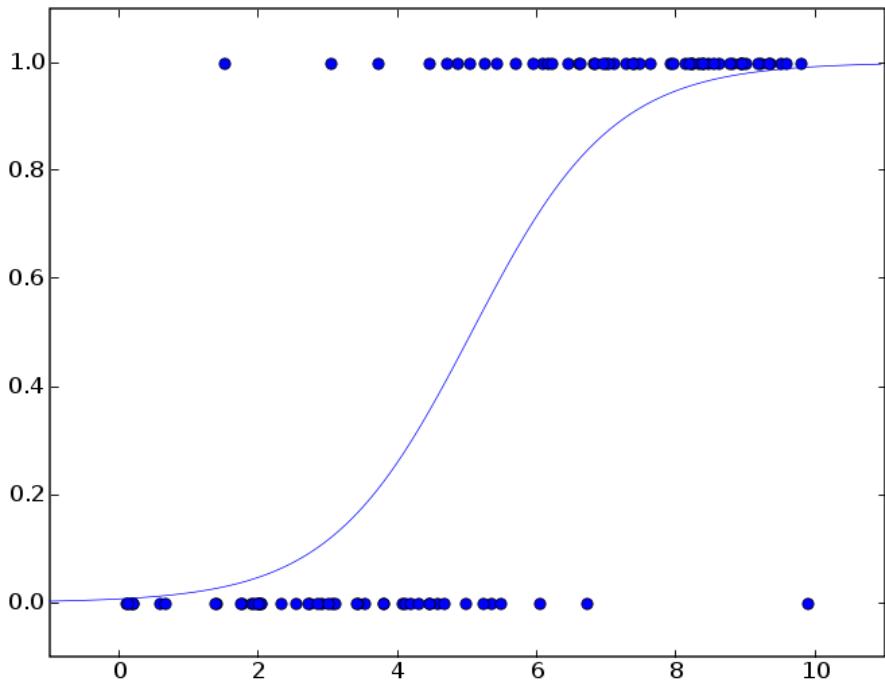
Set up $y = f(x) = \beta_1 x + \beta_2$
Learn β by minimizing the least square error

Linear Regression

- Linear regression: $Y = wX + b$
 - Data modeled to fit a **straight line**
 - Often uses the least-square method to fit the line
 - Two regression coefficients, w and b , specify the line and are to be estimated by using the data at hand

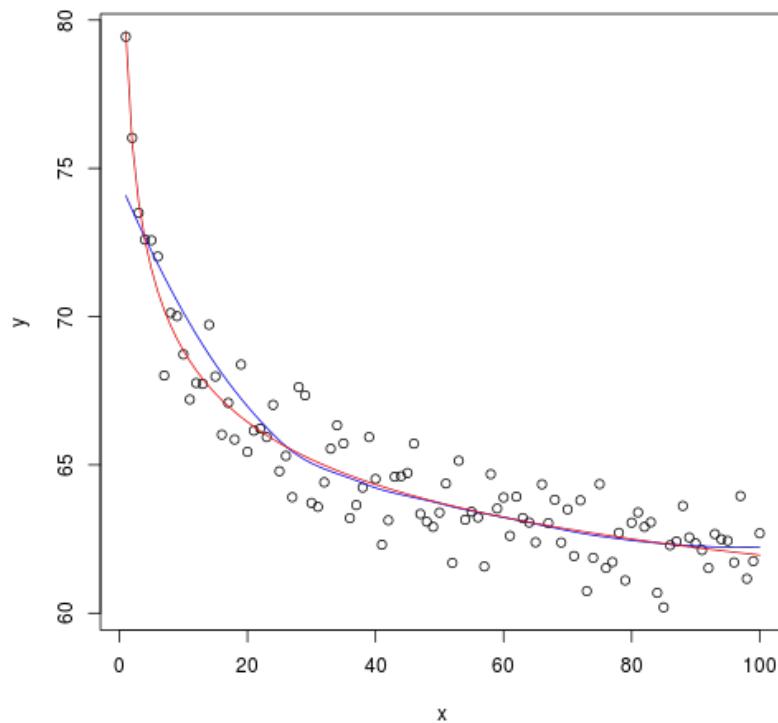


Logistic Regression



Nonlinear Regression

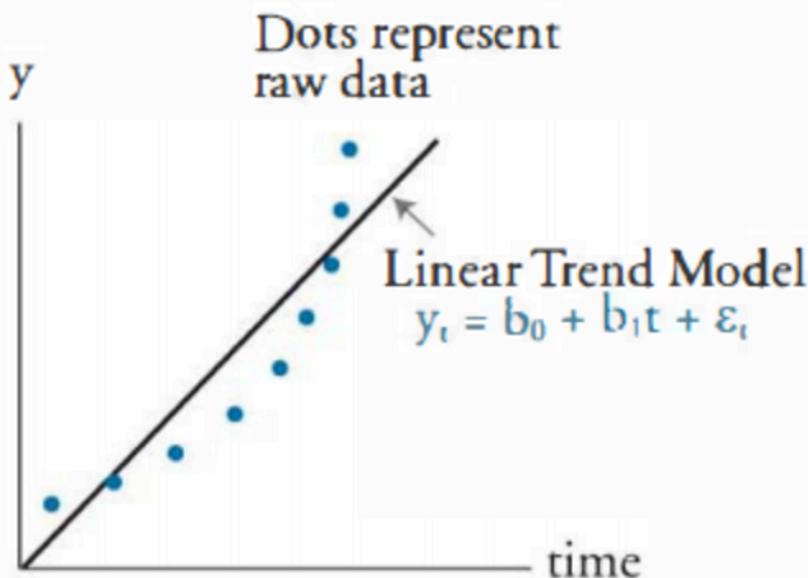
- Nonlinear regression:
 - Data are modeled by a function which is a **nonlinear** combination of the model parameters and depends on one or more independent variables



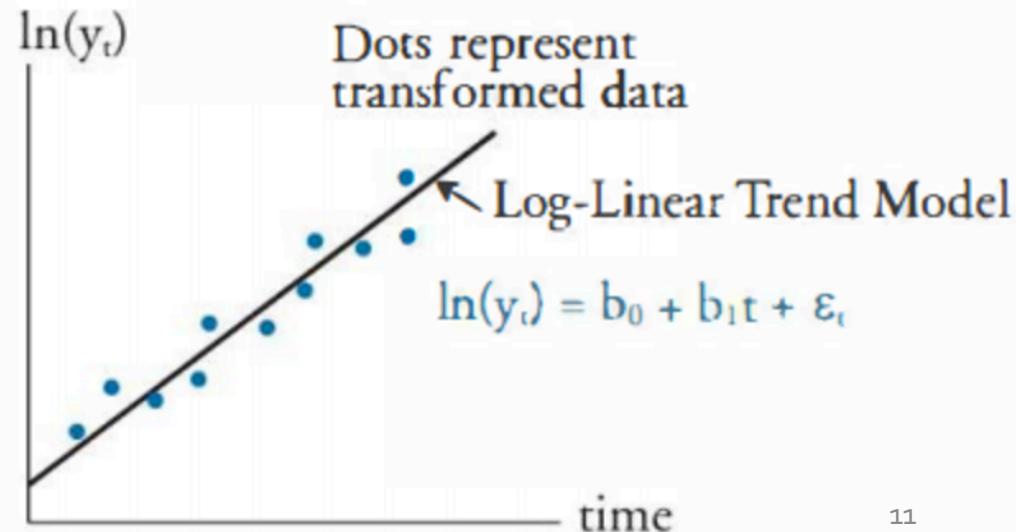
Log-Linear Model

- Log-linear model
 - A math model that takes the form of a **function whose logarithm** is a **linear** combination of the parameters of the model
- Q: How about Log-Log model?*

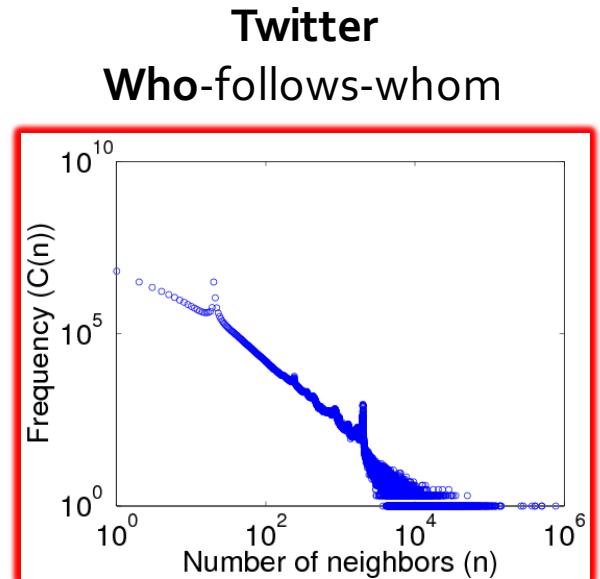
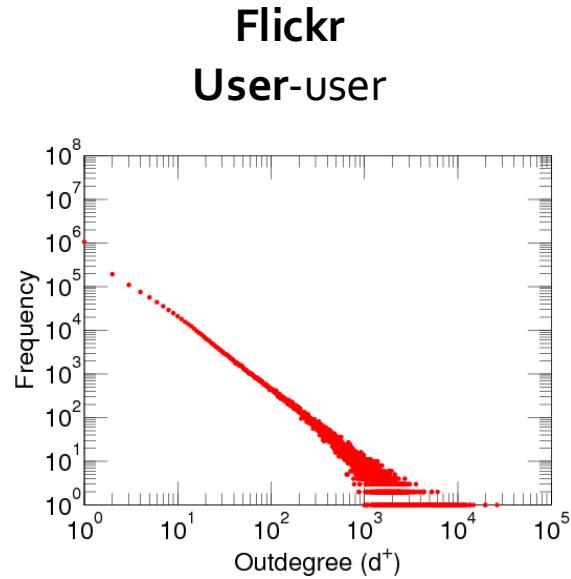
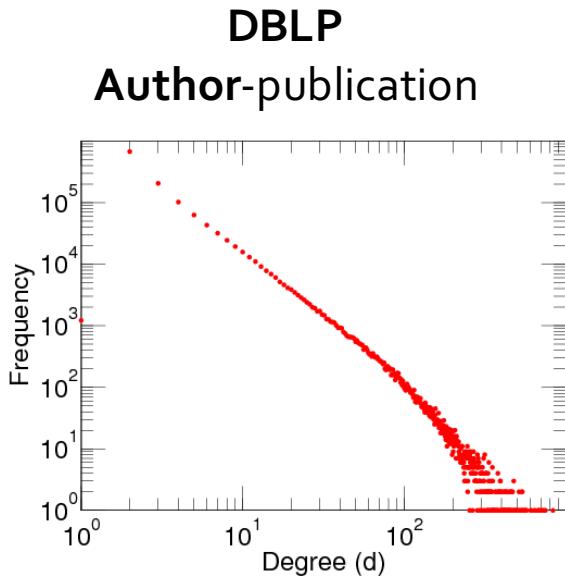
Linear Trend Model



Log-Linear Trend Model



Log-Log: Power Law

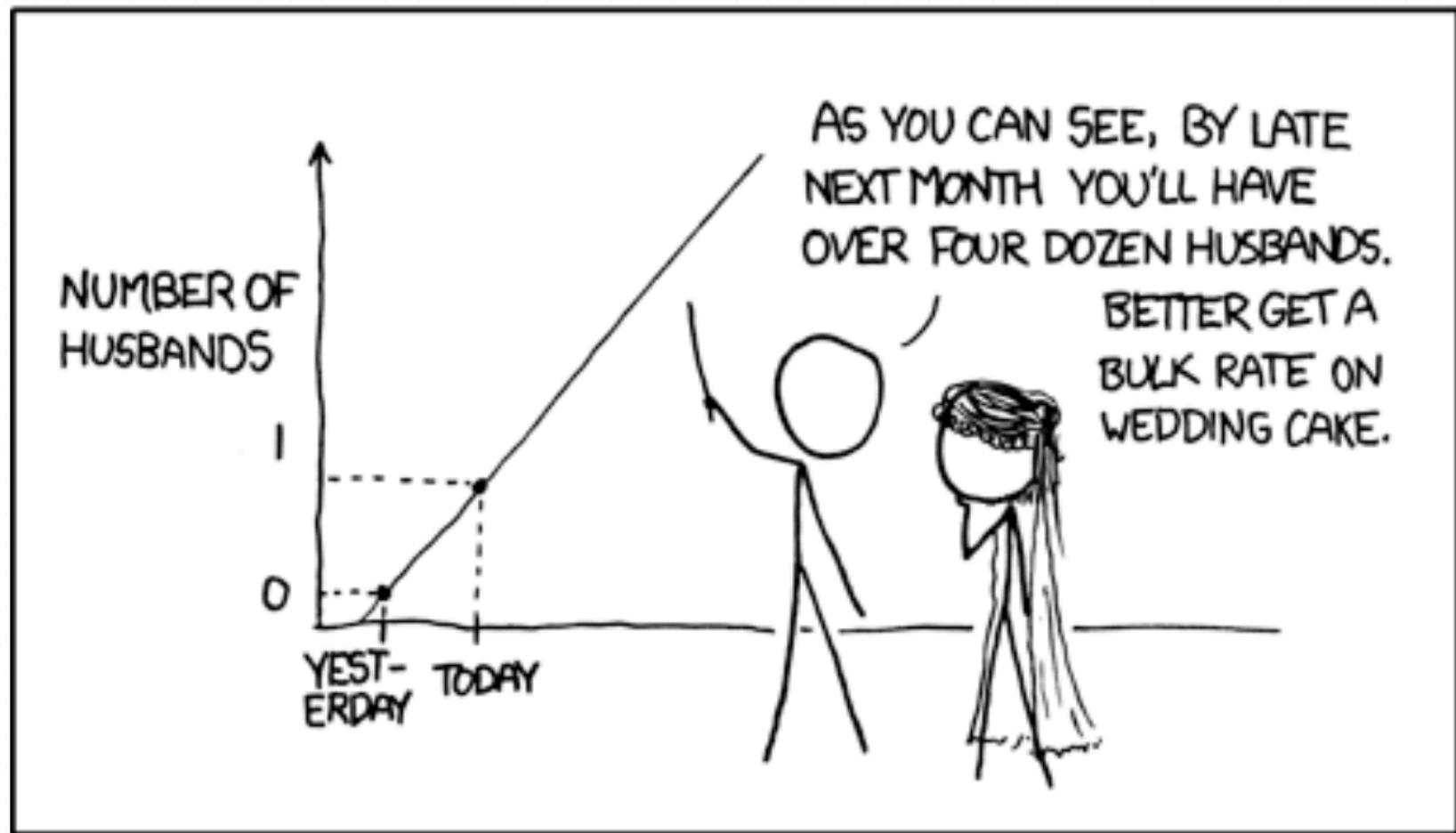


[konect.uni-koblenz.de/networks/]

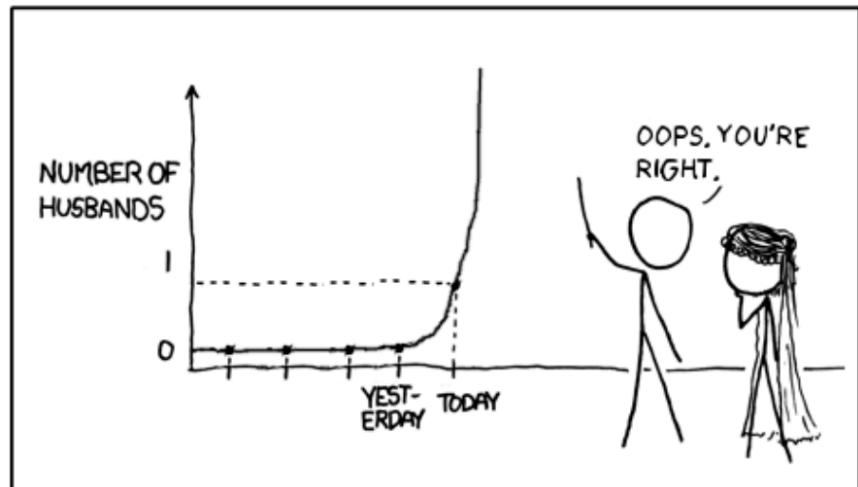
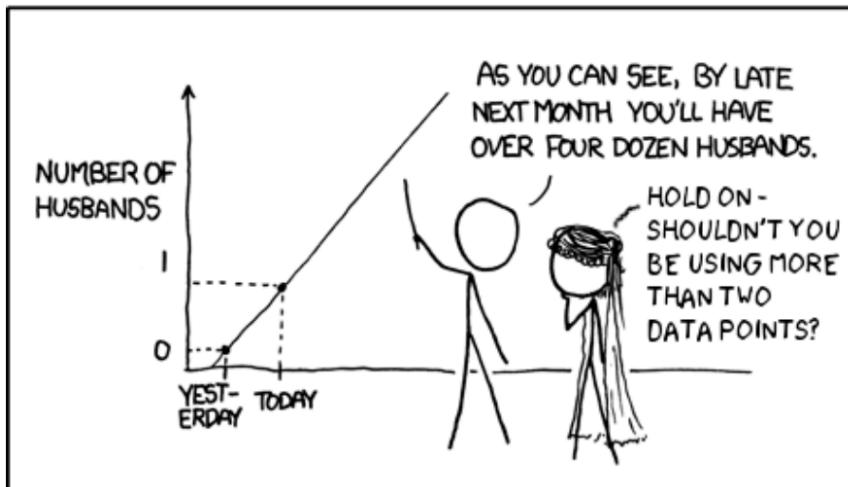
Out-degree distributions (network):

Power-law distributions in networks [Faloutsos et al.
SIGCOMM'99; Chung et al. PNAS'02]

Caution: Extrapolation



Caution: Extrapolation

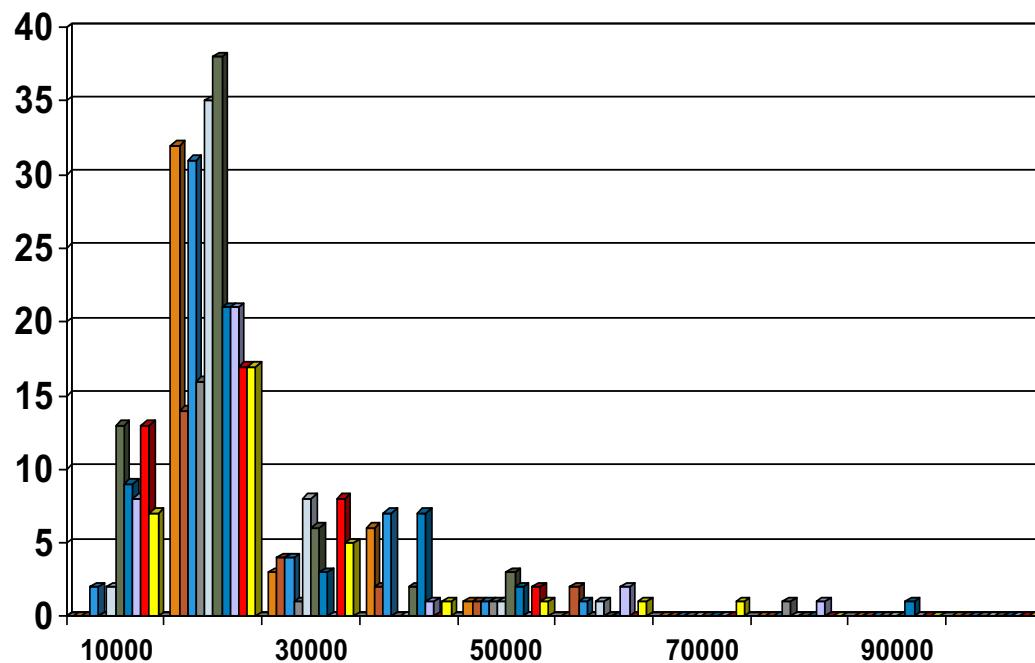


Today: Data Reduction

- Describe **numerosity reduction** (reducing #instances)
 - Parametric methods: Fit some model and estimate model parameters
 - Regression: Describe linear/non-linear regression models
 - Nonparametric methods
 - Histograms
 - Clustering
 - Sampling: Describe stratified sampling
- Describe **dimensionality reduction** (reducing #features)
 - Feature selection
 - Heuristic search
 - Feature extraction
 - Principal component analysis (PCA)
 - Singular value decomposition (SVD)

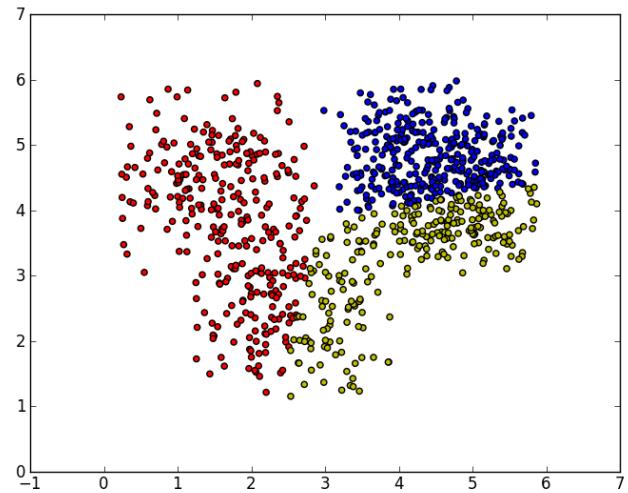
Histograms

- Divide data into buckets and store average (sum) for each bucket



Clustering

- Partition data set into clusters based on similarity, and store cluster representation (e.g., **centroid and diameter**) only
- Can be very effective if data is clustered but not if data is “smeared”
- There are many choices of clustering definitions and clustering algorithms



Sampling

- Sampling: obtaining a small sample s to represent the whole data set N
- Key principle: Choose a **representative** subset of the data
 - Simple random sampling may have very poor performance in the presence of skew

Simple random sampling:

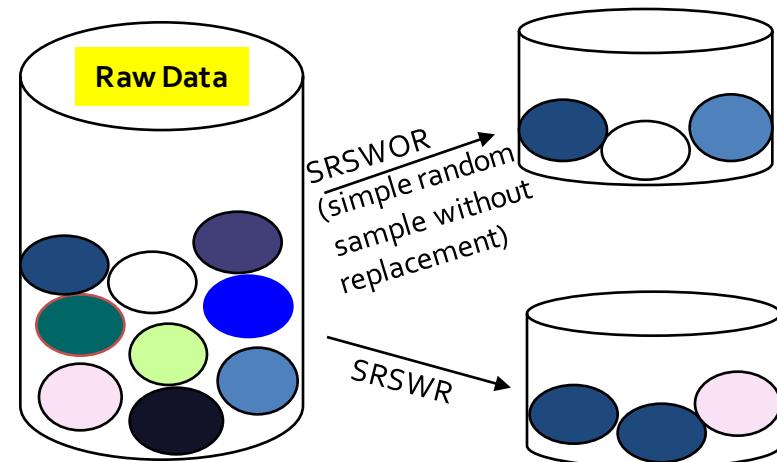
Equal probability of selecting any particular item

Sampling without replacement:

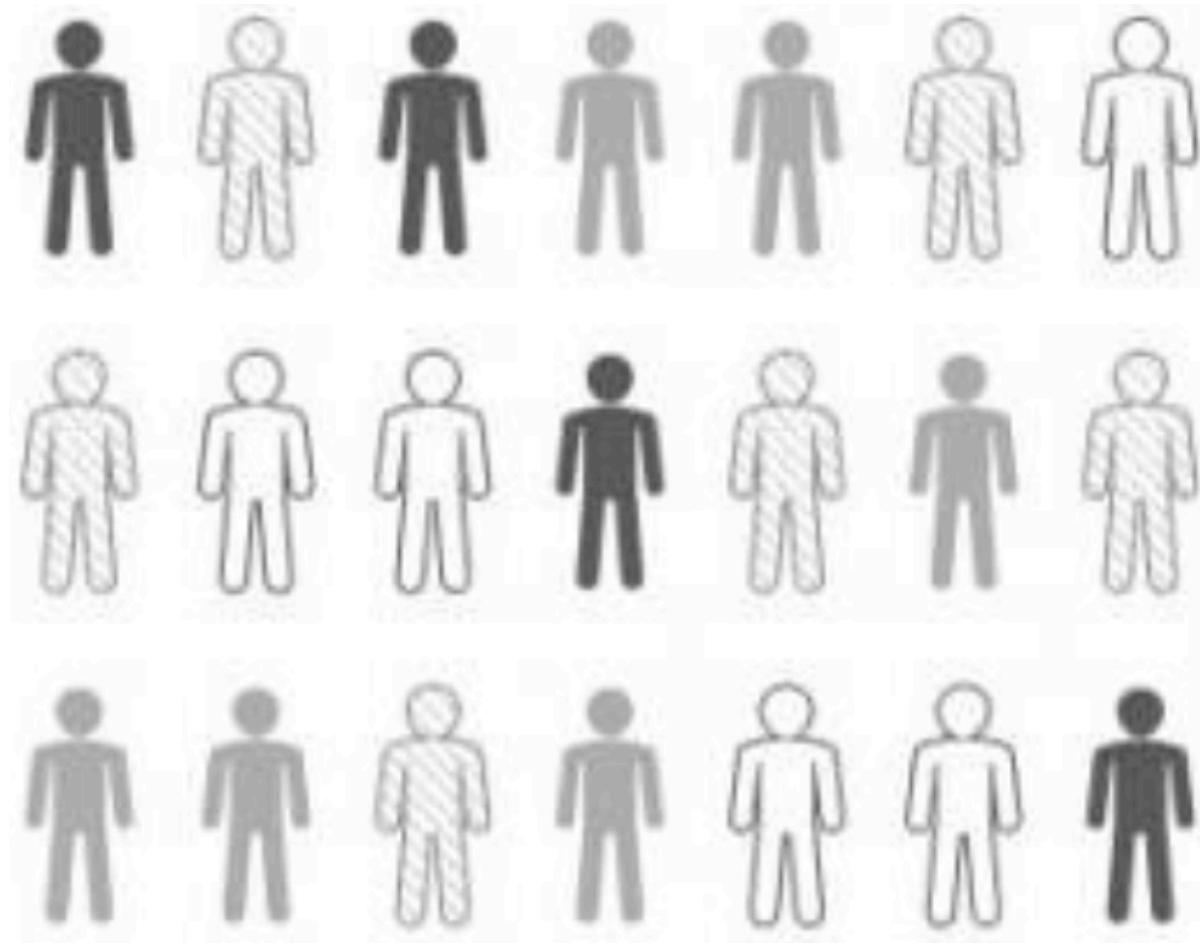
Once an object is selected, it is removed from the population

Sampling with replacement:

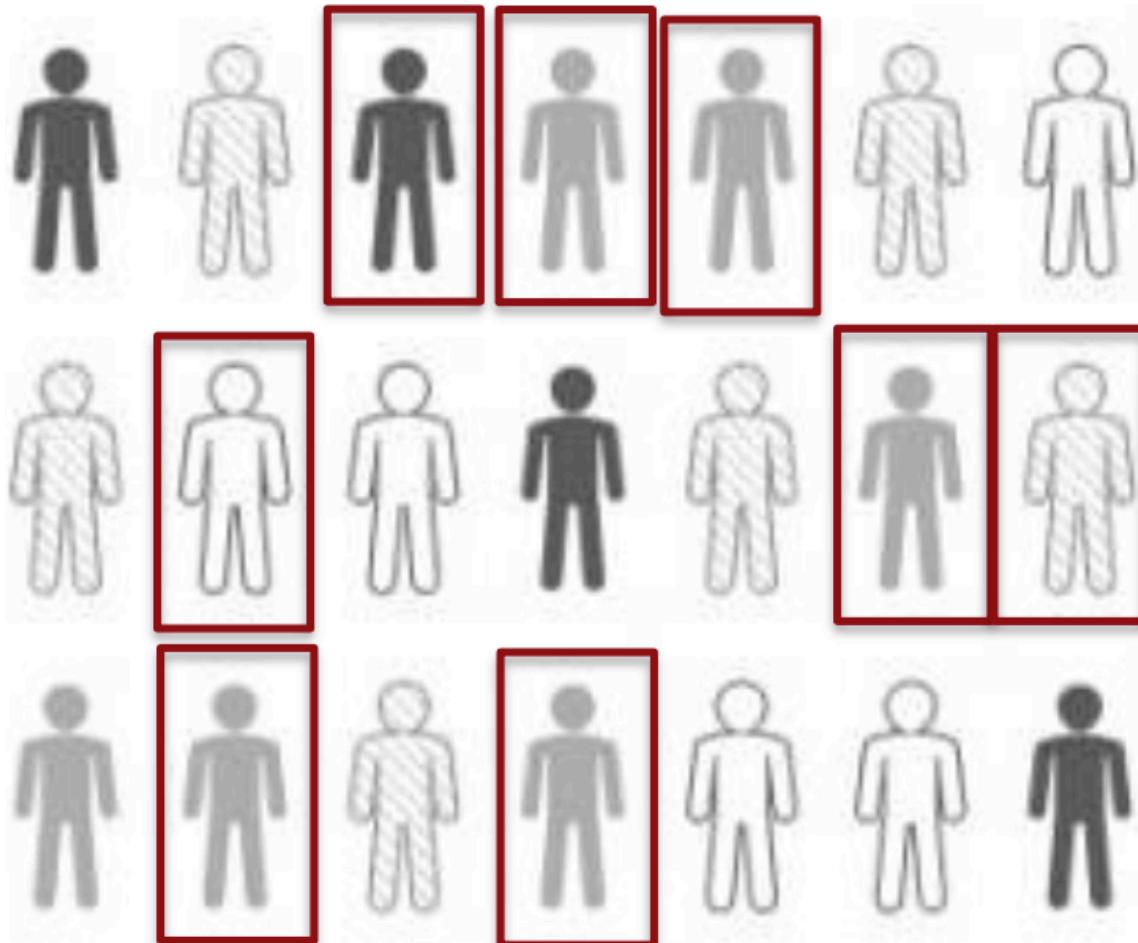
A selected object is not removed from the population



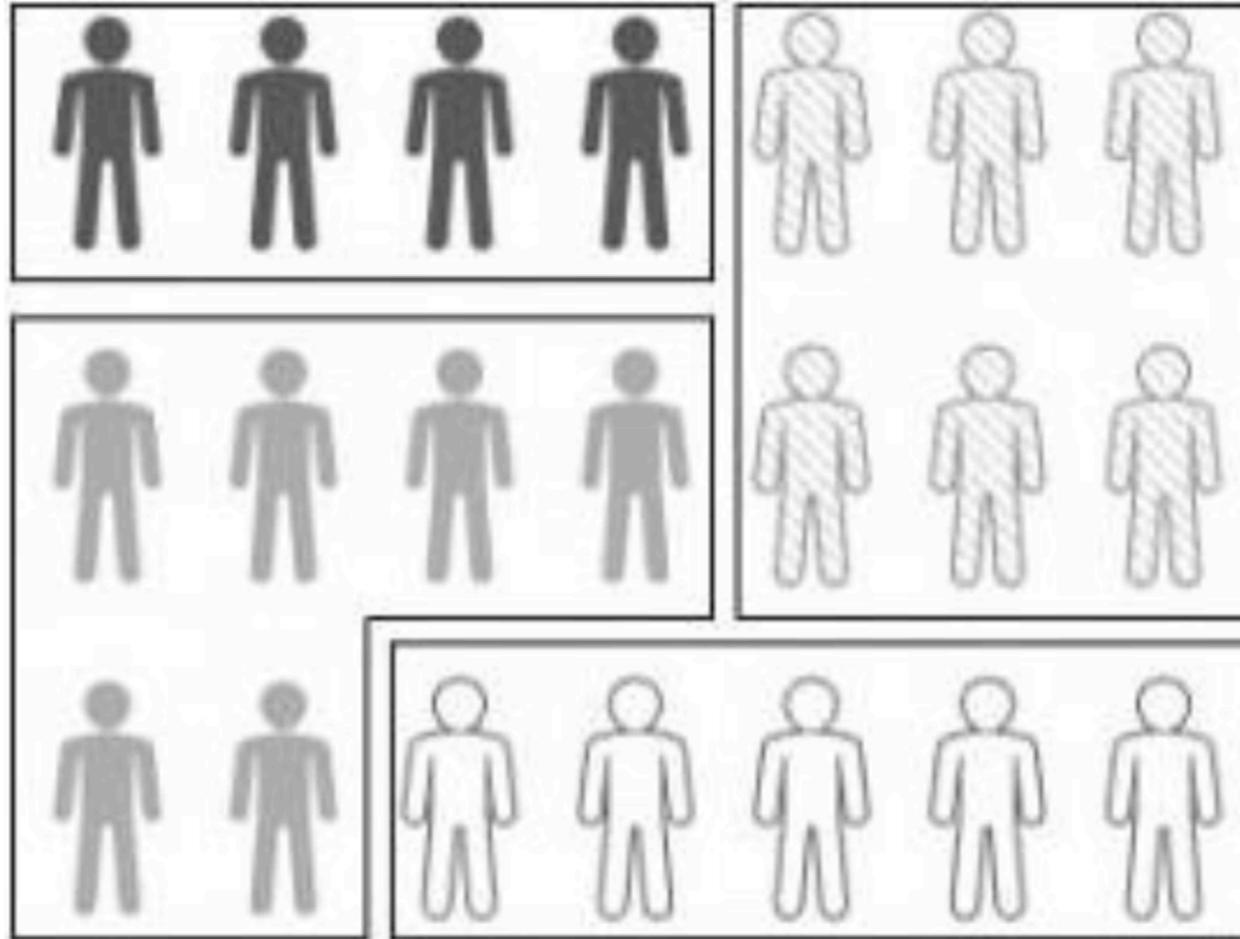
What if the data is *imbalanced*?



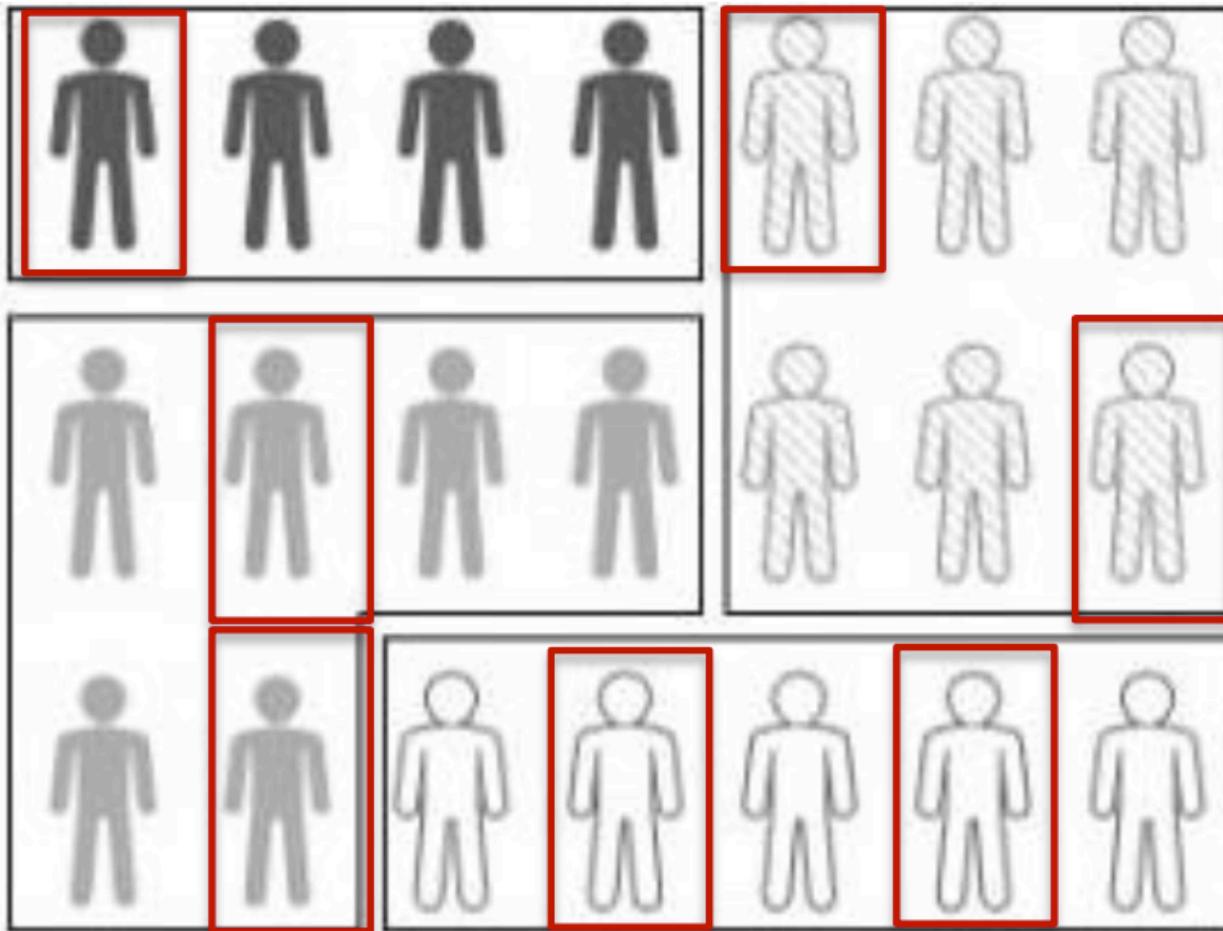
What if the data is *imbalanced*?



What if the data is *imbalanced*?

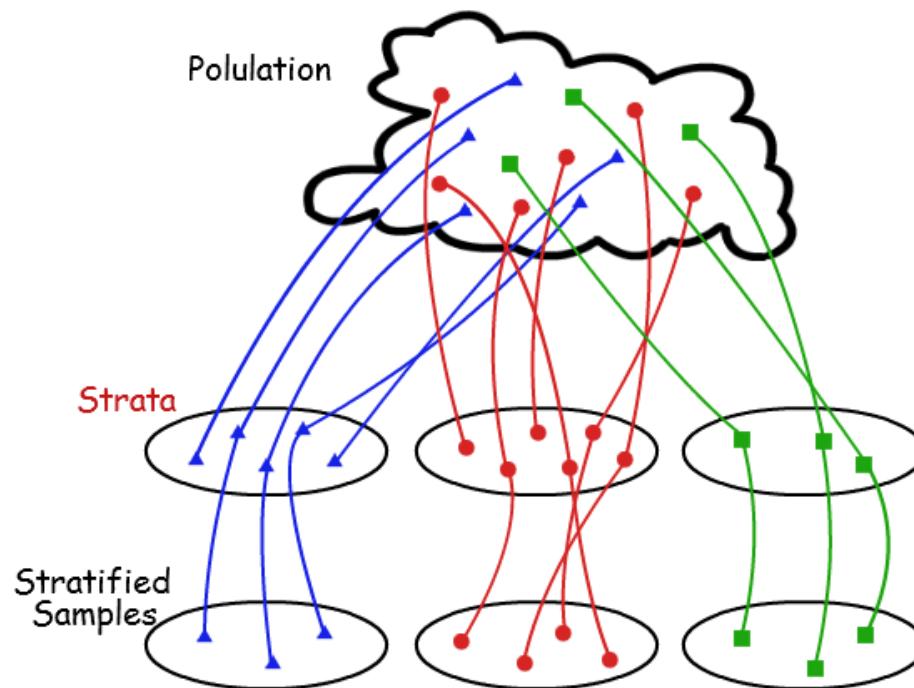


What if the data is *imbalanced*?



Stratified Sampling

- **Stratified sampling**
 - Partition (or cluster) the data set, and draw samples from each partition (proportionally, i.e., approximately the same percentage of the data)



Today: Data Reduction

- Describe **numerosity reduction** (reducing #instances)
 - Parametric methods: Fit some model and estimate model parameters
 - Regression: Describe linear/non-linear regression models
 - Nonparametric methods
 - Histograms
 - Clustering
 - Sampling: Describe stratified sampling
- Describe **dimensionality reduction** (reducing #features)
 - Feature selection
 - Heuristic search
 - Feature extraction
 - Principal component analysis (PCA)
 - Singular Value Decomposition (SVD)

Dimensionality Reduction Methodologies

- **Feature selection:** Find a subset of the original variables (or features, attributes)
 - Heuristic search in attribute subset selection
 - ...
- **Feature extraction:** Transform the data in the high-dimensional space to a space of fewer dimensions
 - Principal component analysis
 - ...

Attribute Subset Selection

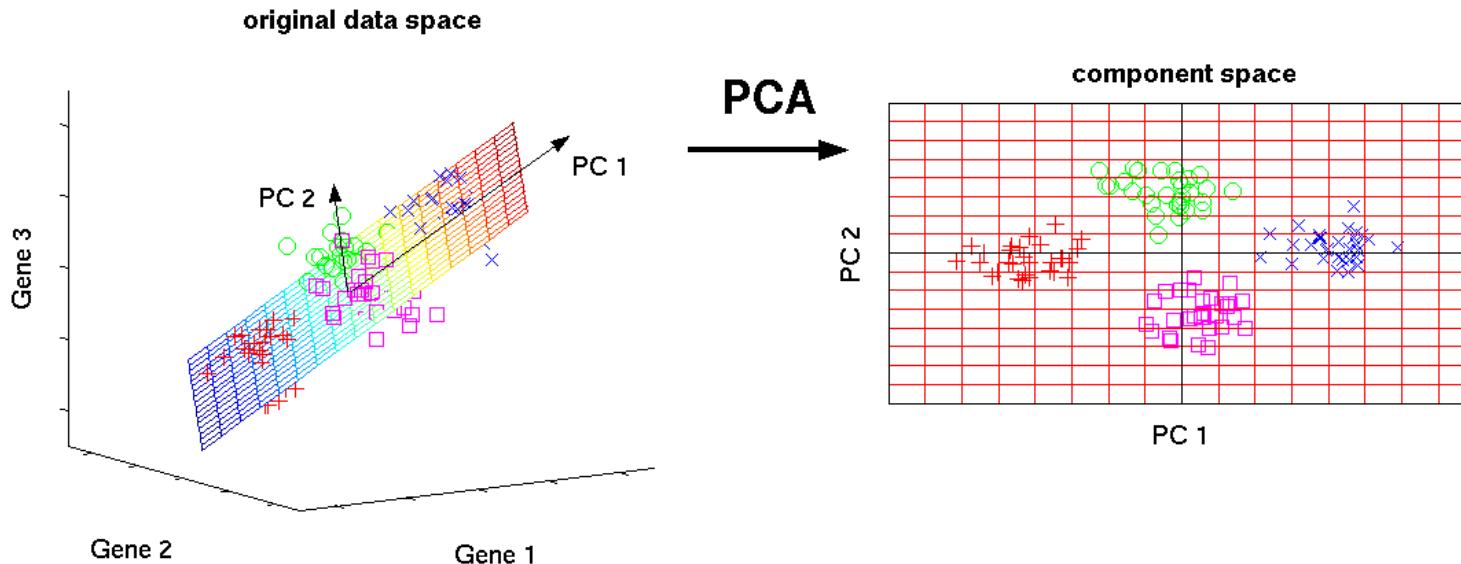
- Redundant attributes
 - Duplicate much or all of the information contained in one or more other attributes
 - Purchase price of a product and the amount of sales tax paid
- Irrelevant attributes
 - Contain no information that is useful for the data mining task at hand
 - A student's ID is often irrelevant to the task of predicting his/her GPA

Heuristic Search in Attribute Selection

- There are 2^d possible attribute combinations of d attributes
- Typical heuristic attribute selection methods:
 - **Method 1:** Best single attribute under the attribute independence assumption: choose by significance tests
 - **Method 2:** Best step-wise feature selection:
 - The best single-attribute is picked first
 - Then next best attribute condition to the first, ...
 - **Method 3:** Step-wise attribute elimination:
 - Repeatedly eliminate the worst attribute
 - **Method 4:** Best combined attribute selection and elimination

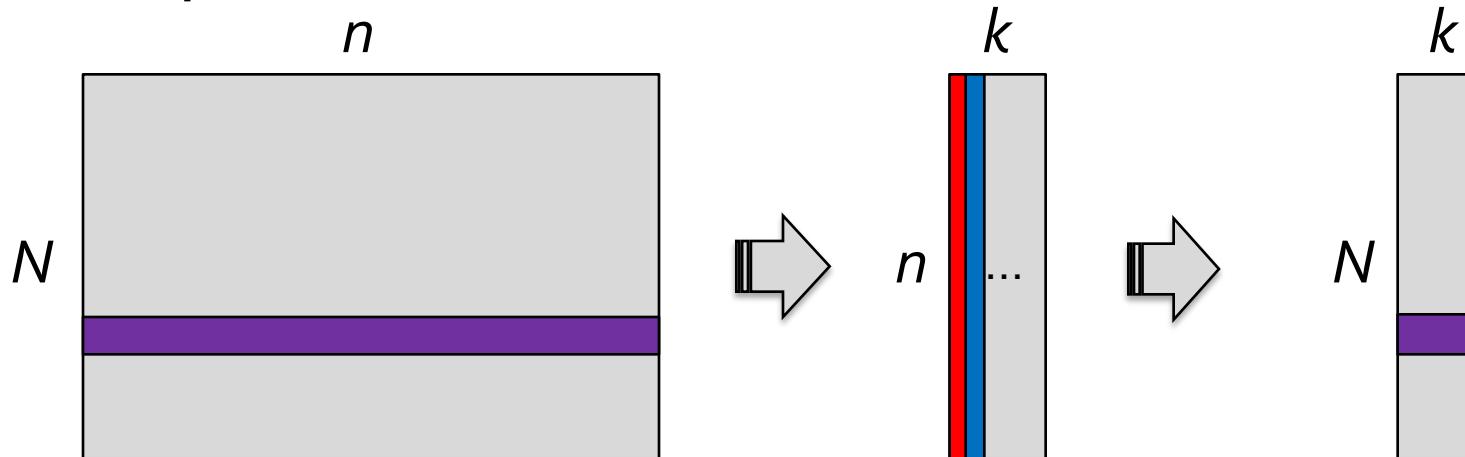
Principal Component Analysis (PCA)

- PCA: A statistical procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called ***principal components***
- The original data are projected onto a **much smaller space**, resulting in dimensionality reduction (e.g., n=3 to k=2)



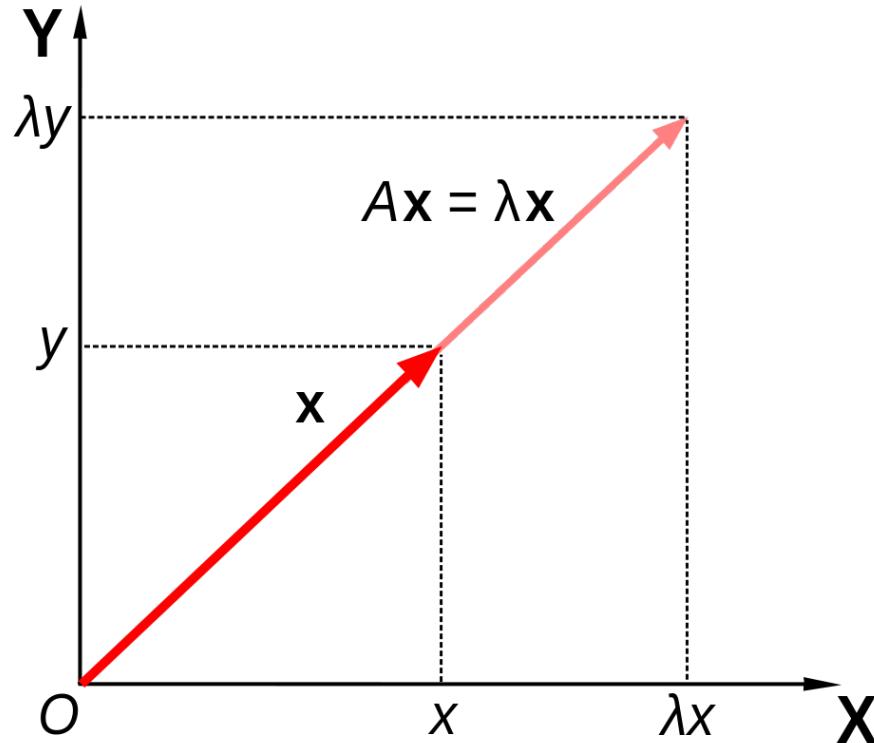
PCA (cont.)

- Given N data vectors from n -dimensions, find $k \leq n$ orthogonal vectors (principal components) best used to represent data
 - Normalize input data: Each attribute falls within the same range
 - Compute k orthonormal (unit) vectors, i.e., principal components
normalized eigenvector
- Each input data (vector) is a linear combination of the k principal components

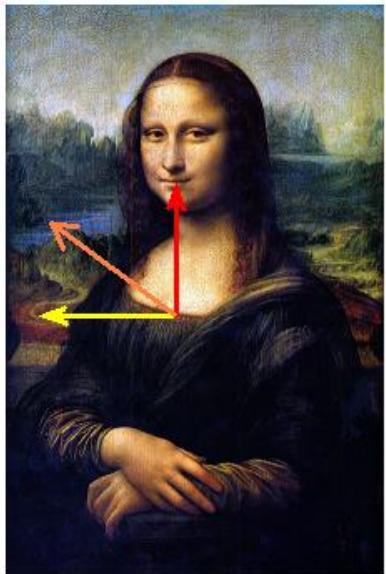


Eigenvectors (cont.)

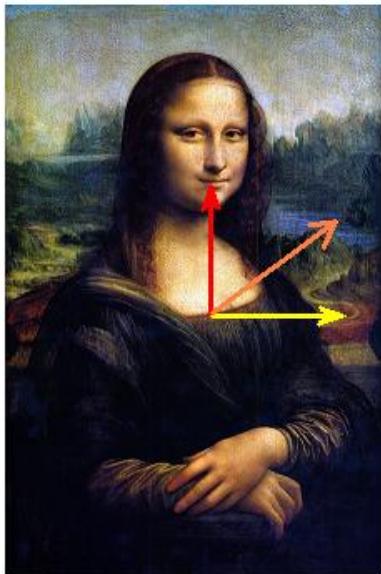
- For a square matrix \mathbf{A} ($n \times n$), find the eigenvector \mathbf{x} ($n \times 1$).
 - \mathbf{A} represents the linear transformation (from n to n)
- Matrix \mathbf{A} acts by stretching the vector \mathbf{x} , not changing its direction, so \mathbf{x} is an eigenvector of \mathbf{A} .



Eigenvectors (cont.)



A_1



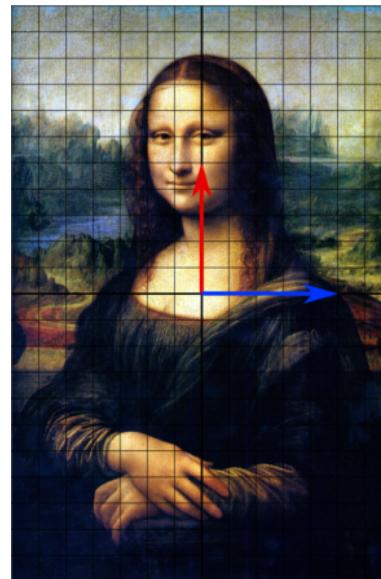
Which vectors are eigenvectors?

- Red
- Orange
- Yellow

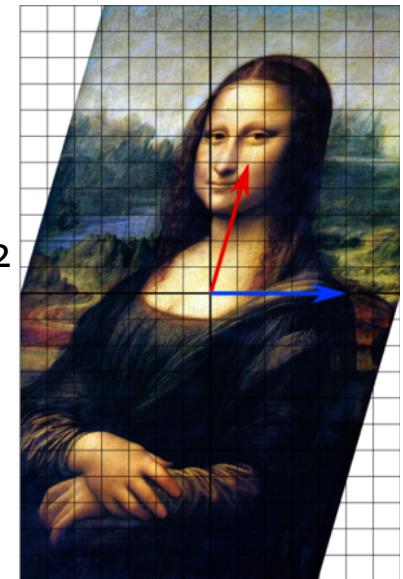
What are the eigenvalues?

Which vectors are eigenvectors?

- Red
- Blue

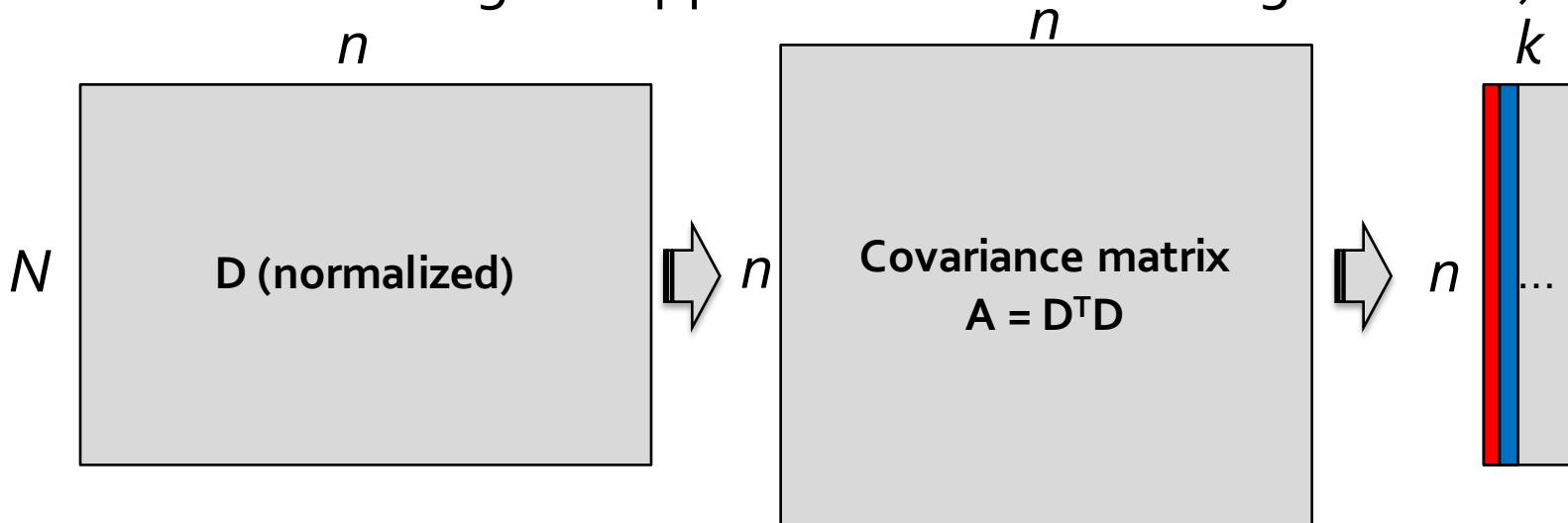


A_2



PCA and Eigenvectors

- For ***Square Matrix***: Data matrix to Covariance matrix
- The principal components are sorted in order of **decreasing “significance” or strength**
- **From n to k :** Since the components are sorted, the size of the data can be reduced by eliminating the weak components (i.e., using the strongest principal components, to reconstruct a good approximation of the original data)



PCA and Eigenvectors (cont.)

- Method: Find the **eigenvectors of covariance (square) matrix**, and these eigenvectors define the new space

$$\begin{aligned}\mathbf{Ax} = \lambda\mathbf{x} &\Leftrightarrow \mathbf{Ax} - \lambda\mathbf{x} = \mathbf{0} \\ &\Leftrightarrow \mathbf{Ax} - \lambda\mathbf{I}\mathbf{x} = \mathbf{0} \\ &\Leftrightarrow (\mathbf{A} - \lambda\mathbf{I})\mathbf{x} = \mathbf{0}.\end{aligned}$$

The equation $\mathbf{Ax} = \lambda\mathbf{x}$ has nonzero solutions for the vector x if and only if the matrix $\mathbf{A} - \lambda\mathbf{I}$ has zero determinant.

Example: Find the eigenvalues of the matrix $\mathbf{A} = \begin{bmatrix} 2 & 2 \\ 5 & -1 \end{bmatrix}$.

Ex. Eigenvalues

Example: Find the eigenvalues of the matrix $\mathbf{A} = \begin{bmatrix} 2 & 2 \\ 5 & -1 \end{bmatrix}$.

The eigenvalues are those λ for which $\det(\mathbf{A} - \lambda\mathbf{I}) = 0$. Now

$$\begin{aligned}\det(\mathbf{A} - \lambda\mathbf{I}) &= \det\left(\begin{bmatrix} 2 & 2 \\ 5 & -1 \end{bmatrix} - \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right) \\ &= \det\left(\begin{bmatrix} 2 & 2 \\ 5 & -1 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix}\right) \\ &= \begin{vmatrix} 2 - \lambda & 2 \\ 5 & -1 - \lambda \end{vmatrix} \\ &= (2 - \lambda)(-1 - \lambda) - 10 \\ &= \lambda^2 - \lambda - 12.\end{aligned}$$

The eigenvalues of \mathbf{A} are the solutions of the quadratic equation $\lambda^2 - \lambda - 12 = 0$, namely $\lambda_1 = -3$ and $\lambda_2 = 4$.

Ex. Eigenvectors

First, we work with $\lambda = -3$. The equation $\mathbf{Ax} = \lambda\mathbf{x}$ becomes $\boxed{\mathbf{Ax} = -3\mathbf{x}}$. Writing

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

and using the matrix \mathbf{A} from above, we have

$$\mathbf{Ax} = \boxed{\begin{bmatrix} 2 & 2 \\ 5 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}} = \begin{bmatrix} 2x_1 + 2x_2 \\ 5x_1 - x_2 \end{bmatrix},$$

while

$$-3\mathbf{x} = \begin{bmatrix} -3x_1 \\ -3x_2 \end{bmatrix}.$$

Setting these equal, we get

$$\boxed{\begin{bmatrix} 2x_1 + 2x_2 \\ 5x_1 - x_2 \end{bmatrix}} = \begin{bmatrix} -3x_1 \\ -3x_2 \end{bmatrix} \Rightarrow 2x_1 + 2x_2 = -3x_1 \quad \text{and} \quad 5x_1 - x_2 = -3x_2$$

$$\Rightarrow 5x_1 = -2x_2$$

$$\Rightarrow \boxed{x_1 = -\frac{2}{5}x_2.}$$

$$\boxed{\mathbf{u}_1 = \begin{bmatrix} 2 \\ -5 \end{bmatrix}}$$

Ex. Eigenvectors (cont.)

Similarly, we can find eigenvectors associated with the eigenvalue $\lambda = 4$ by solving

$$\boxed{\mathbf{Ax} = 4\mathbf{x}}$$

$$\begin{bmatrix} 2x_1 + 2x_2 \\ 5x_1 - x_2 \end{bmatrix} = \begin{bmatrix} 4x_1 \\ 4x_2 \end{bmatrix} \Rightarrow 2x_1 + 2x_2 = 4x_1 \quad \text{and} \quad 5x_1 - x_2 = 4x_2 \\ \Rightarrow x_1 = x_2.$$

Hence the set of eigenvectors associated with $\lambda = 4$ is spanned by

$$\boxed{\mathbf{u}_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}}.$$

Ex. Eigenvalues (cont.)

Example: Find the eigenvalues and associated eigenvectors of the matrix

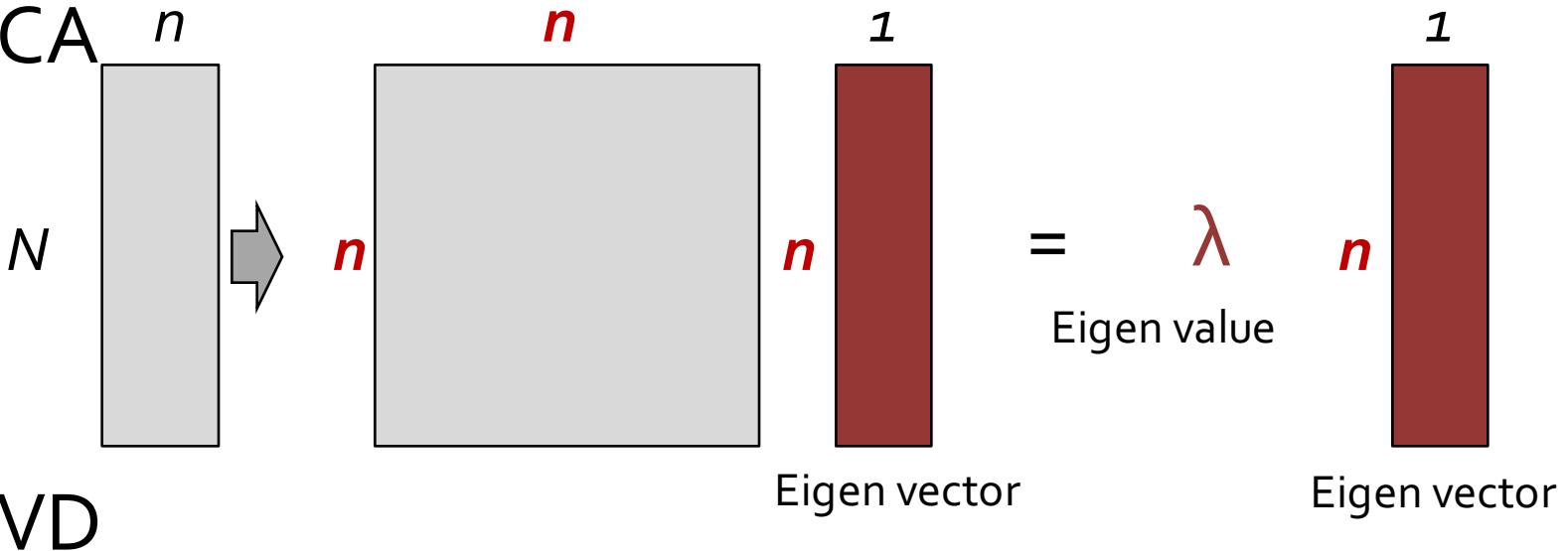
$$\mathbf{A} = \begin{bmatrix} 7 & 0 & -3 \\ -9 & -2 & 3 \\ 18 & 0 & -8 \end{bmatrix}.$$

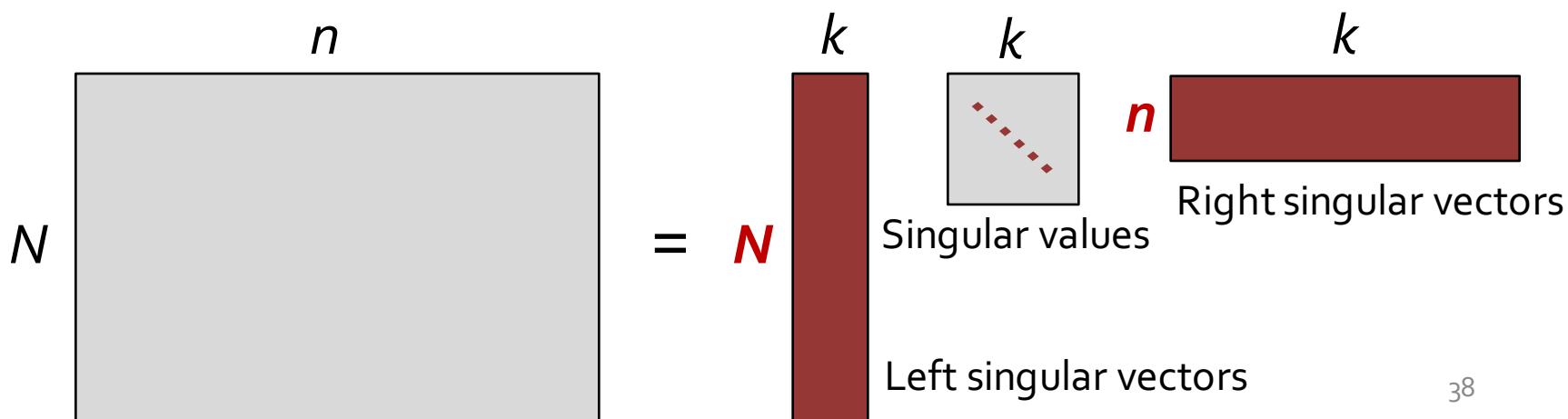
First we compute $\det(\mathbf{A} - \lambda\mathbf{I})$ via a cofactor expansion along the second column:

$$\begin{aligned} \left| \begin{array}{ccc} 7-\lambda & 0 & -3 \\ -9 & -2-\lambda & 3 \\ 18 & 0 & -8-\lambda \end{array} \right| &= (-2-\lambda)(-1)^4 \left| \begin{array}{cc} 7-\lambda & -3 \\ 18 & -8-\lambda \end{array} \right| \\ &= -(2+\lambda)[(7-\lambda)(-8-\lambda) + 54] \\ &= -(\lambda+2)(\lambda^2 + \lambda - 2) \\ &= -(\lambda+2)^2(\lambda-1). \end{aligned}$$

Thus \mathbf{A} has two distinct eigenvalues, $\lambda_1 = -2$ and $\lambda_3 = 1$. (Note that we might say $\lambda_2 = -2$, since, as a root, -2 has multiplicity two. This is why we labelled the eigenvalue 1 as λ_3 .)

Singular Value Decomposition (SVD)

- PCA
- 
- A diagram illustrating PCA decomposition. On the left, a gray vertical rectangle labeled N by n is shown. An arrow points from it to a gray square labeled n by n . This square is followed by an equals sign. To the right of the equals sign is a red vertical rectangle labeled 1 by n , with the word "Eigen value" written below it. Another equals sign follows, and to its right is another red vertical rectangle labeled 1 by n .
- SVD


$$\begin{matrix} & n \\ \begin{matrix} N \end{matrix} & \end{matrix} = \begin{matrix} & k \\ \begin{matrix} N \end{matrix} & \end{matrix} \begin{matrix} & k \\ \begin{matrix} \text{Singular values} \end{matrix} & \end{matrix} \begin{matrix} & k \\ \begin{matrix} \text{Right singular vectors} \end{matrix} & \end{matrix}$$

A diagram illustrating SVD decomposition. On the left, a gray vertical rectangle labeled N by n is shown. An equals sign follows. To the right of the equals sign is a red vertical rectangle labeled k by n , with the label "Left singular vectors" below it. Next is a gray square labeled k by k with a dashed red diagonal line, representing the singular value matrix. To the right of the square is another red vertical rectangle labeled k by k , with the label "Right singular vectors" below it.

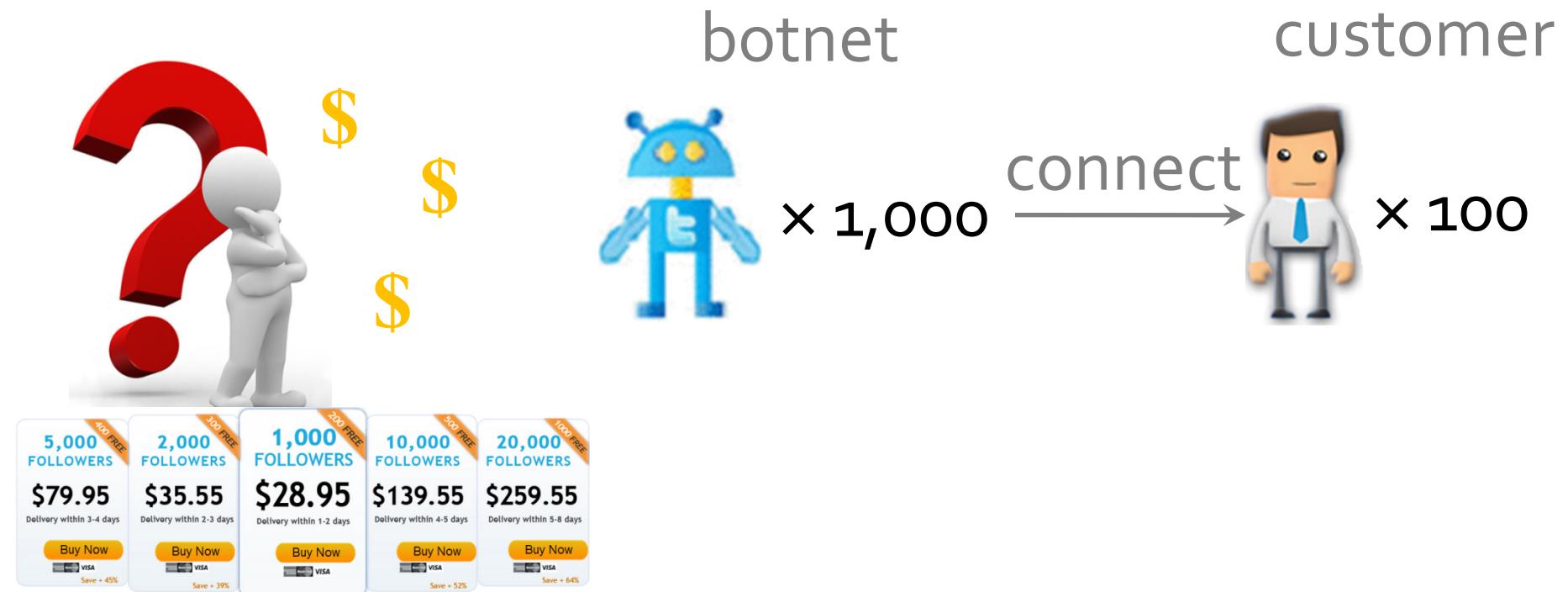
Extra: Use SVD to Find Strange Behaviors

- Sell followers: “Become a Twitter Rockstar”

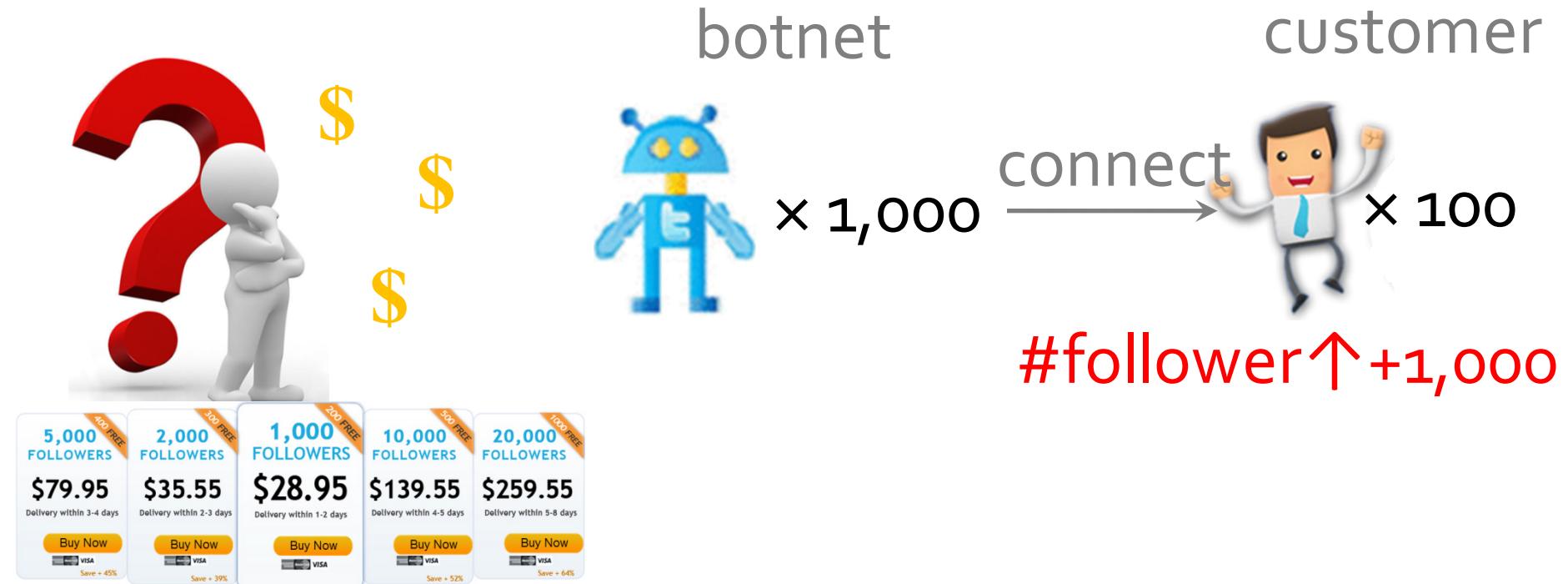


Jiang, M., Cui, P., Beutel, A., Faloutsos, C. and Yang, S., 2014, May. Inferring strange behavior from connectivity pattern in social networks. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining* (pp. 126-138). Springer, Cham.

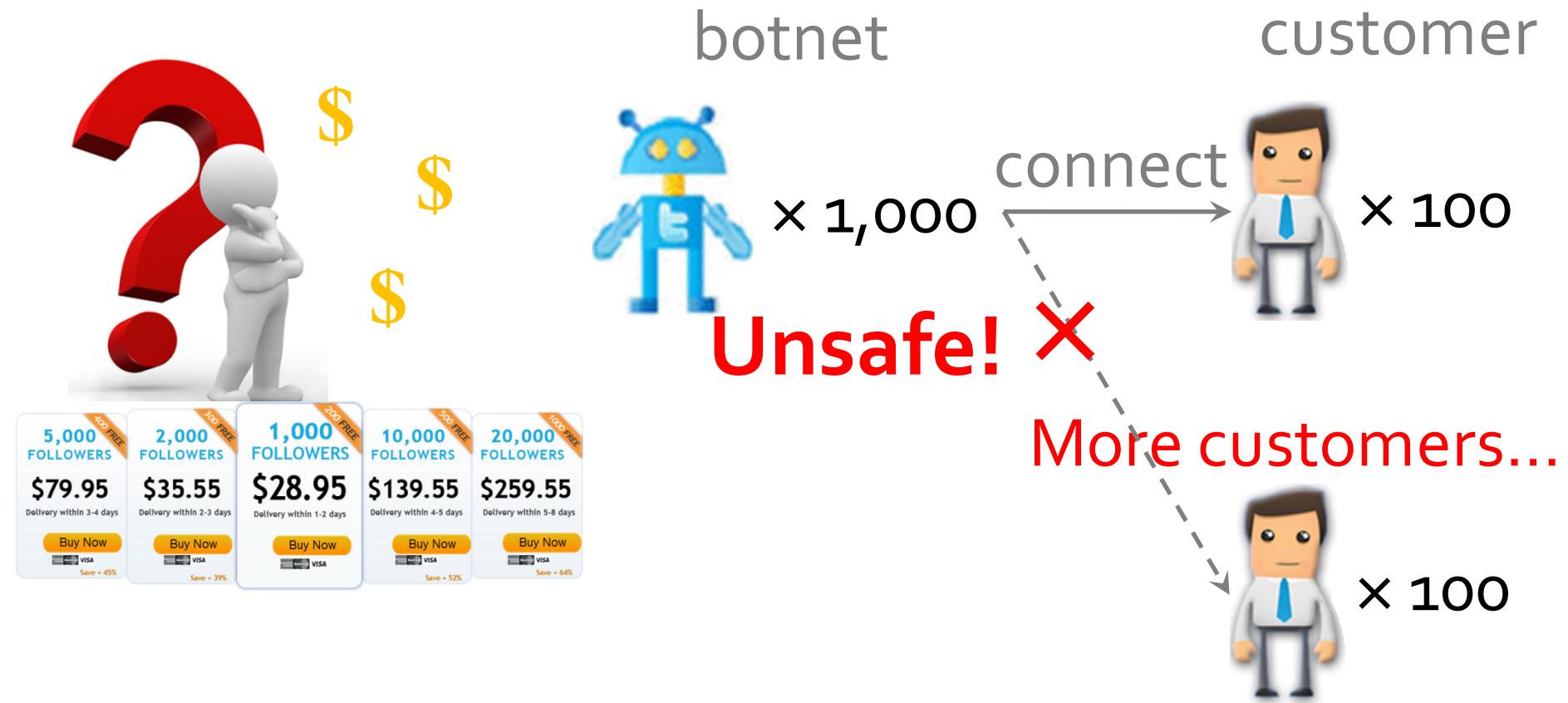
A Strange Behavior



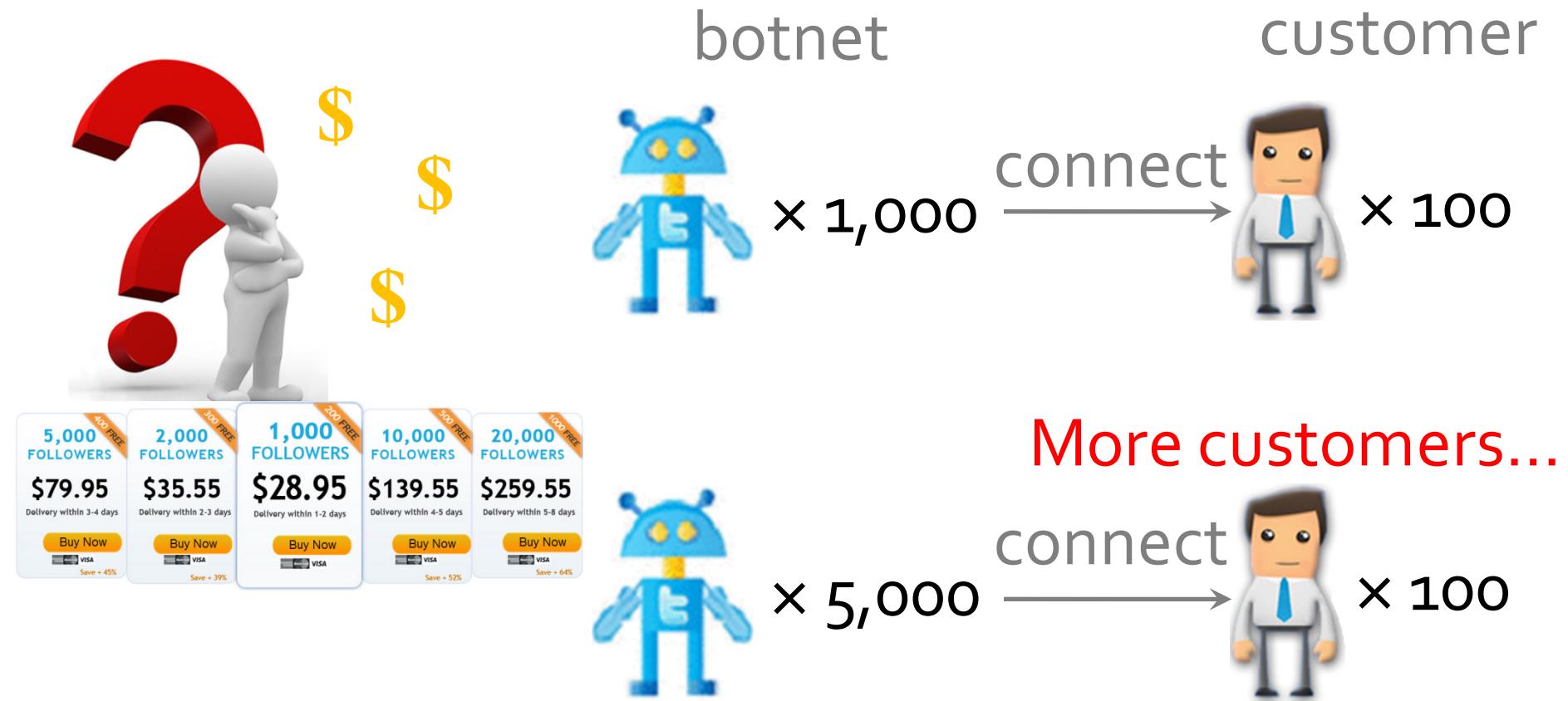
A Strange Behavior



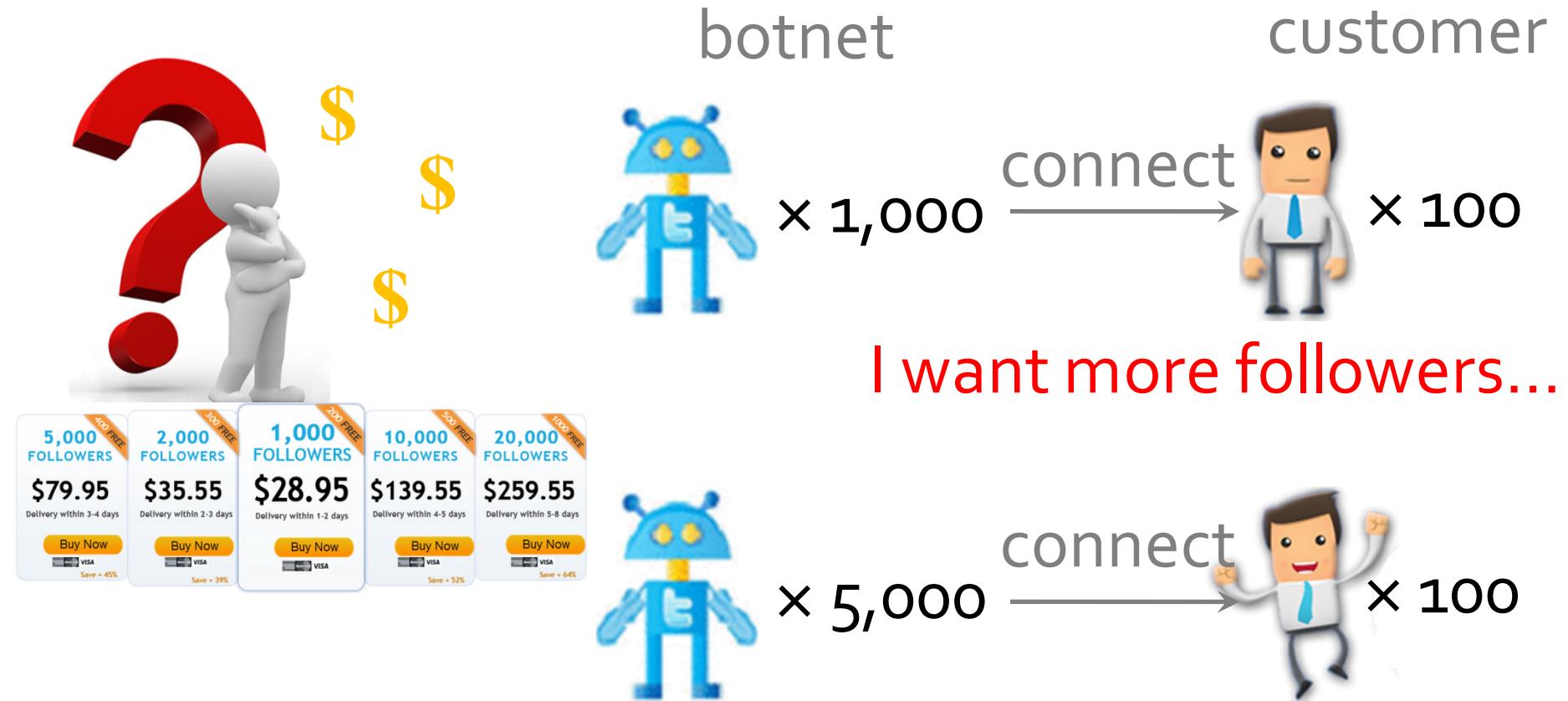
A Strange Behavior



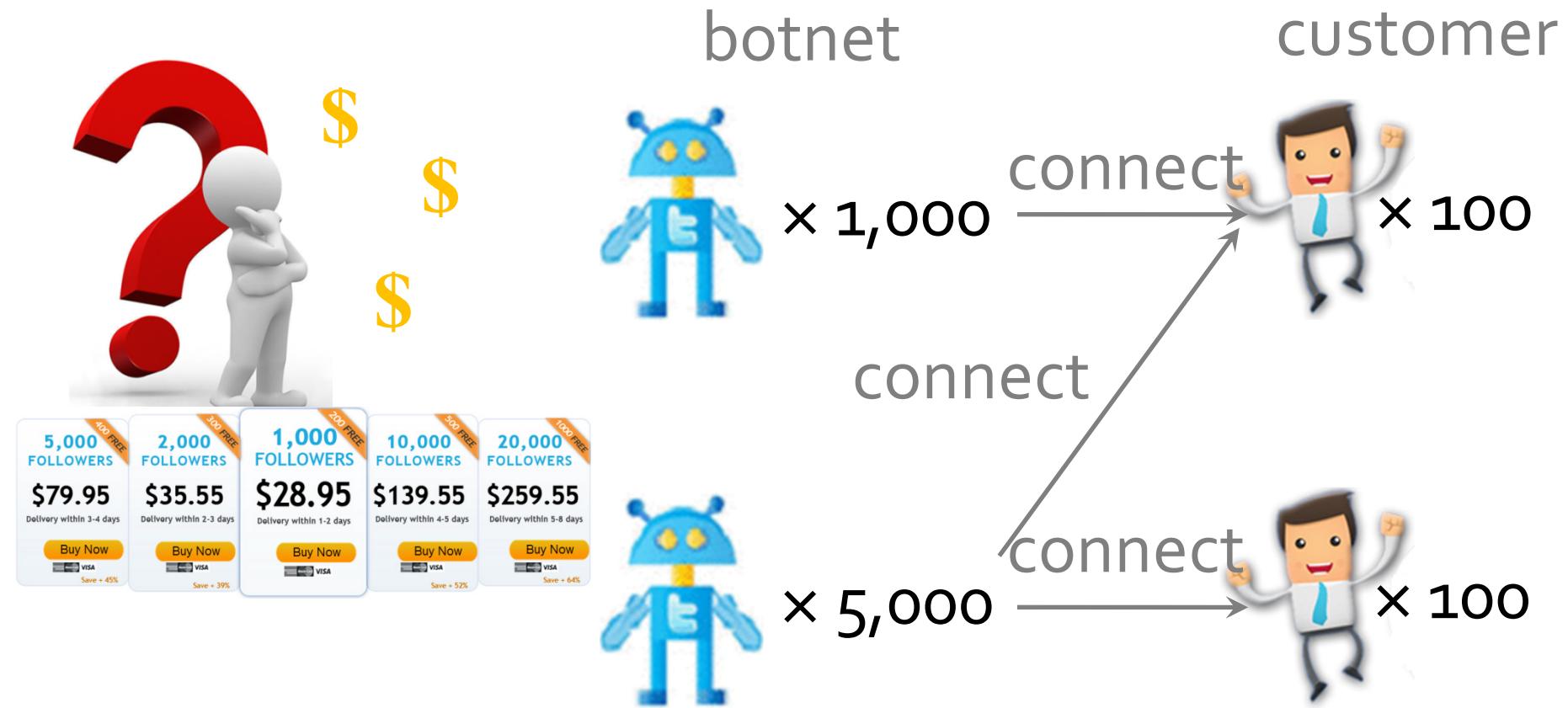
A Strange Behavior



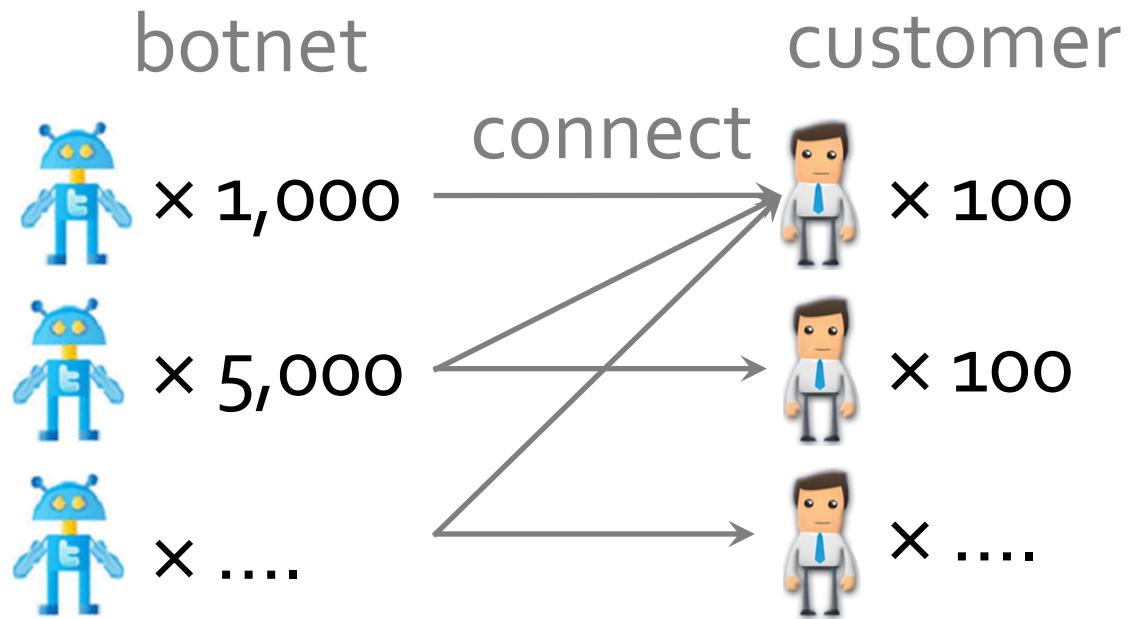
A Strange Behavior



A Strange Behavior



A Strange Behavior



More groups of customers
More groups of botnets
More companies....

A Strange Behavior



botnet



connect

customer



Detect dense bipartite cores!
How can we evade detection?
Some other activity!

A Strange Behavior



botnet



customer

connect

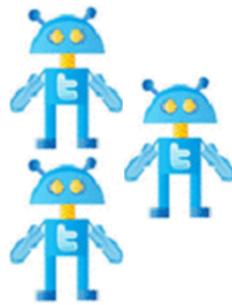


“Camouflage”:
may connect
to popular idols
to look normal

A Strange Behavior



botnet



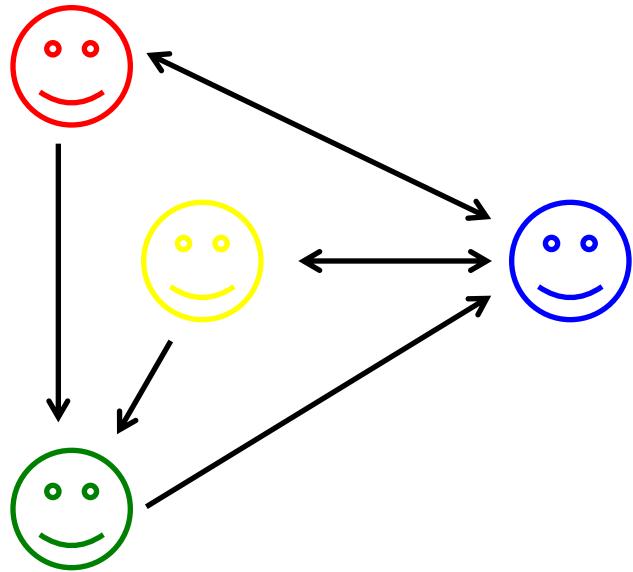
customer

connect

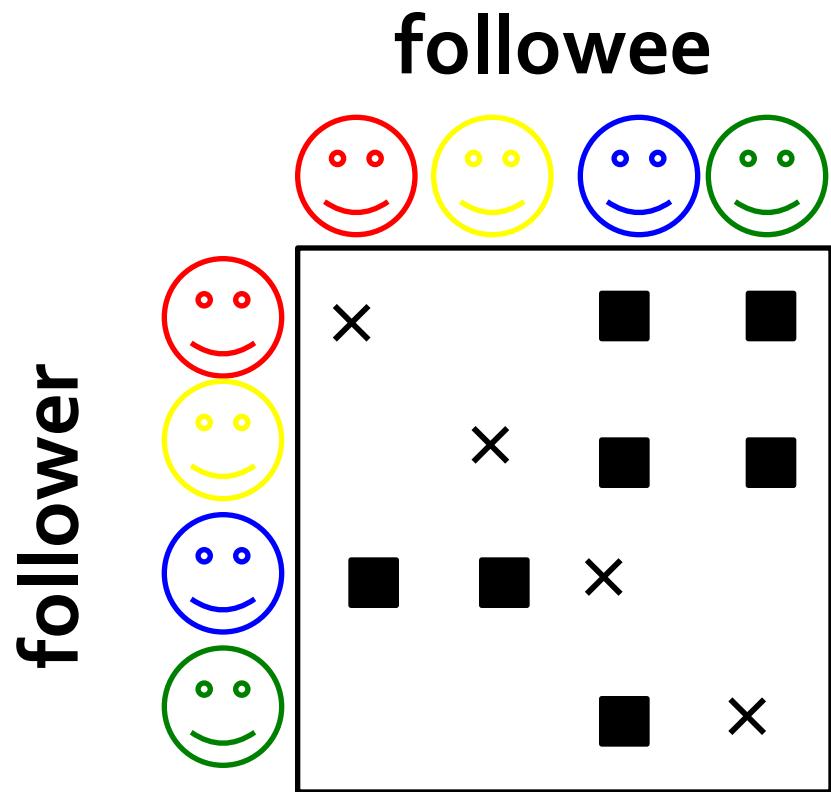


"Fame":
may have
a few honest
followers

Adjacency Matrix Reminder



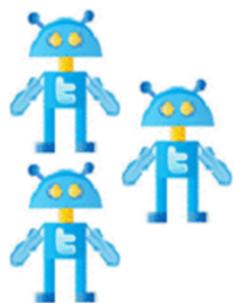
Graph Structure



Adjacency Matrix

Strange → Lockstep Behavior

botnet

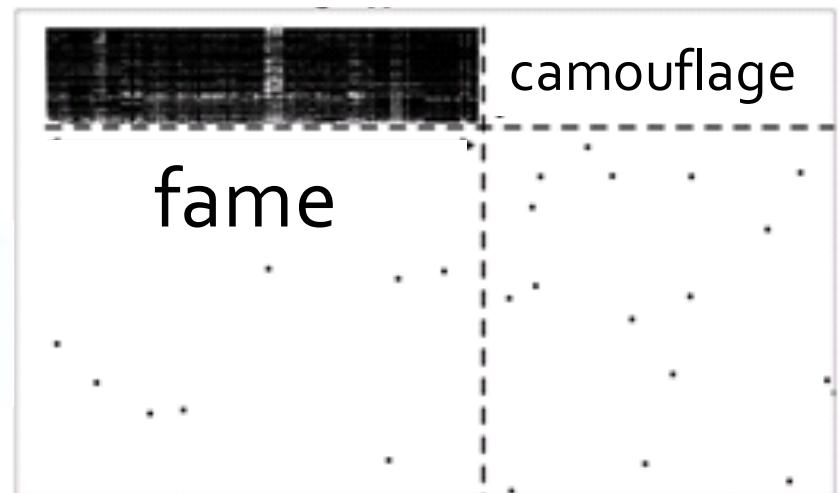


customer

connect

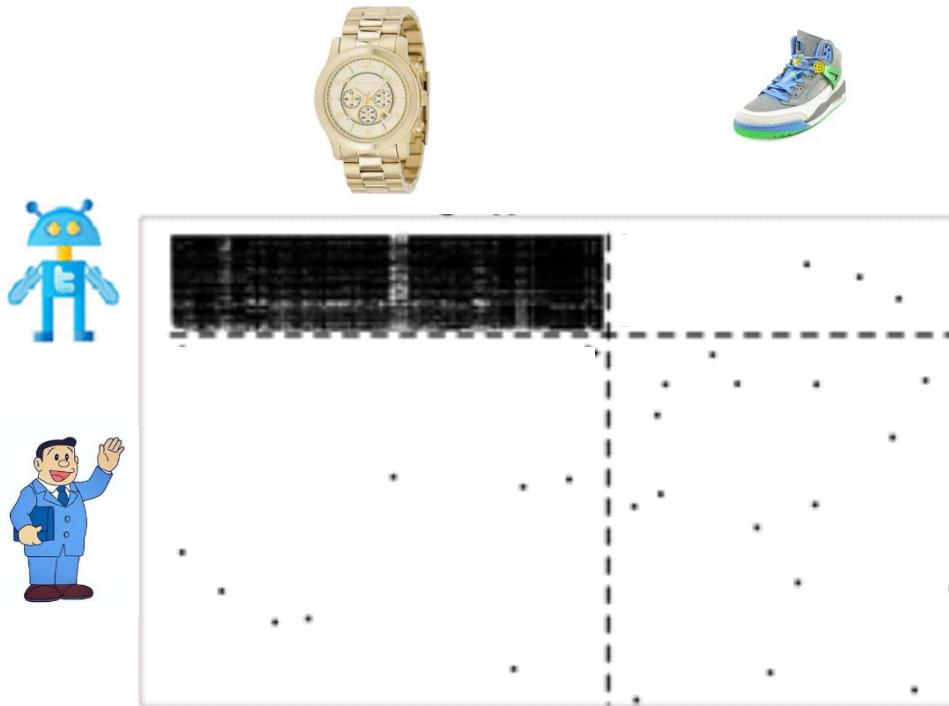


- Groups
- Acting together
- Little other activity



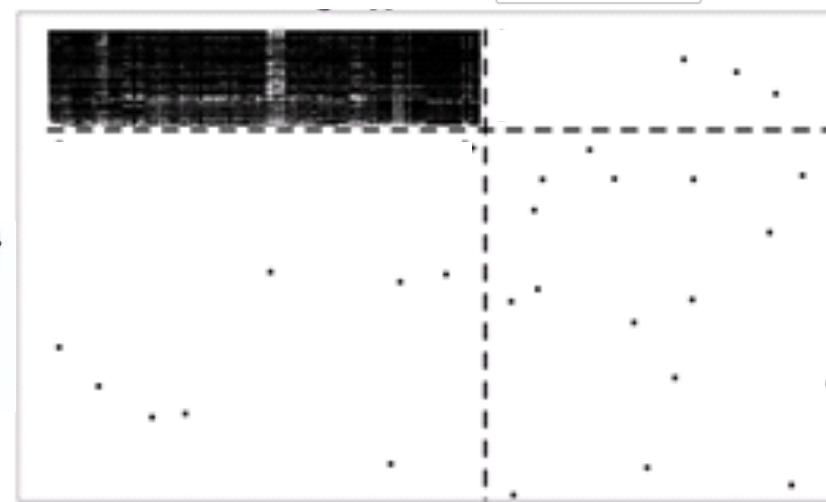
More Applications

- eBay reviews



More Applications

- Facebook “Likes”

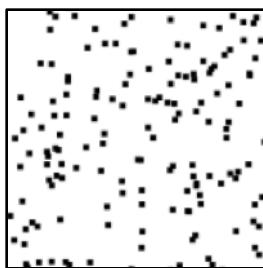


Barack Obama
40,591,963 likes • 1,280,959 talking about this

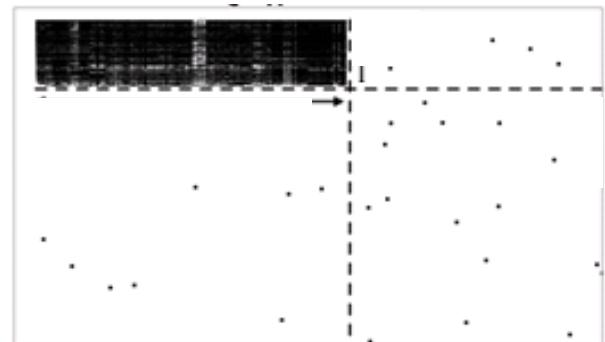


Problem Definition

- Given adjacency matrix

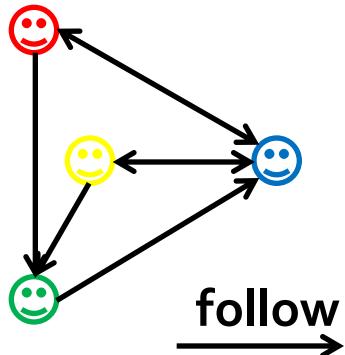


reordering
→

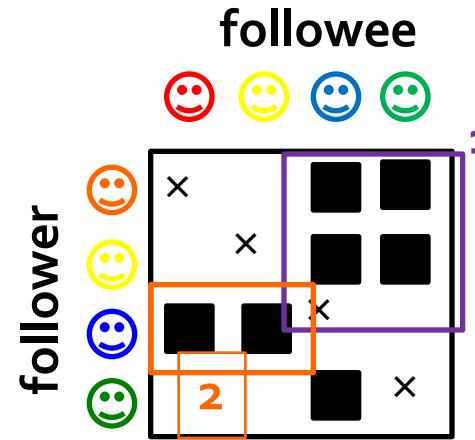


- Find Strange = “Lockstep” Behavior

SVD Reminder

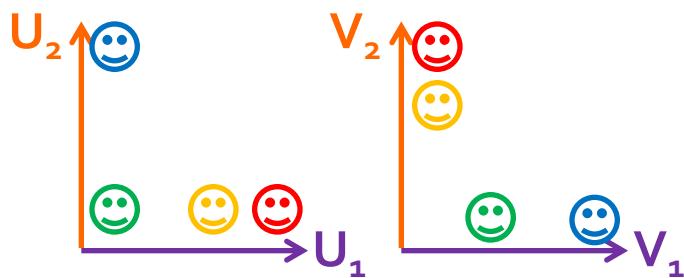


Graph Structure

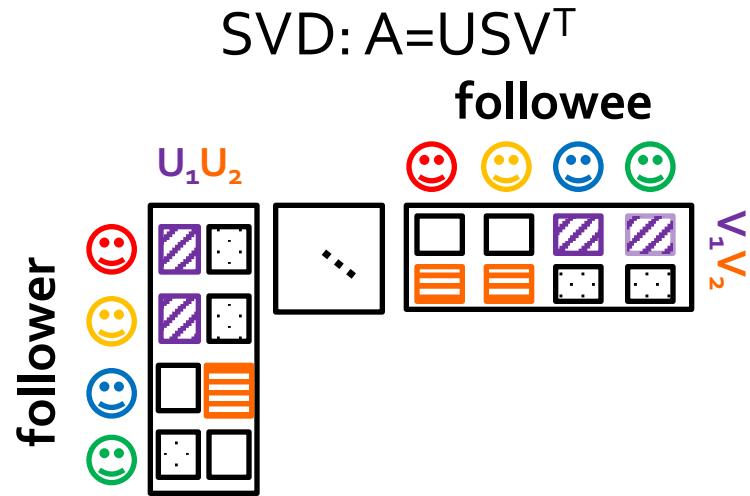


Adjacency Matrix

Pairs of singular vectors:

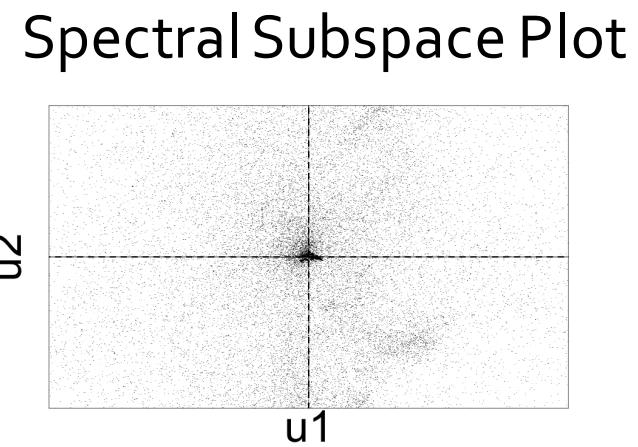
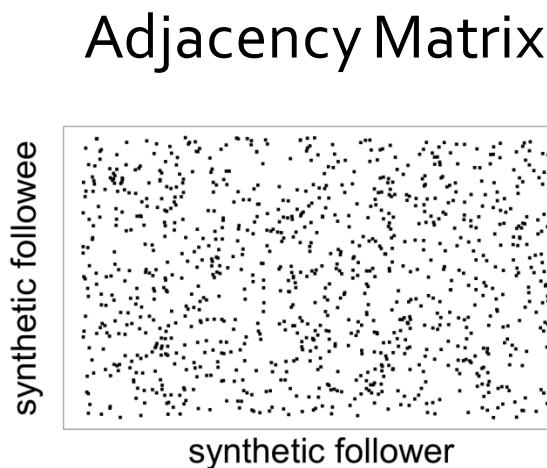


"Spectral Subspace Plot"



Lockstep and Spectral Subspace Plot

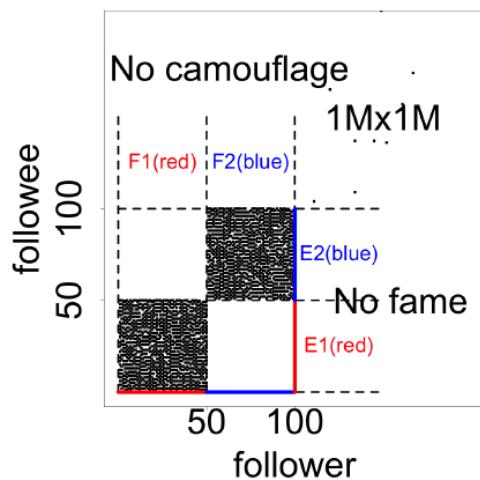
- Case #o: No lockstep behavior in random power law graph of 1M nodes, 3M edges
- Random \longleftrightarrow “Scatter”



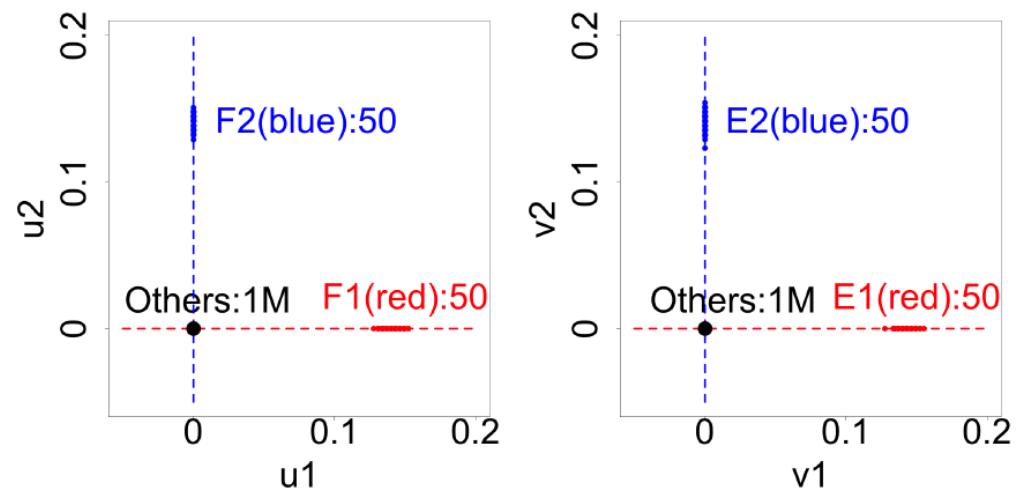
Lockstep and Spectral Subspace Plot

- Case #1: non-overlapping lockstep
- “Blocks” \longleftrightarrow “Rays”

Adjacency Matrix



Spectral Subspace Plot

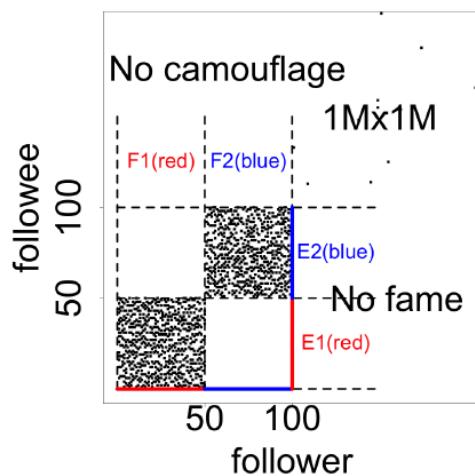


Rule 1 (short “rays”): two blocks, high density (90%), no “camouflage”, no “fame”

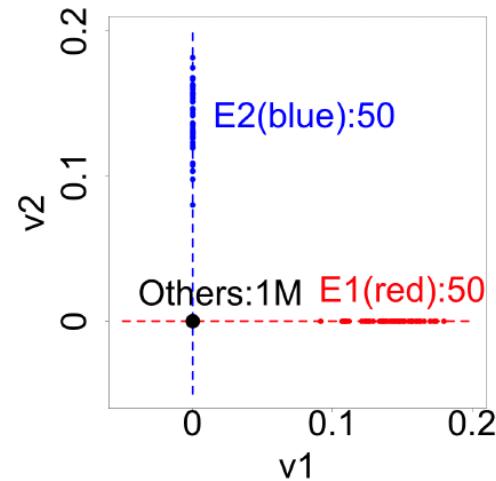
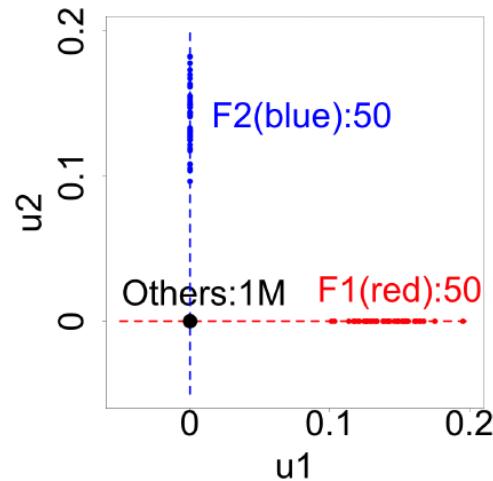
Lockstep and Spectral Subspace Plot

- Case #2: non-overlapping lockstep
- “Blocks; low density” \longleftrightarrow Elongation

Adjacency Matrix



Spectral Subspace Plot

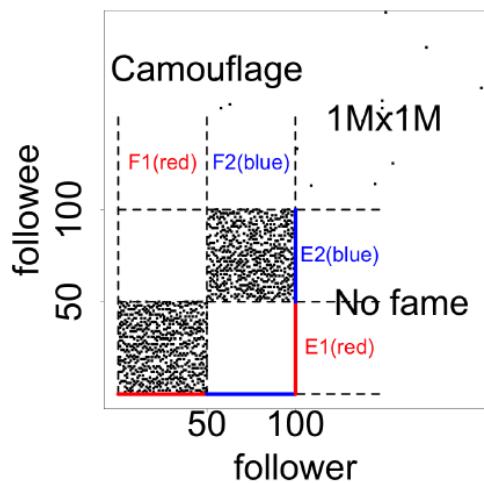


Rule 2 (long “rays”): two blocks, low density (50%), no “camouflage”, no “fame”

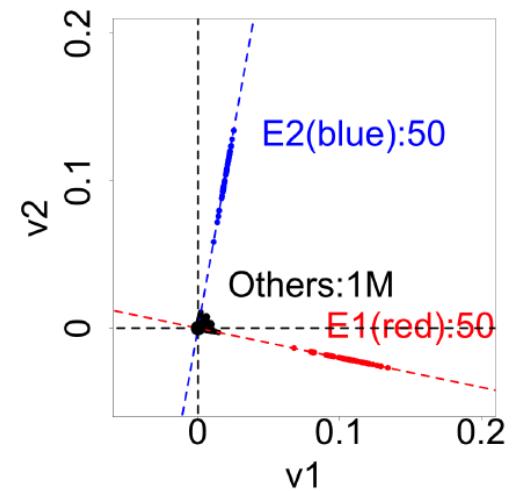
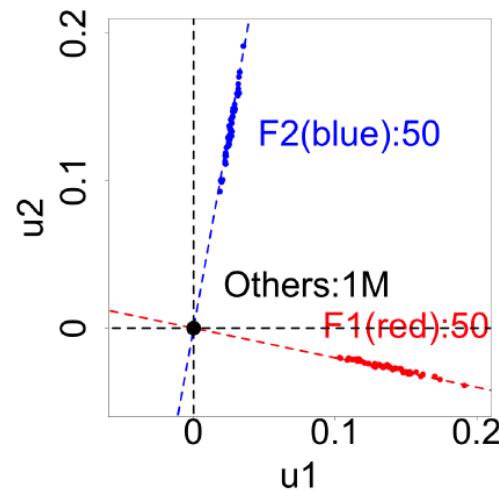
Lockstep and Spectral Subspace Plot

- Case #3: non-overlapping lockstep
- “Camouflage” (or “Fame”) \longleftrightarrow Tilting “Rays”

Adjacency Matrix



Spectral Subspace Plot

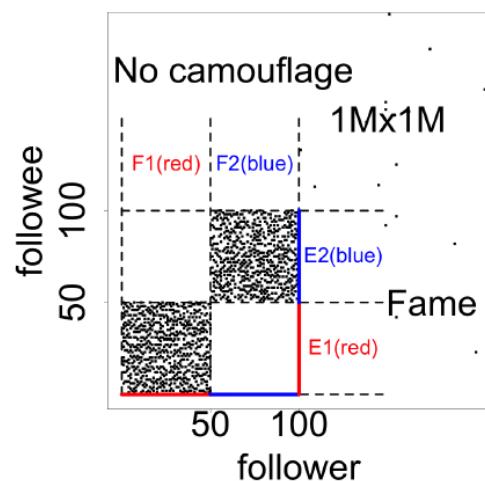


Rule 3 (tilting “rays”): two blocks, with “camouflage”, no “fame”

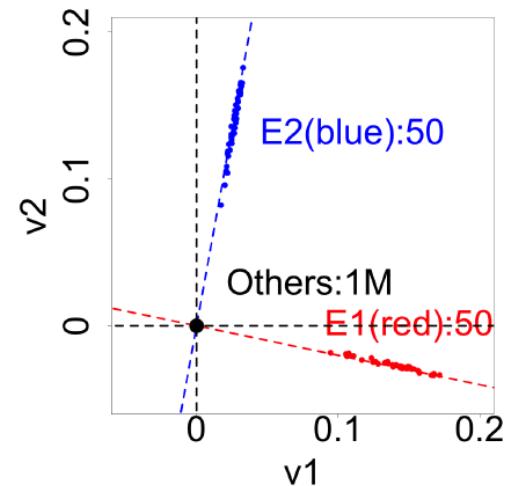
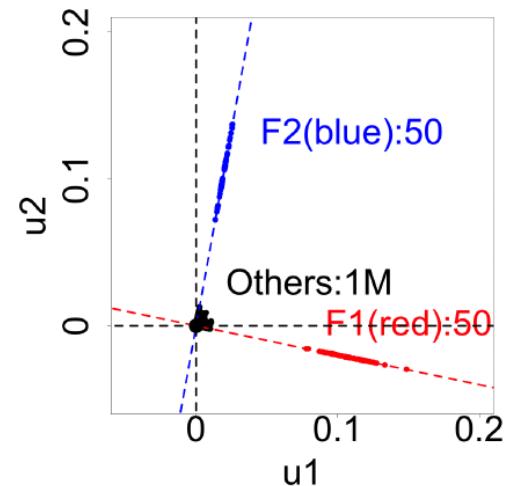
Lockstep and Spectral Subspace Plot

- Case #3: non-overlapping lockstep
 - “Camouflage” (or “Fame”) \longleftrightarrow Tilting “Rays”

Adjacency Matrix



Spectral Subspace Plot



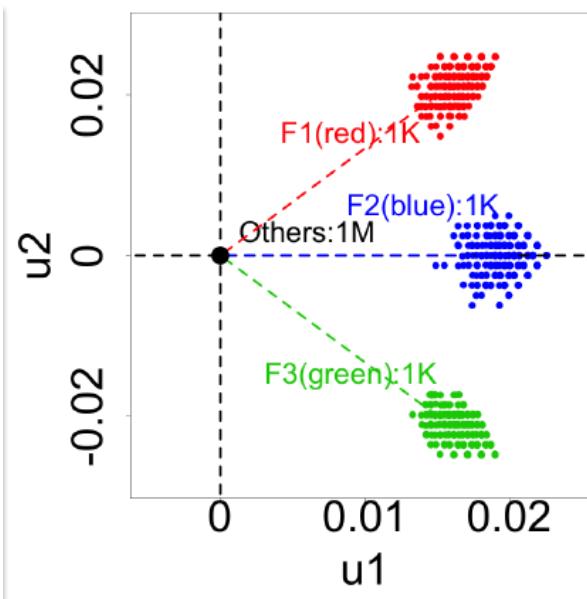
Rule 3 (tilting “rays”): two blocks, no “camouflage”, with “fame”

Lockstep and Spectral Subspace Plot

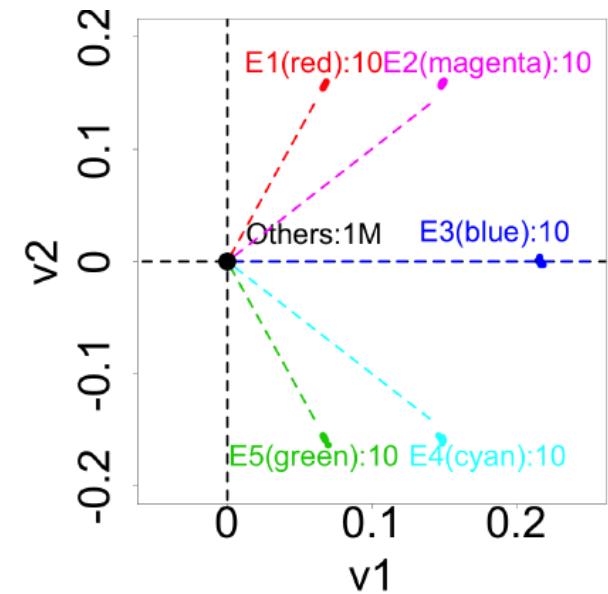
- Case #4: ? lockstep
- "?" ← → "Pearls"

Adjacency Matrix

?



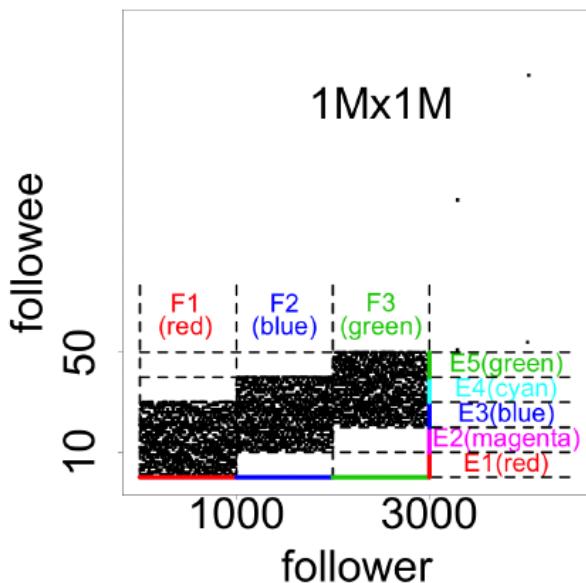
Spectral Subspace Plot



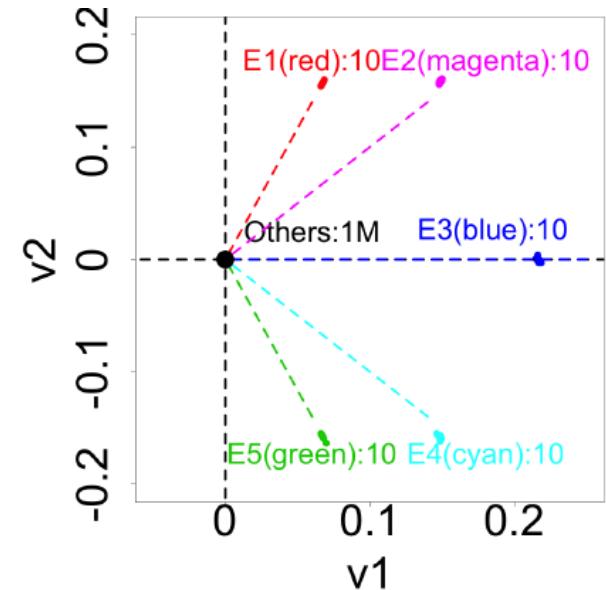
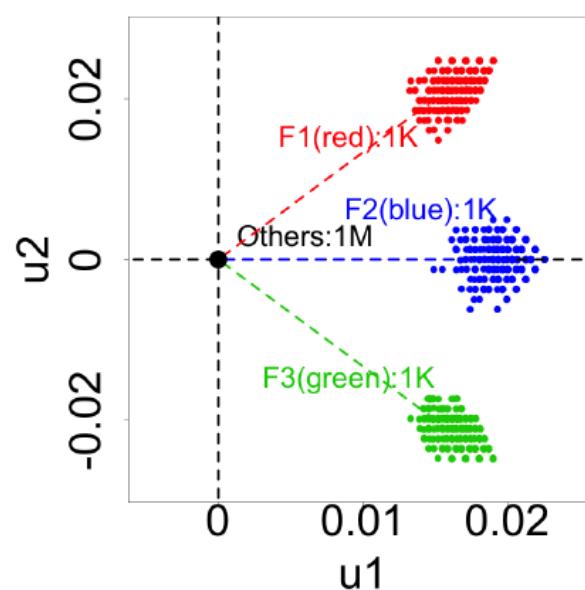
Lockstep and Spectral Subspace Plot

- Case #4: overlapping lockstep
- “Staircase” \longleftrightarrow “Pearls”

Adjacency Matrix



Spectral Subspace Plot



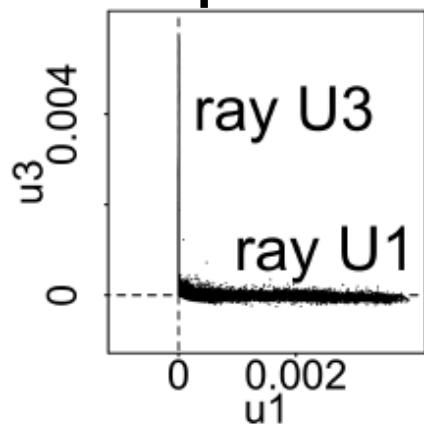
Rule 4 (“pearls”): a “staircase” of three partially overlapping blocks.

Algorithm

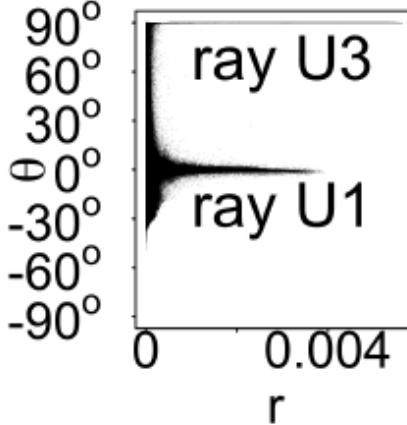
- Step 1: Seed selection
 - Spot “Rays” and “Pearls”
 - Catch seed followers
- Step 2: Belief Propagation
 - Blame followees with strange followers
 - Blame followers with strange followees

Automatically Spot “Rays” and “Pearls”

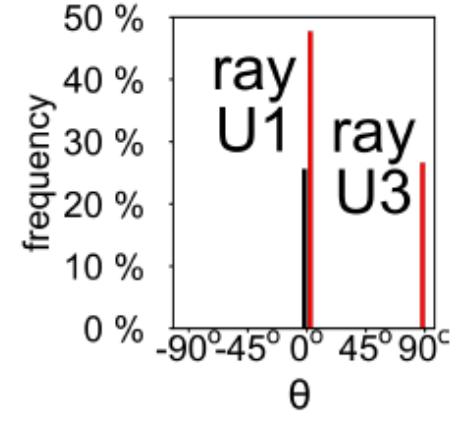
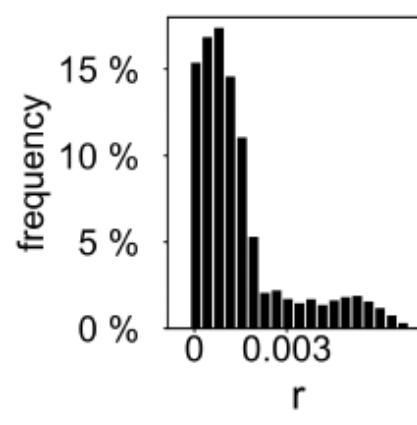
Spectral
Subspace Plot



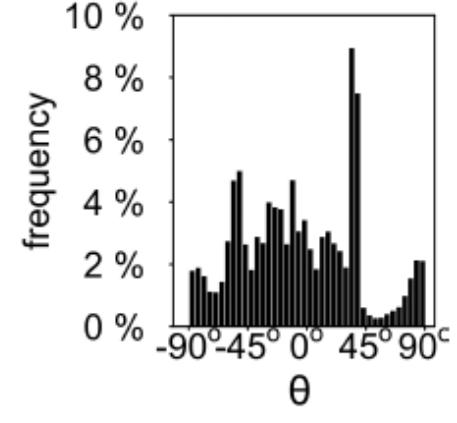
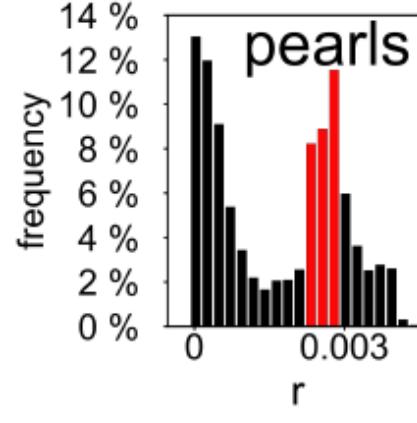
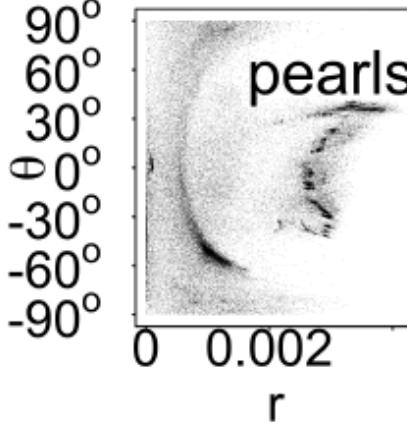
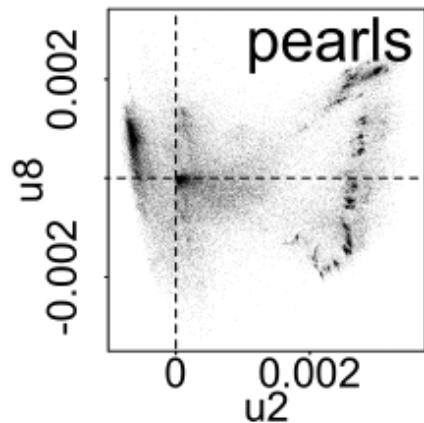
Polar Coordinate
Transform



Histograms



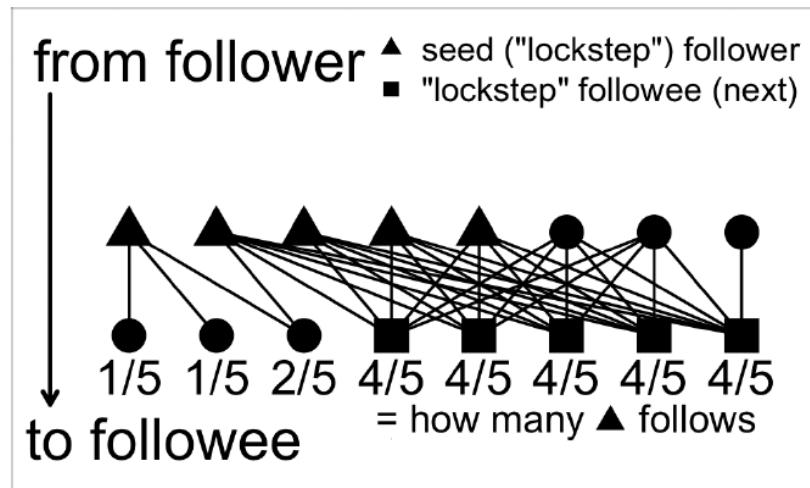
“rays” show two apparent spikes on θ frequency at 0° and 90°



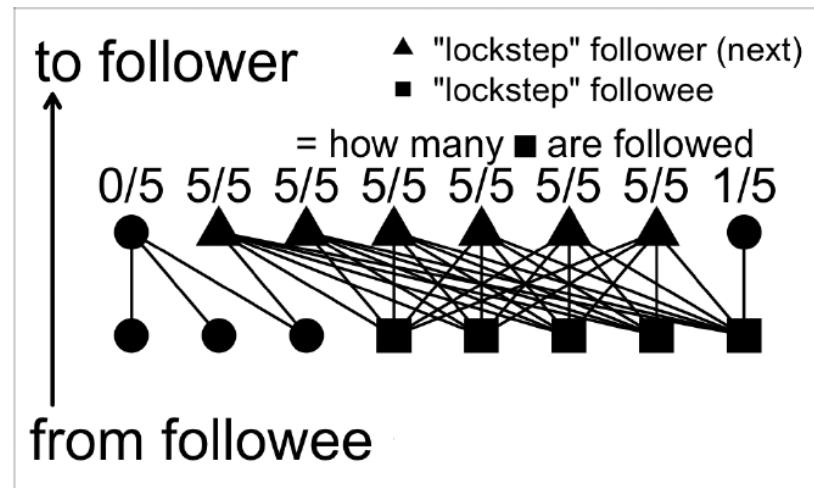
“pearls” show a spike on r frequency at a much-greater-than-zero value

BP-based Algorithm

- Blame followees with strange followers
- Blame followers with strange followees



(a) select “lockstep” followees:
from (seed) followers to followees



(b) select “lockstep” followers:
from followees to followers

Dataset

- Tencent Weibo
- 117 million nodes (users)
- 3.33 billion directed edges

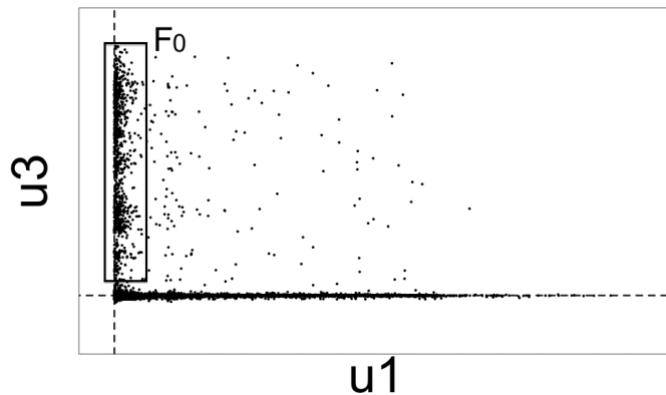


A screenshot of the Tencent Weibo website. The interface is in Chinese. At the top, there's a navigation bar with links for '首页' (Home), '爆料新鲜' (Leak Fresh), '微频道' (Micro Channel), '找人' (Find People), '活动' (Activities), and '应用' (Applications). Below the navigation is a search bar with the placeholder '来,说说你在做什么,想什么'. On the right side of the search bar is a button for '广播' (Broadcast). The main content area shows a timeline of tweets. The first tweet is from a user named '樱冢吖' with the text '洗着衣服看了一晚上的圣斗士, 是有多宅啊!!!'. It includes a timestamp '5分钟前' (5 minutes ago) and a link to '阅读(15)'. The second tweet is from a user named '杨茜' with the text '应大家要求写游记啦,论坛已经置顶啦,欢迎去踩~www.16fan.com"超勇敢沙巴游记"'. It also includes a timestamp '29分钟前' (29 minutes ago) and a link to '阅读(59)'. The third tweet is from a user named 'WIS护肤' with the text '男生脸上的痘印吧,当然是要把干净!力荐WIS: http://url.cn/Nucyw2 没痘印才叫帅!' and an emoji of a smiling face with sweat. This tweet includes a timestamp '29分钟前' (29 minutes ago) and a link to '阅读(1)'.

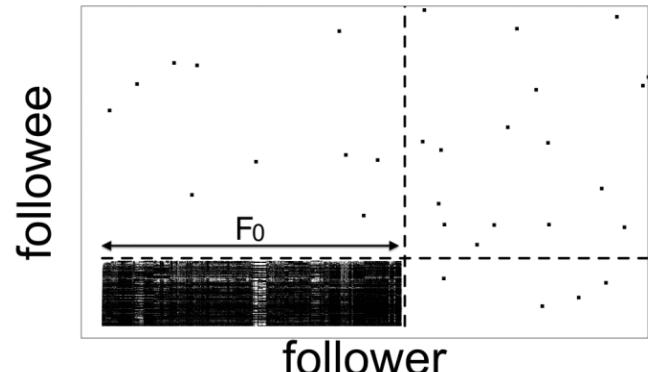
Real Data



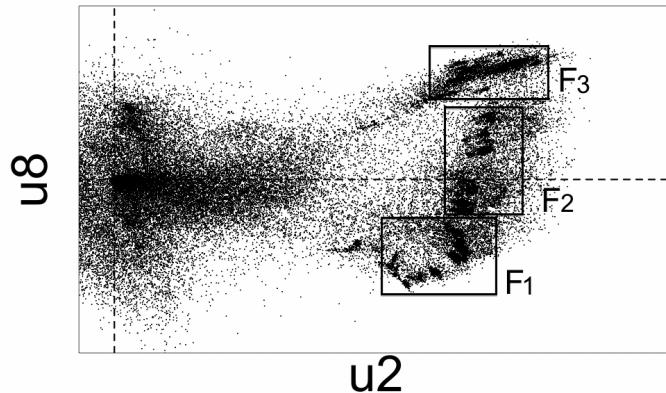
“Rays”



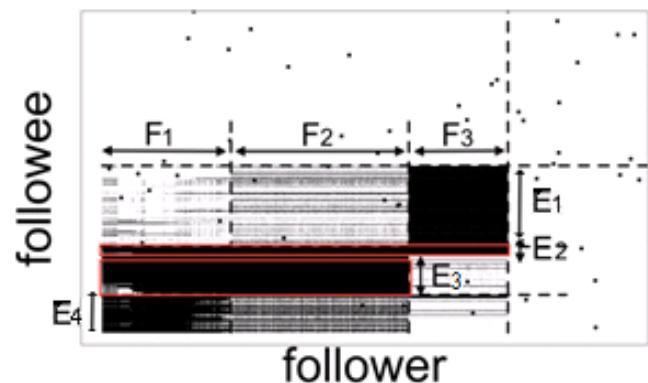
“Block”



“Pearls”



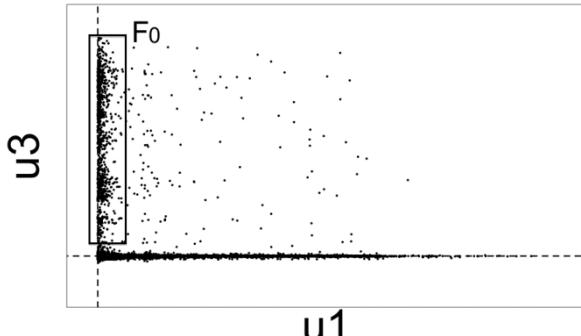
“Staircase”



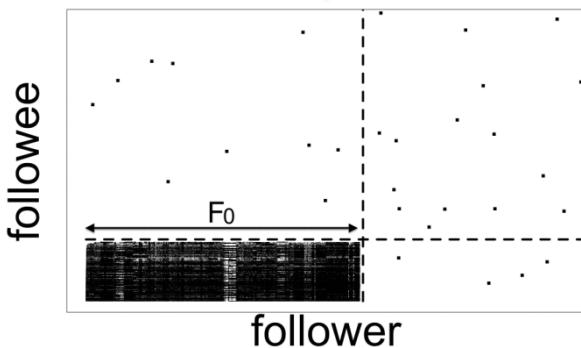
Real Data



“Rays”



“Block”

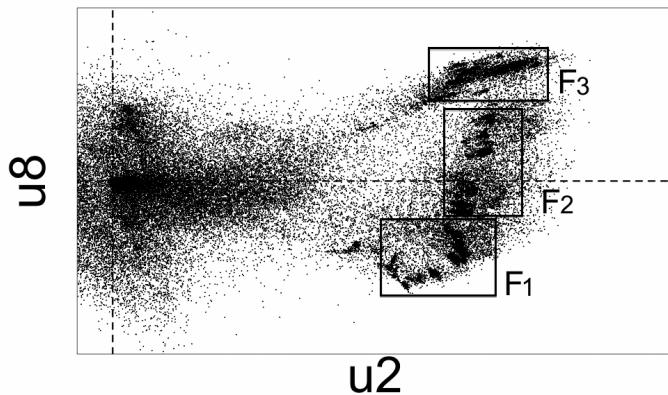


	“ray” F_0
Num. seeds	100
Size of block	$83, 208 \times 30$
Density	81.3%
Camouflage	0.14%
Fame	0.05%
Out-degree	231 ± 109
In-degree	2.0 ± 1.4

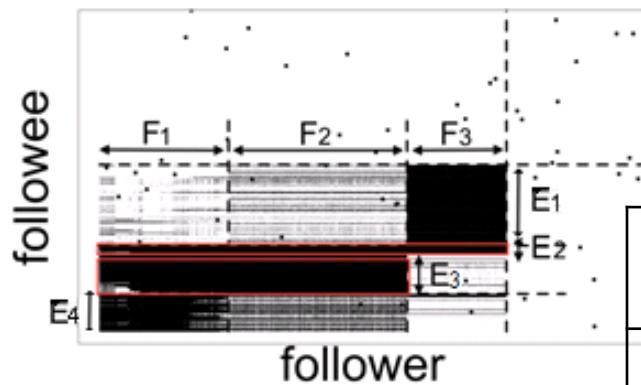
Real Data



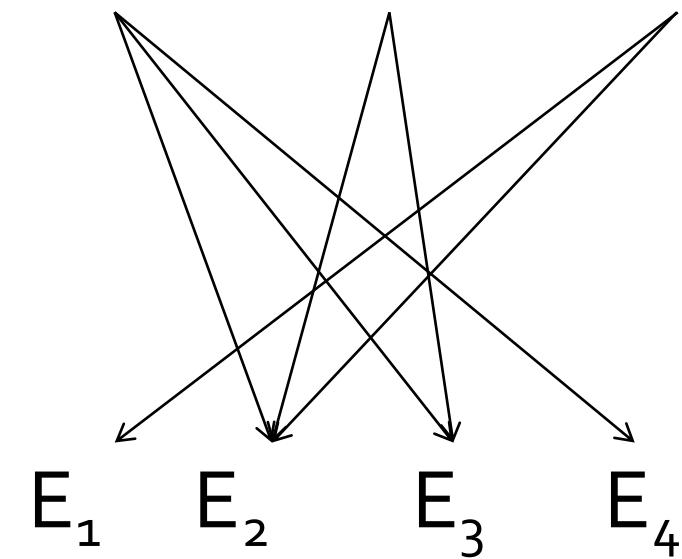
“Pearls”



“Staircase”



3,188
in F_1 7,210
in F_2 2,457
in F_3

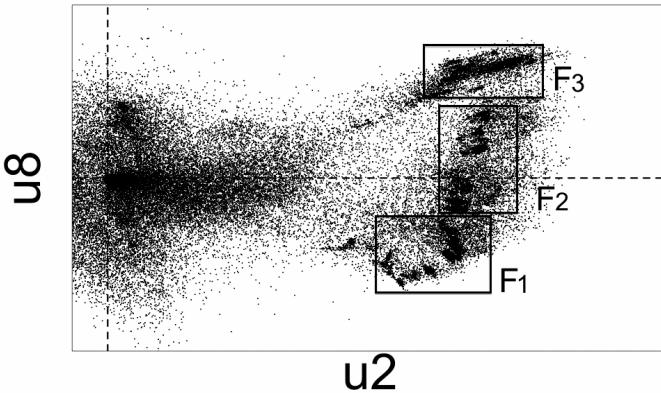


“F-E”	F_1 -...	F_2 -...	F_3 -...
Density	91.3%	92.6%	89.1%

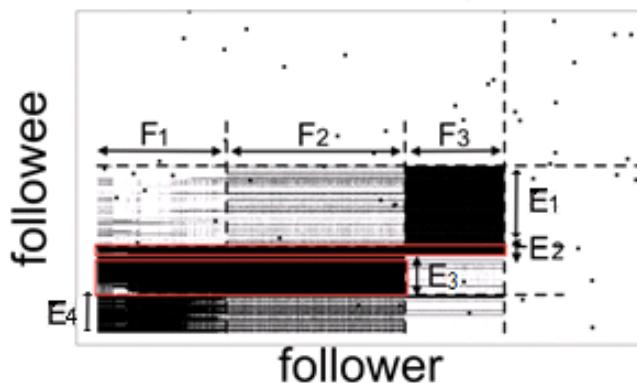
Real Data



“Pearls”



“Staircase”

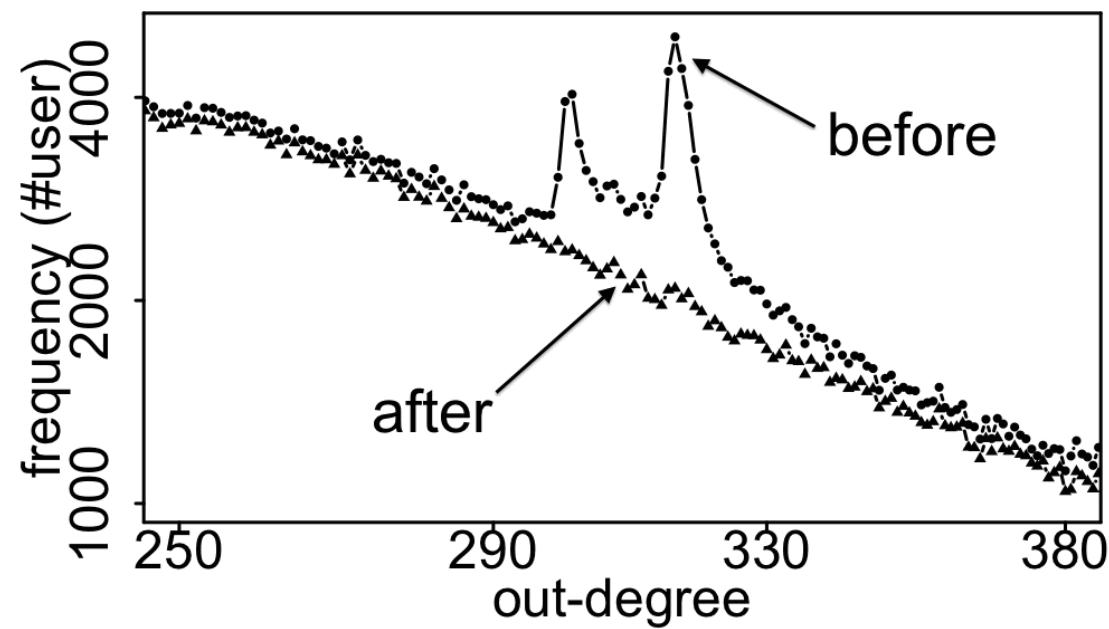
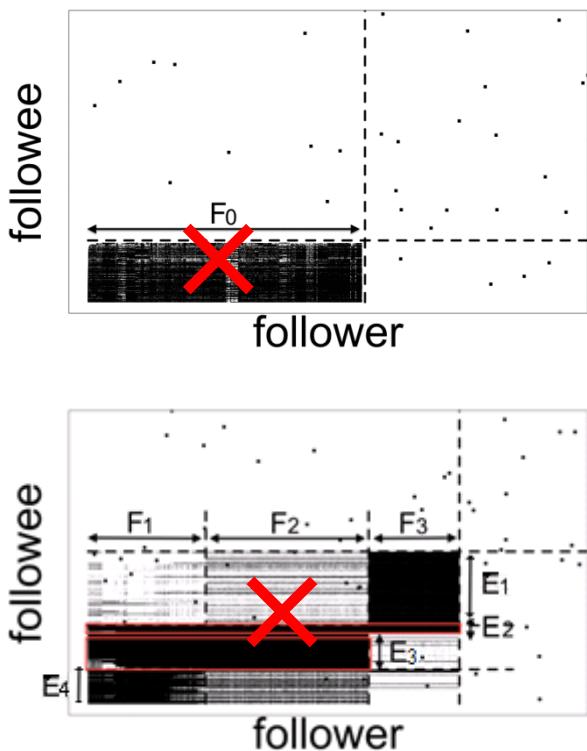


	“pearl” F ₁	“pearl” F ₂	“pearl” F ₃	“pearl” Total
	1,239	107	990	—
3, 188 × 135	7, 210 × 79	2, 457 × 148	10, 052 × 270	
91.3%	92.6%	89.1%	43.1%	
0.06%	0.10%	0.05%	0.07%	
1.93%	1.94%	1.72%	1.73%	
310±7	312±7	304±5	310±7	
8±3	10±3	11±3	12±3	

Real Data



- Spikes on the out-degree distribution



Summary: Data Reduction

- Describe **numerosity reduction** (reducing #instances)
 - Parametric methods: Fit some model and estimate model parameters
 - Regression: Describe linear/non-linear regression models
 - Nonparametric methods
 - Histograms
 - Clustering
 - Sampling: Describe stratified sampling
- Describe **dimensionality reduction** (reducing #features)
 - Feature selection
 - Heuristic search
 - Feature extraction
 - Principal component analysis (PCA)
 - Singular Value Decomposition (SVD)

References

- D. P. Ballou and G. K. Tayi. Enhancing data quality in data warehouse environments. Comm. of ACM, 42:73-78, 1999
- T. Dasu and T. Johnson. Exploratory Data Mining and Data Cleaning. John Wiley, 2003
- T. Dasu, T. Johnson, S. Muthukrishnan, V. Shkapenyuk. Mining Database Structure; Or, How to Build a Data Quality Browser. SIGMOD'02
- H. V. Jagadish et al., Special Issue on Data Reduction Techniques. Bulletin of the Technical Committee on Data Engineering, 20(4), Dec. 1997
- D. Pyle. Data Preparation for Data Mining. Morgan Kaufmann, 1999
- E. Rahm and H. H. Do. Data Cleaning: Problems and Current Approaches. IEEE Bulletin of the Technical Committee on Data Engineering. Vol.23, No.4
- V. Raman and J. Hellerstein. Potters Wheel: An Interactive Framework for Data Cleaning and Transformation, VLDB'2001
- T. Redman. Data Quality: Management and Technology. Bantam Books, 1992
- R. Wang, V. Storey, and C. Firth. A framework for analysis of data quality research. IEEE Trans. Knowledge and Data Engineering, 7:623-640, 1995