

# Sliced Inverse Regression with Interaction (Siri) Detection for non-Gaussian BN learning

Jun S. Liu

Department of Statistics  
Harvard University

Joint work with Bo Jiang

# General: Regression and Classification

	<i>Covariates</i>	<i>Responses</i>
<i>Ind 1</i>	$x_{11}, x_{12}, \dots, x_{1p}$	$Y_1$
<i>Ind 2</i>	$x_{21}, x_{22}, \dots, x_{2p}$	$Y_2$
• • •		• • •
<i>Ind N</i>	$x_{N1}, x_{N2}, \dots, x_{NP}$	$Y_N$



offthemark.com  
ATLANTIC PUBLISHING © 1998 BOB BOYD PUBLISHING

THE LAB WHERE  
THEY STUDY DRUG INTERACTION

# Variable Selection with Interaction

Let  $Y \in R$  be a univariate response variable and  $X \in R^p$  be a vector of  $p$  continuous predictor variables

$$Y = X_1 \times X_2 + \epsilon, \quad \epsilon \sim N(0, \sigma^2), \quad X \sim \text{MVN}(0, I_p)$$

Suppose  $p = 1000$ . How to find  $X_1$  and  $X_2$  ?

One step forward selection :  $\sim 500,000$  interaction terms

# Variable Selection with Interaction

Let  $Y \in R$  be a univariate response variable and  $\mathbf{X} \in R^p$  be a vector of  $p$  continuous predictor variables

$$Y = X_1 \times X_2 + \epsilon, \quad \epsilon \sim N(0, \sigma^2), \quad \mathbf{X} \sim \text{MVN}(0, I_p)$$

Suppose  $p = 1000$ . How to find  $X_1$  and  $X_2$  ?

One step forward selection :  $\sim 500,000$  interaction terms

Is there any *marginal* relationship between  $Y$  and  $X_1$  ?

[Y|X] ? [X|Y] ? Who is behind the *bar*?

How long should I wait before telling him that he's on the outside of the cage, not the inside?



# General: Regression and Classification

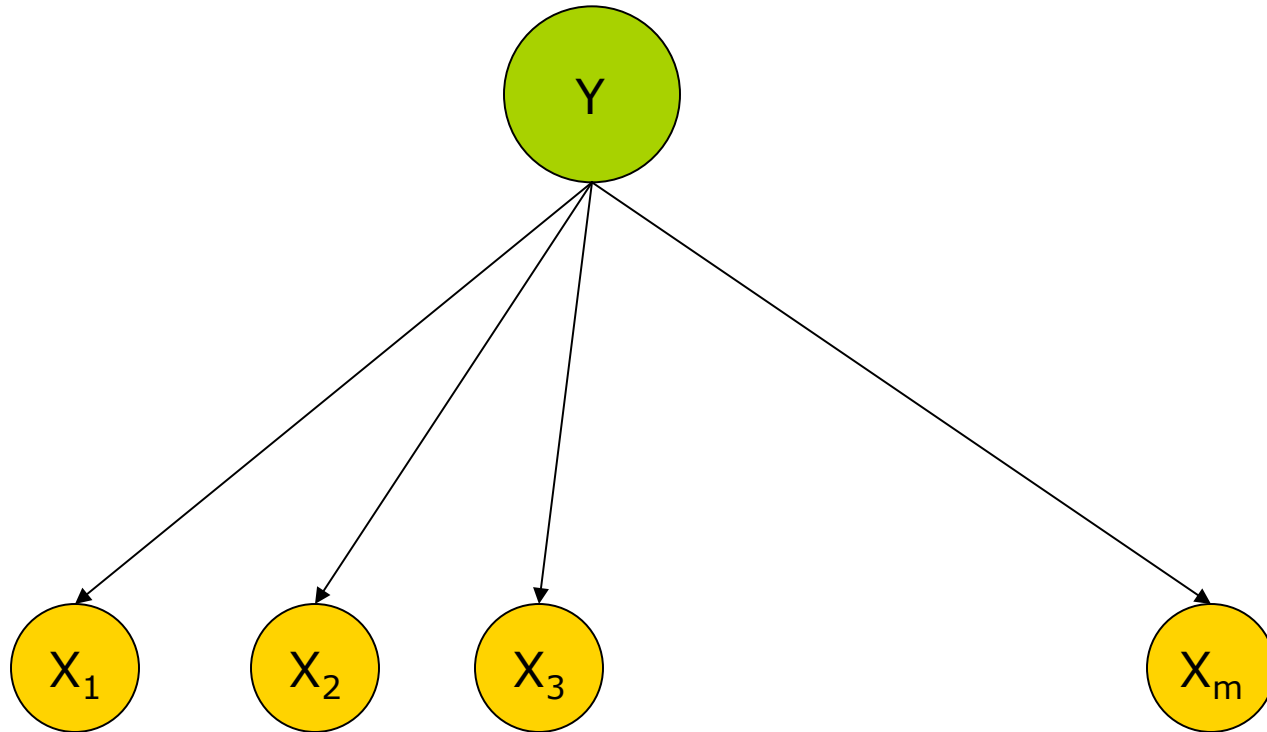
	<i>Covariates</i>	<i>Responses</i>
<i>Ind 1</i>	$\mathbf{x}_{11}, \mathbf{x}_{12}, \dots, \mathbf{x}_{1p}$	$Y_1$
<i>Ind 2</i>	$\mathbf{x}_{21}, \mathbf{x}_{22}, \dots, \mathbf{x}_{2p}$	$Y_2$
• • •		• • •
<i>Ind N</i>	$\mathbf{x}_{N1}, \mathbf{x}_{N2}, \dots, \mathbf{x}_{NP}$	$Y_N$

$$P(Y | \mathbf{X}) = P(\mathbf{X} | Y)P(Y) / P(\mathbf{X})$$

How to model this?

# Naïve Bayes model

---



$$P(Y|X_1, \dots, X_m) = \frac{P(Y) \prod_{j=1}^m P(X_j|Y)}{P(X_1, \dots, X_m)}.$$

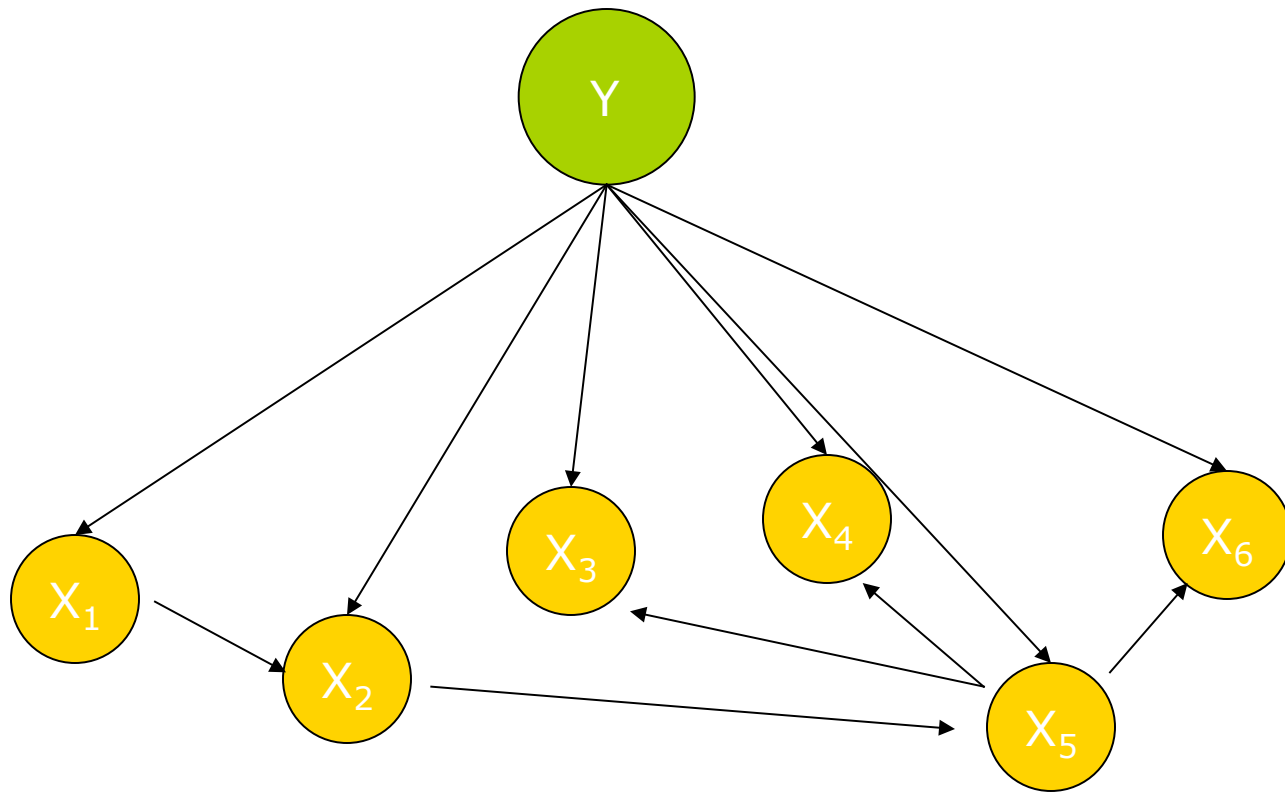


# (Augmented) Naïve Bayes Model

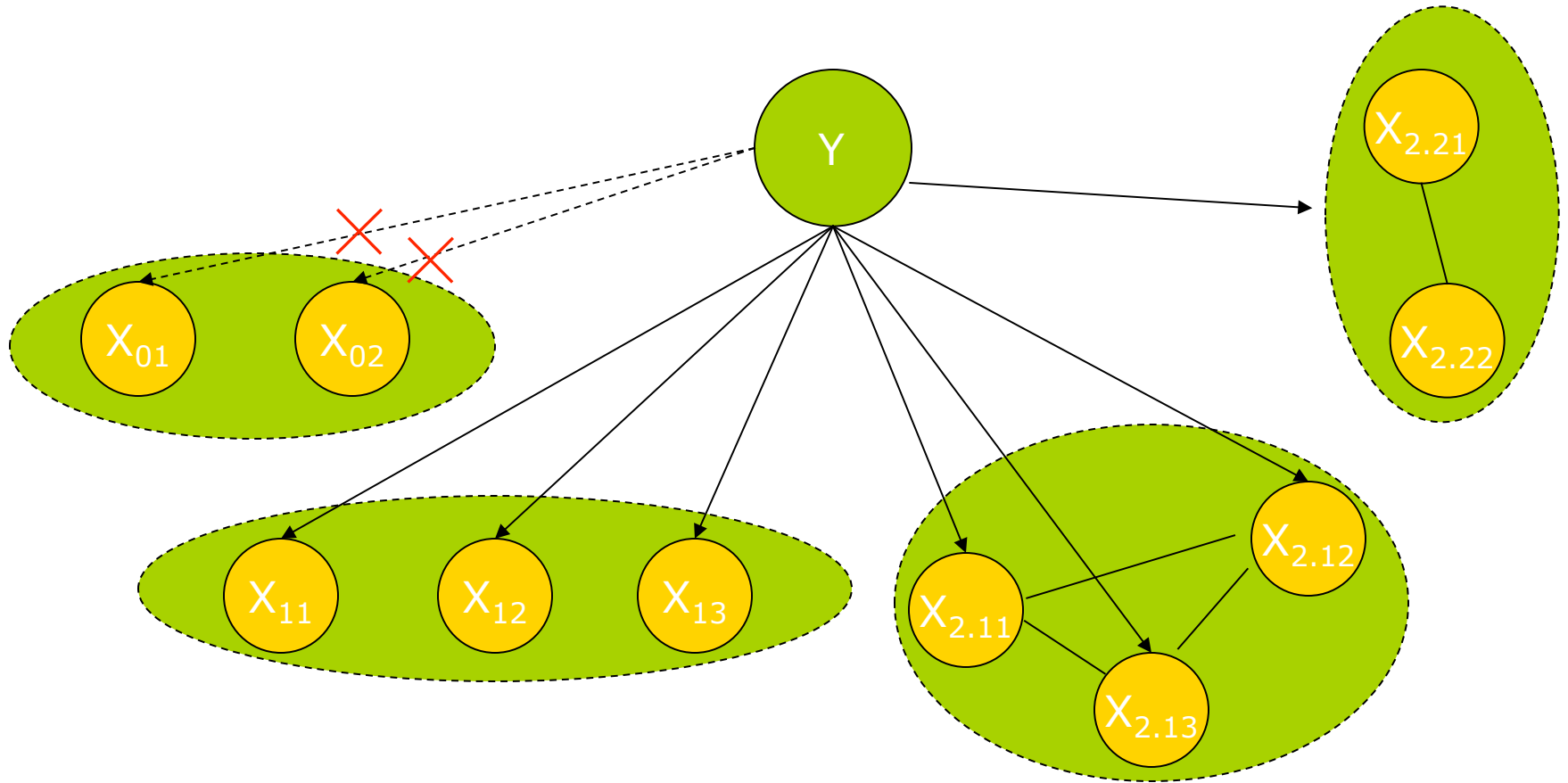
- BEAM: Bayesian Epistasis Association Mapping (Zhang and Liu 2007): discrete univariate response and discrete predictors
- (Augmented) Naïve Bayes Classifier with Variable Selection and Interaction Detection (Yuan Yuan et al.): discrete univariate response and continuous (but discretized) predictors
- Bayesian Partition Model for eQTL study (Zhang et al. 2010): continuous multivariate responses and discrete predictors
- Sliced Inverse Regression with Interaction Detection (SIRI): continuous univariate response and continuous predictors

# Tree-Augmented Naïve Bayes

---



# Augmented Naïve Bayes



# How about continuous covariates?

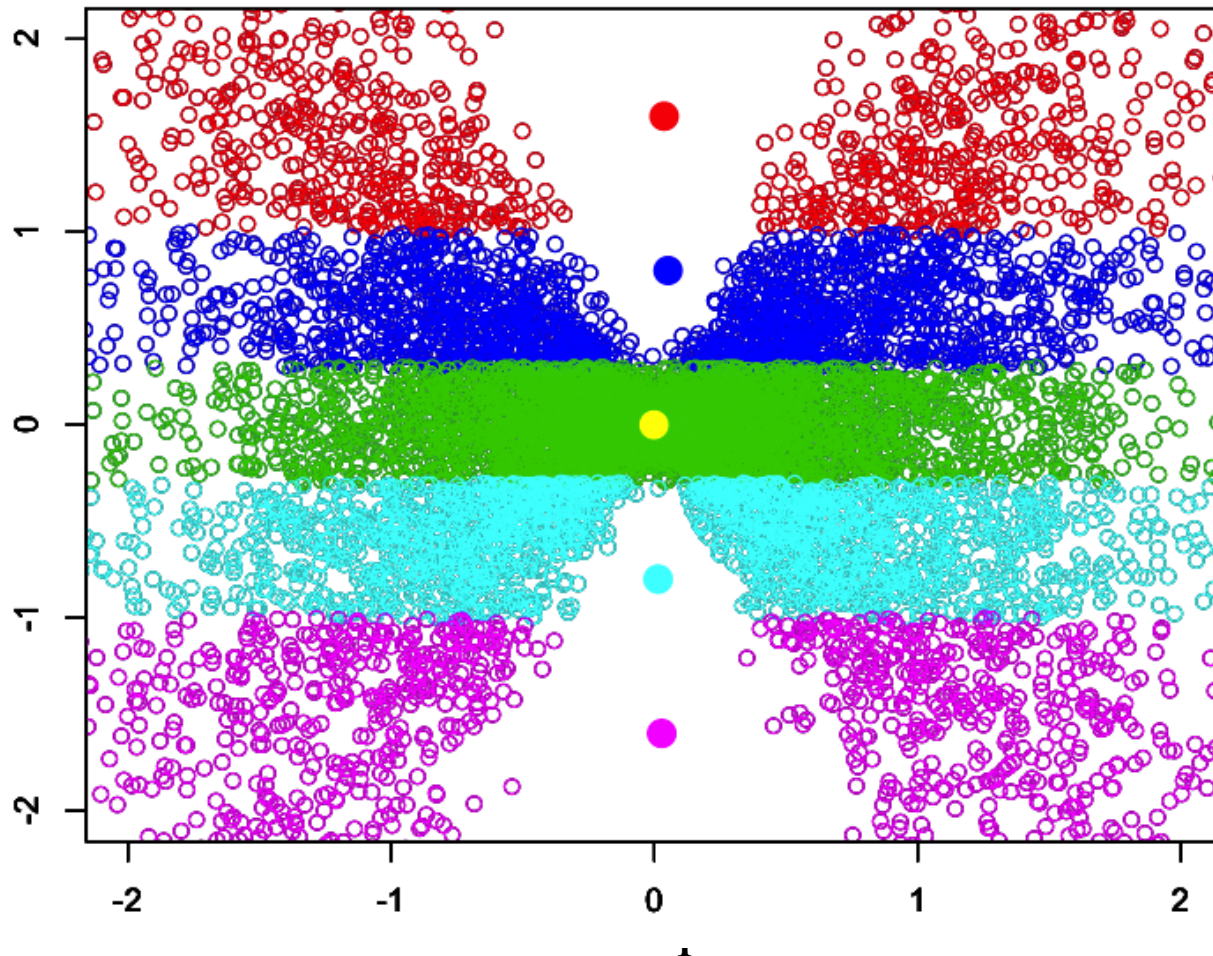
---

- We may discretize  $Y$ , and discretize each  $X$
- Or discretize  $Y$ , assuming joint Gaussian distributions on  $X$ ?
- Sound familiar?

# An observation:

$$Y = X_1 \times X_2 + \epsilon, \quad \epsilon \sim N(0, \sigma^2), \quad \mathbf{X} \sim \text{MVN}(0, I_p)$$

$y$



# Sliced Inverse Regression (SIR, Li 1991)

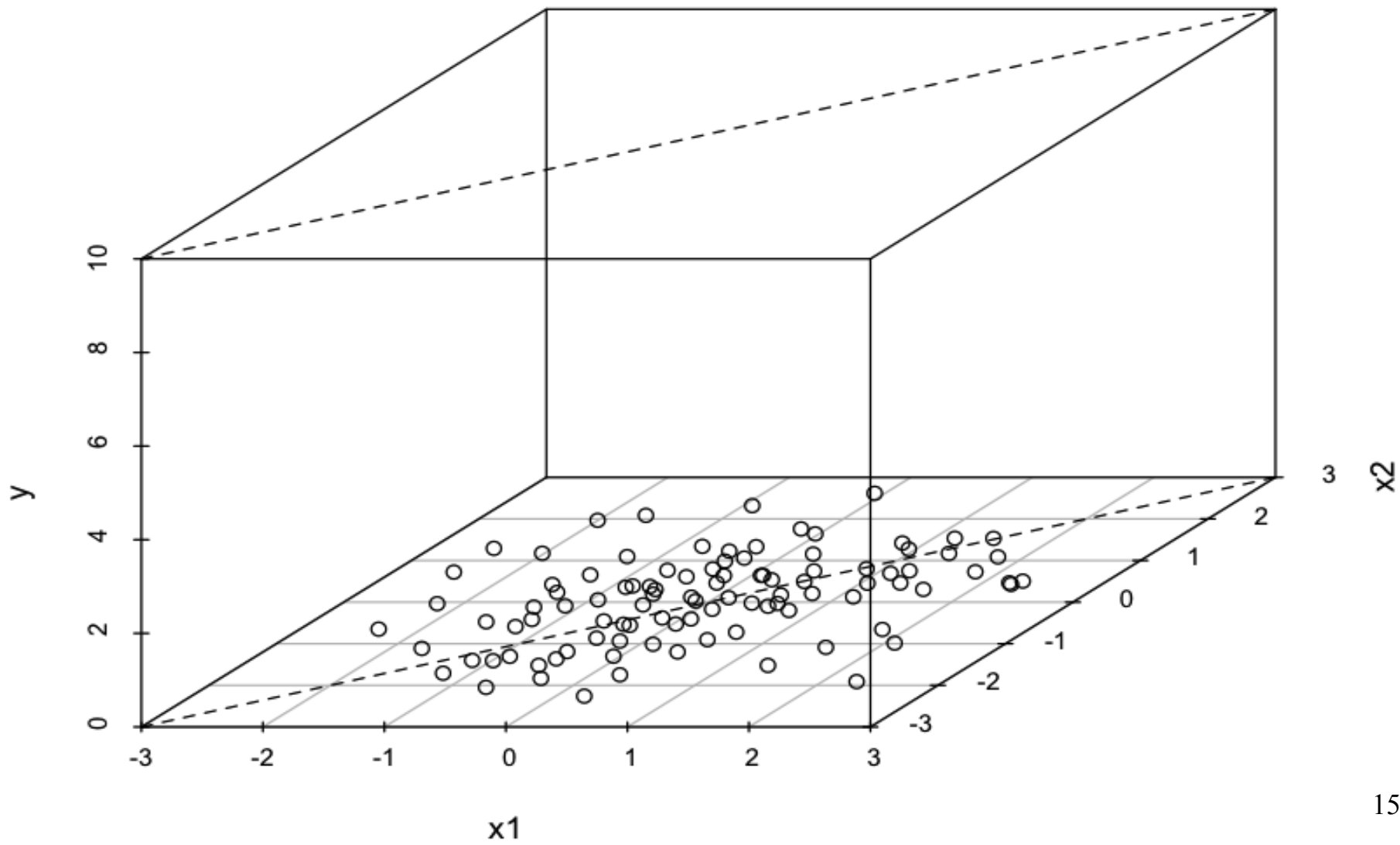
SIR is a tool for dimension reduction in multivariate statistics

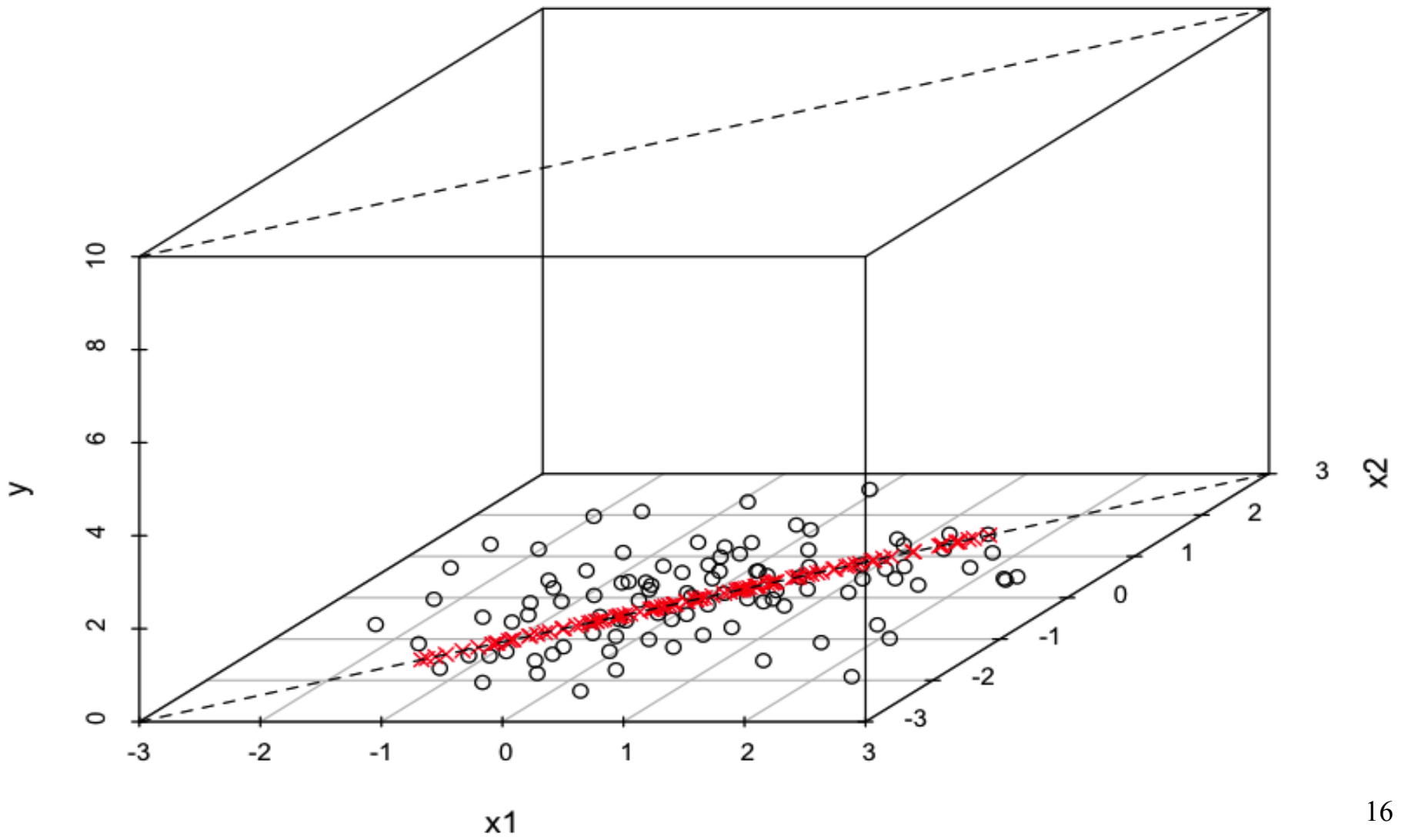
Let  $Y \in R$  be a univariate response variable and  $\mathbf{X} \in R^p$  be a vector of  $p$  continuous predictor variables

$$Y = f(\boldsymbol{\beta}_1^T \mathbf{X}, \dots, \boldsymbol{\beta}_K^T \mathbf{X}, \epsilon)$$

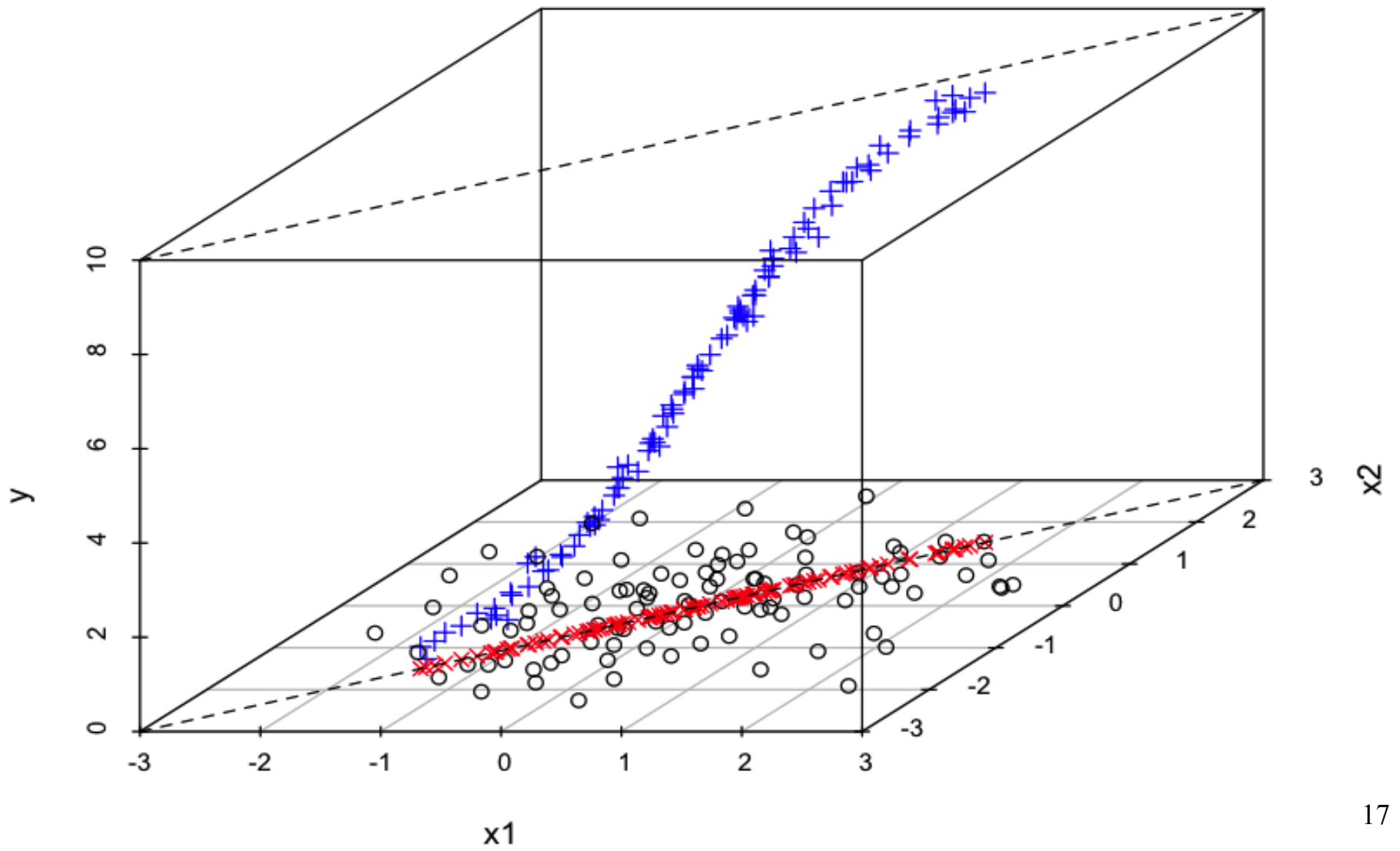
$f$  is an unknown function and  $\epsilon$  is the error with finite variance

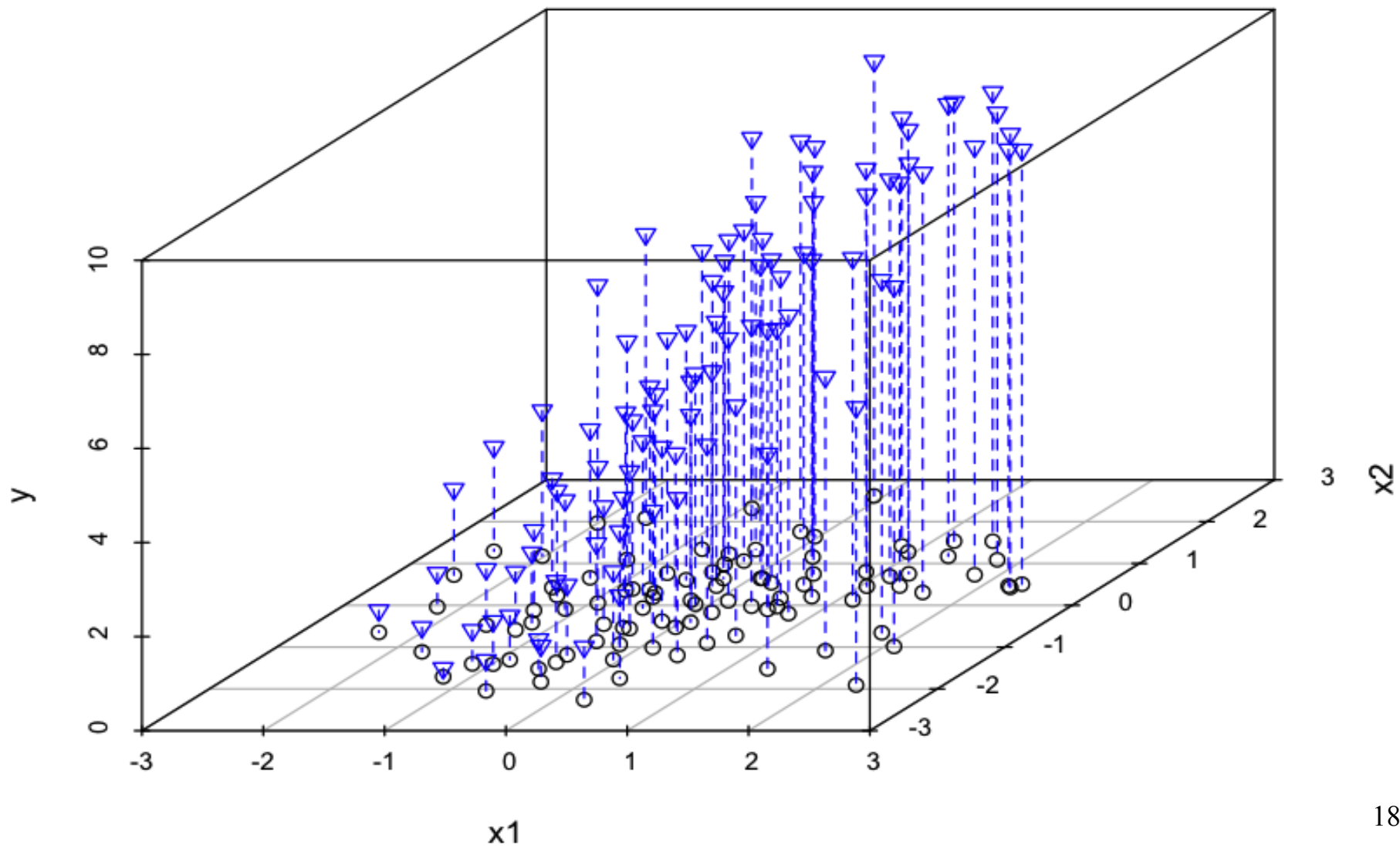
How to identify unknown projection vectors  $\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_K$ ?

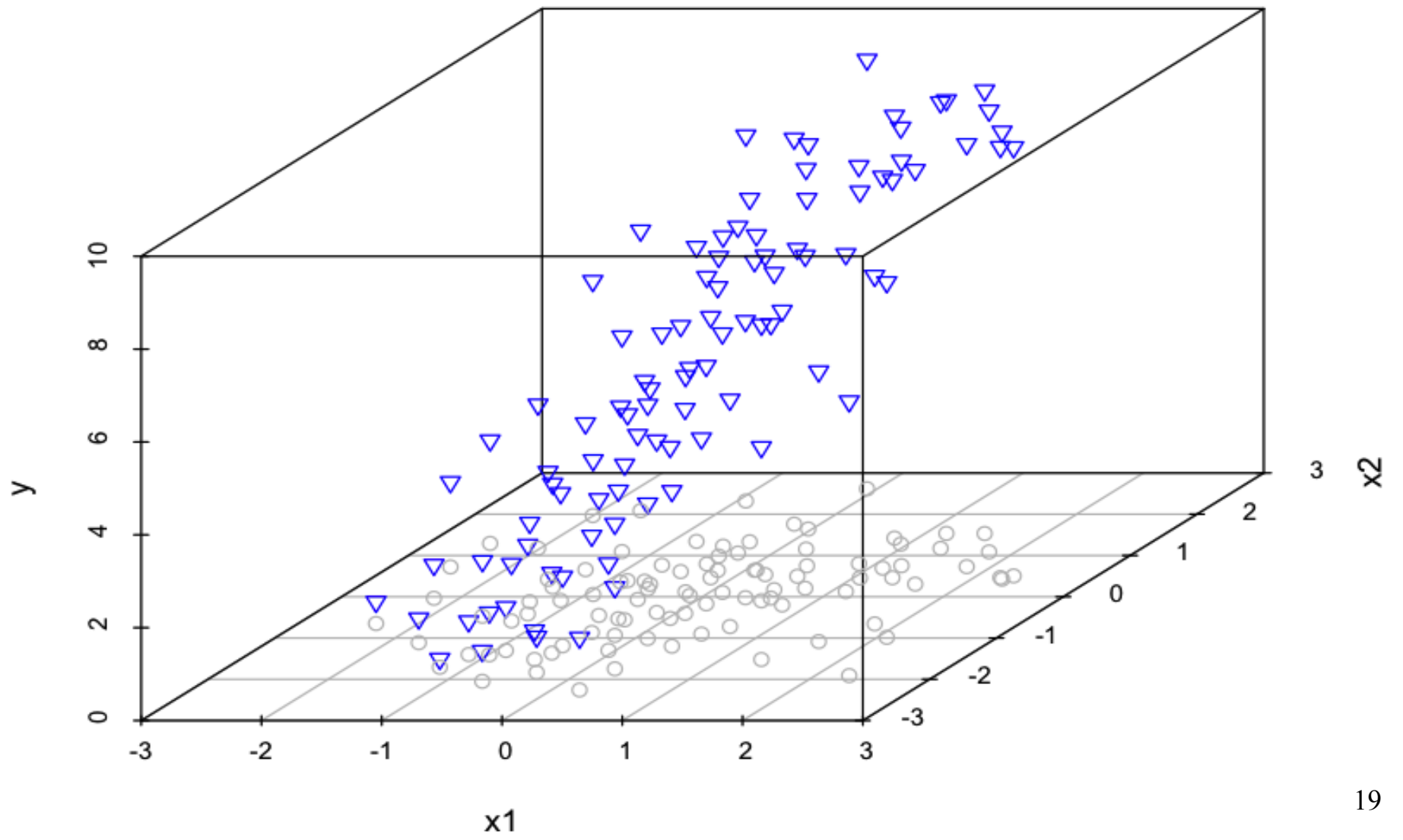


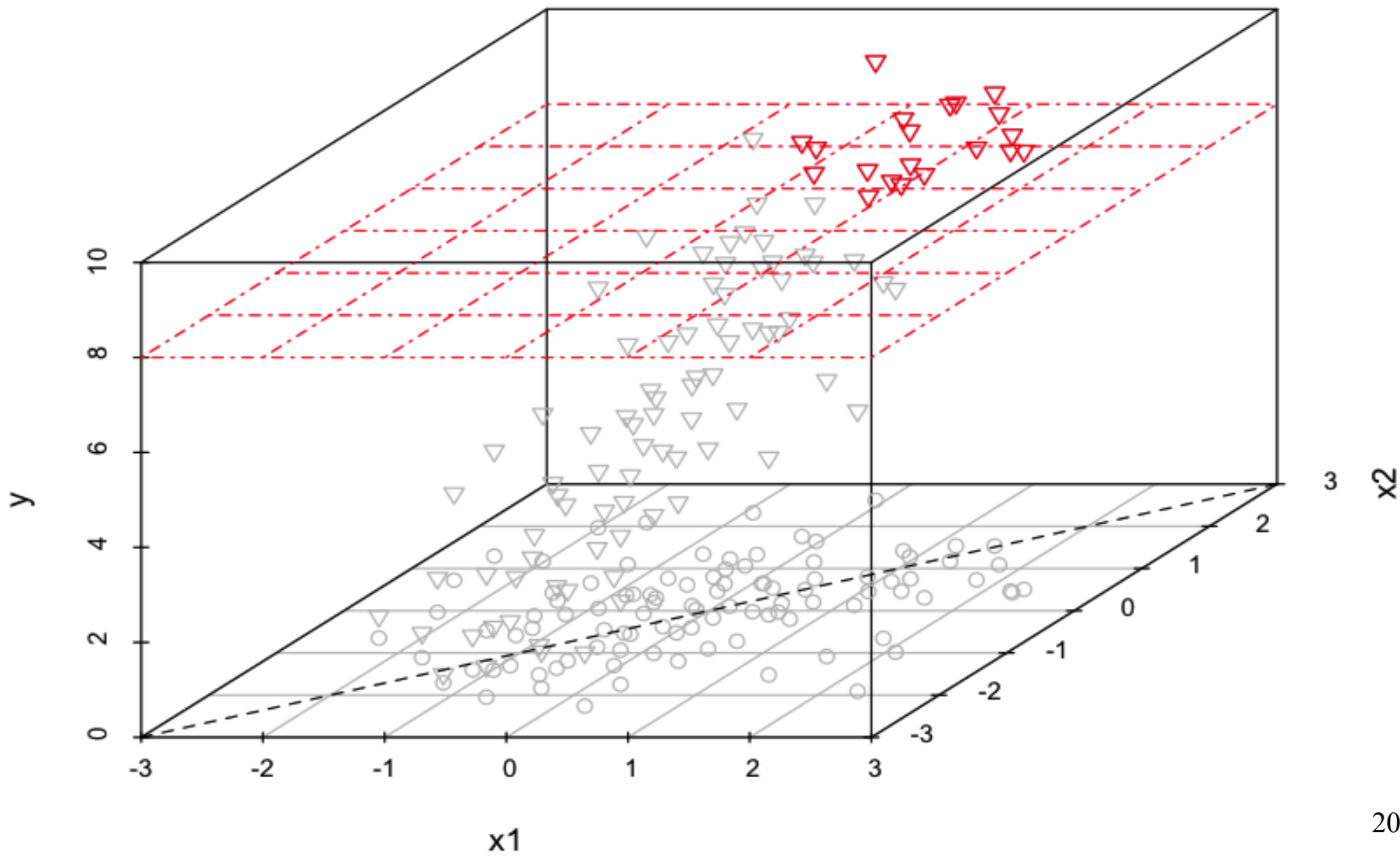


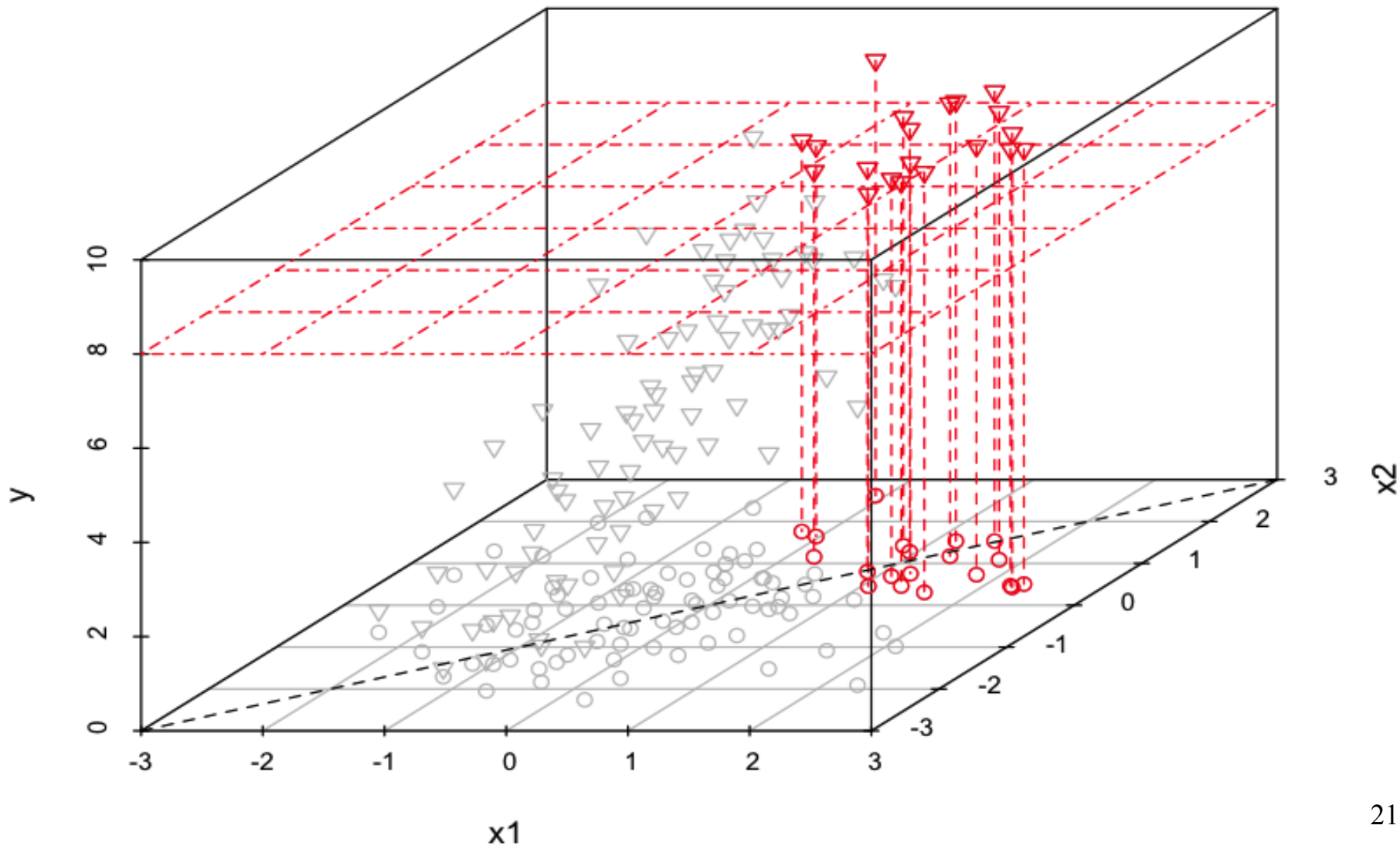


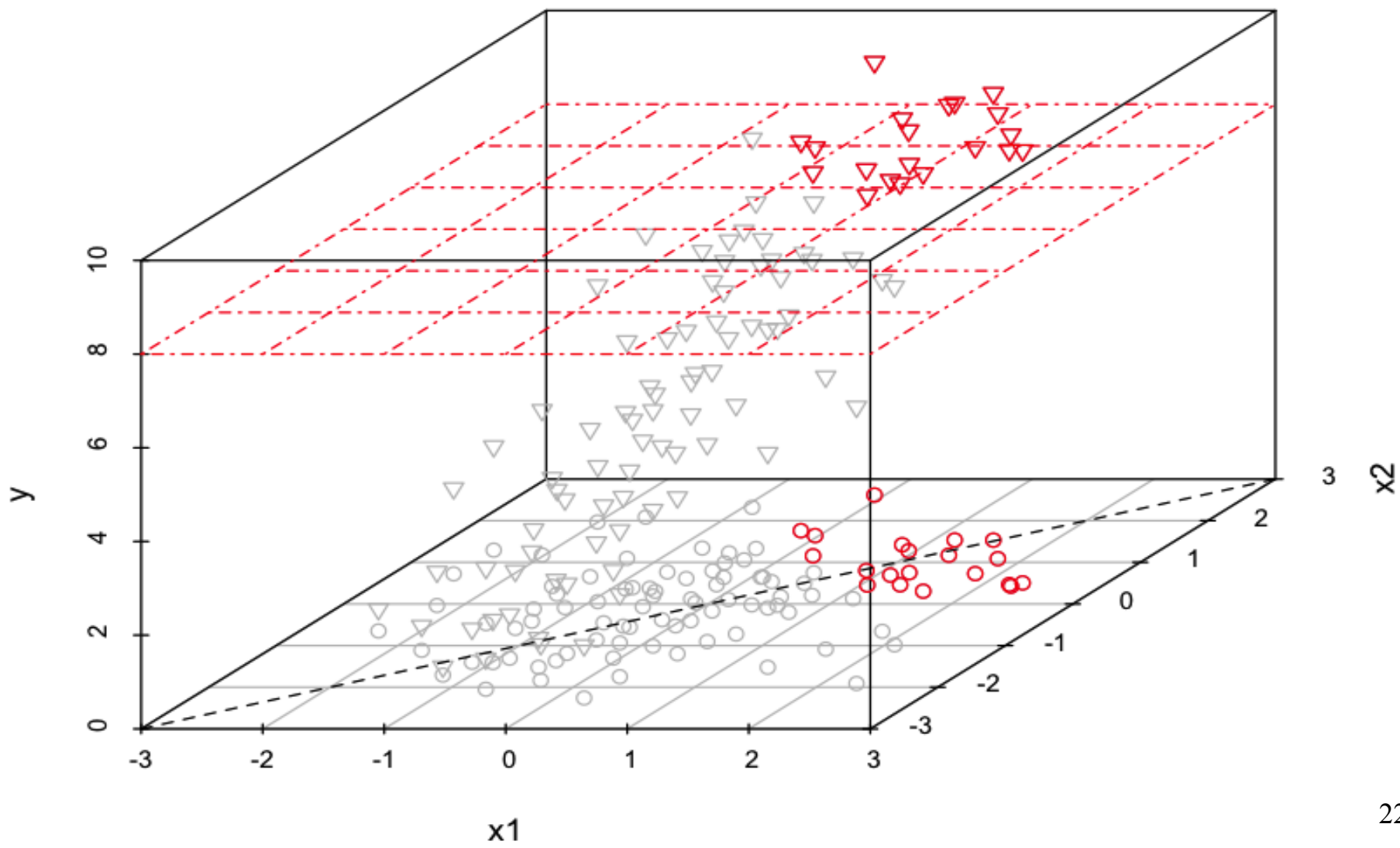


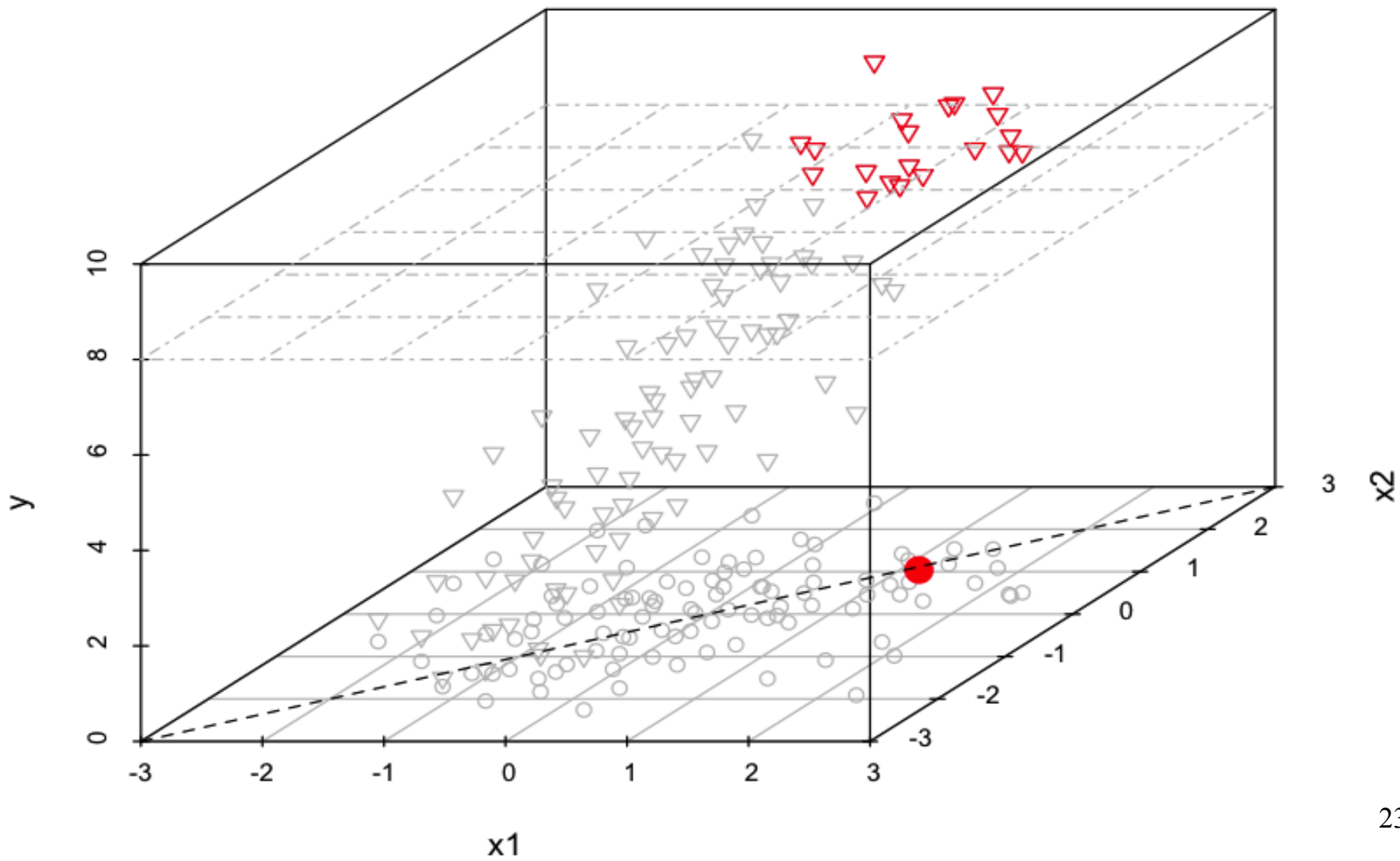


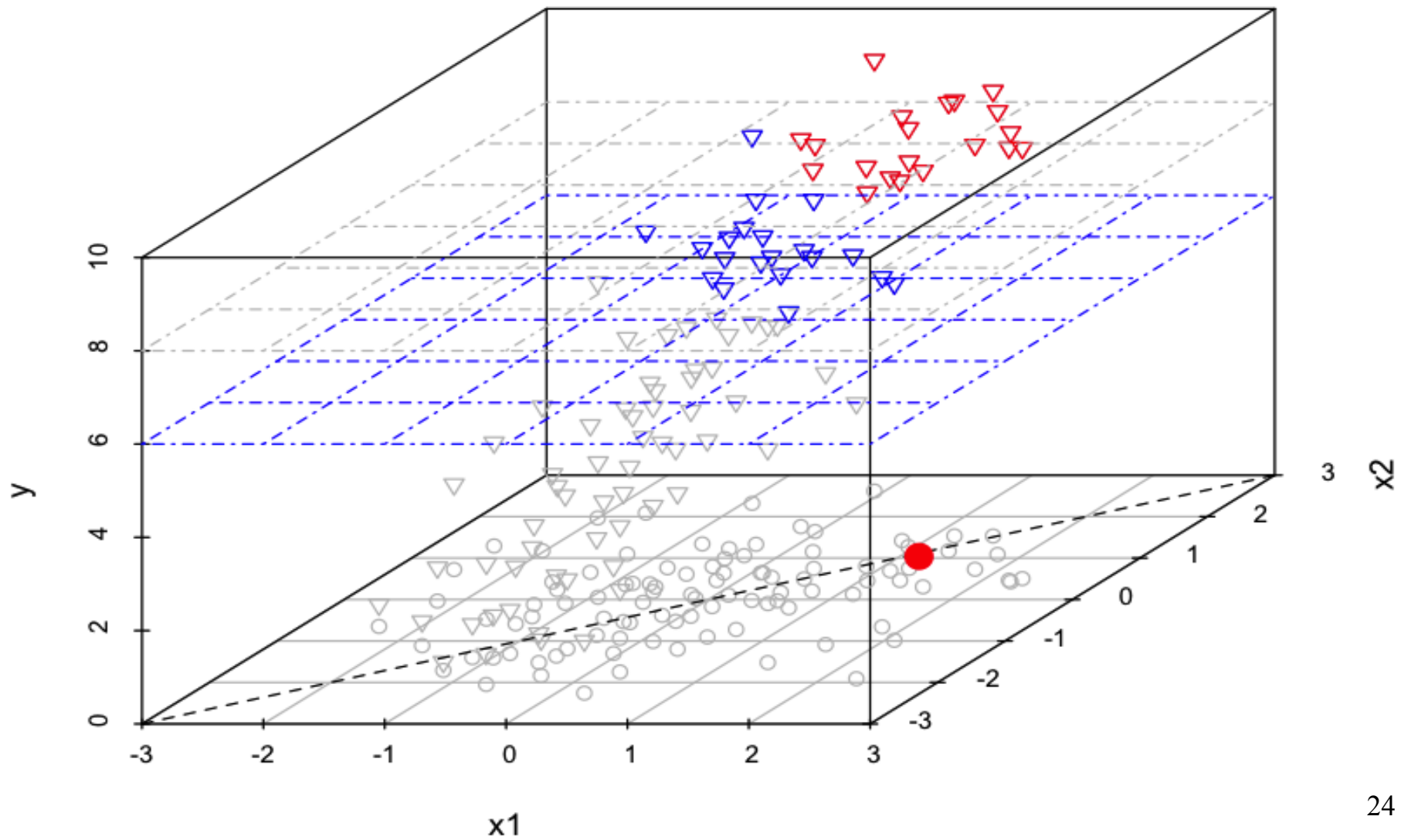




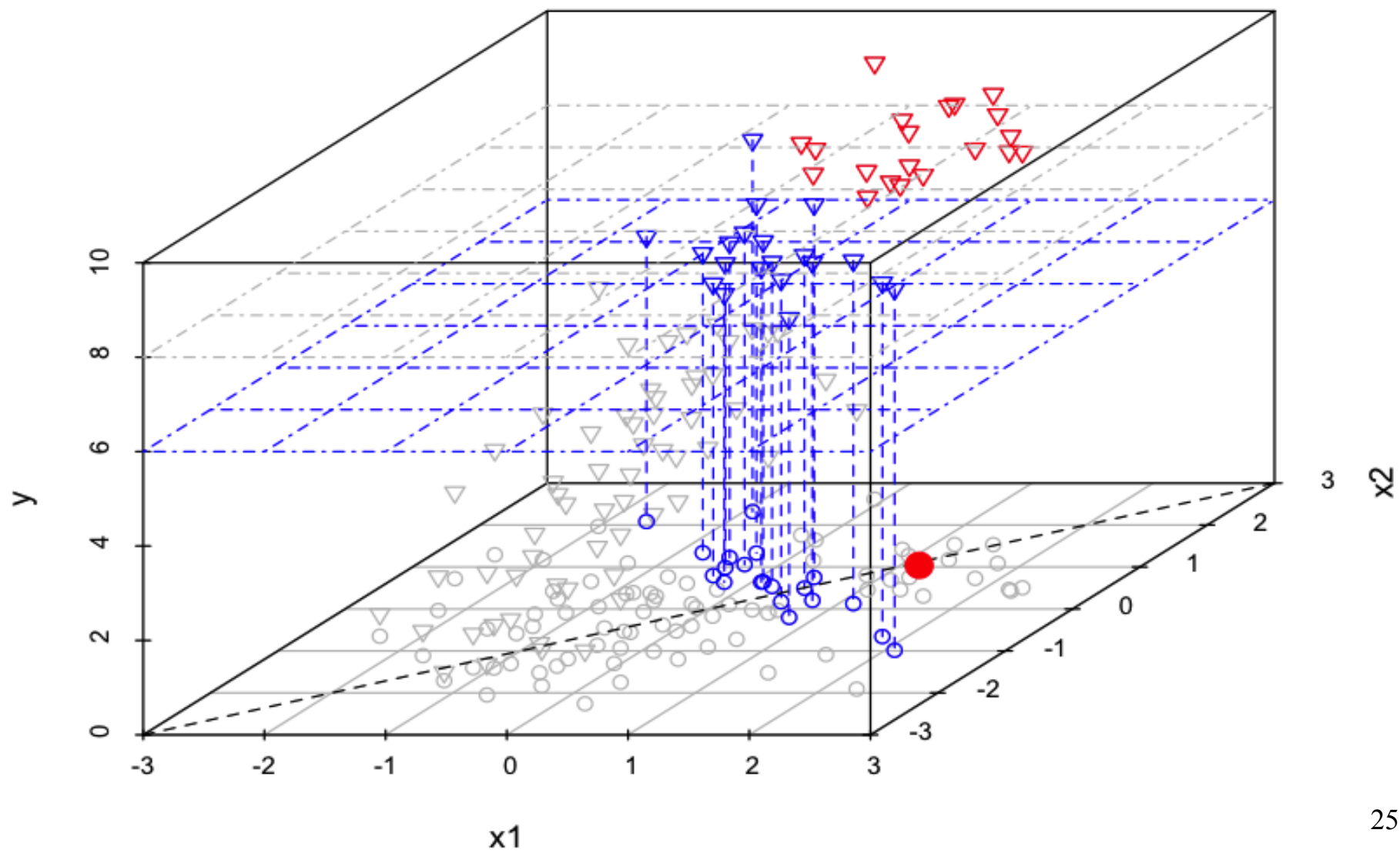


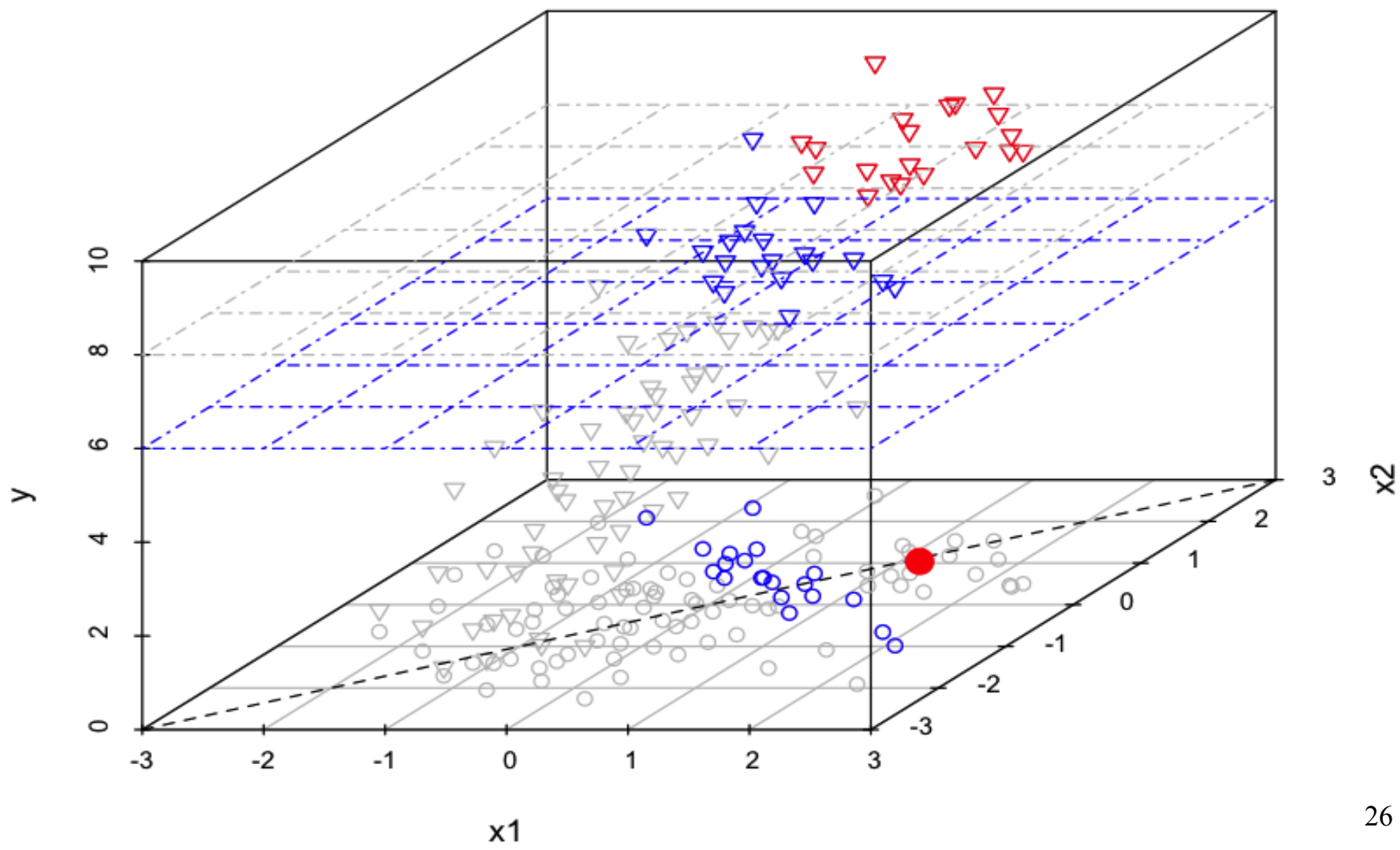


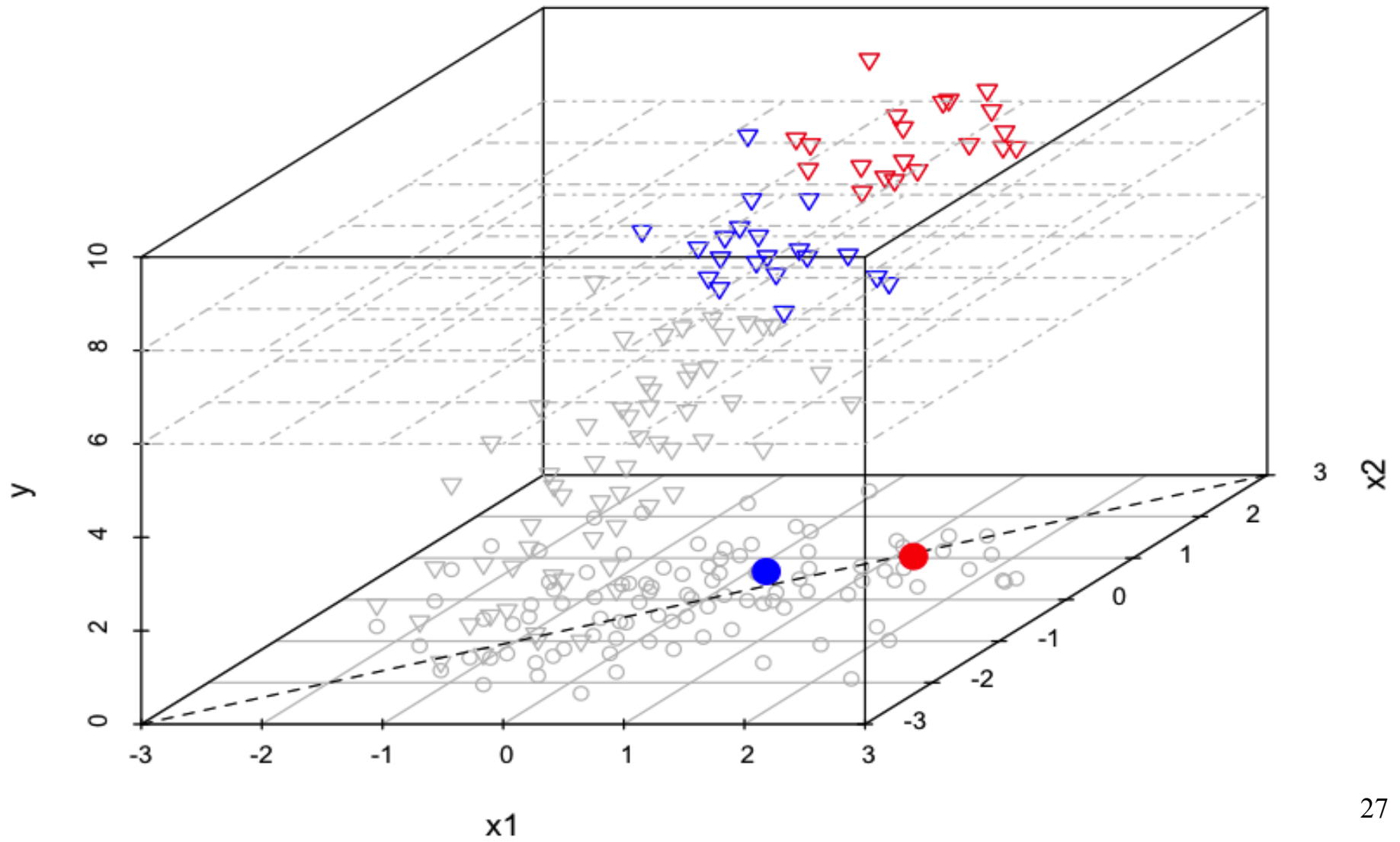


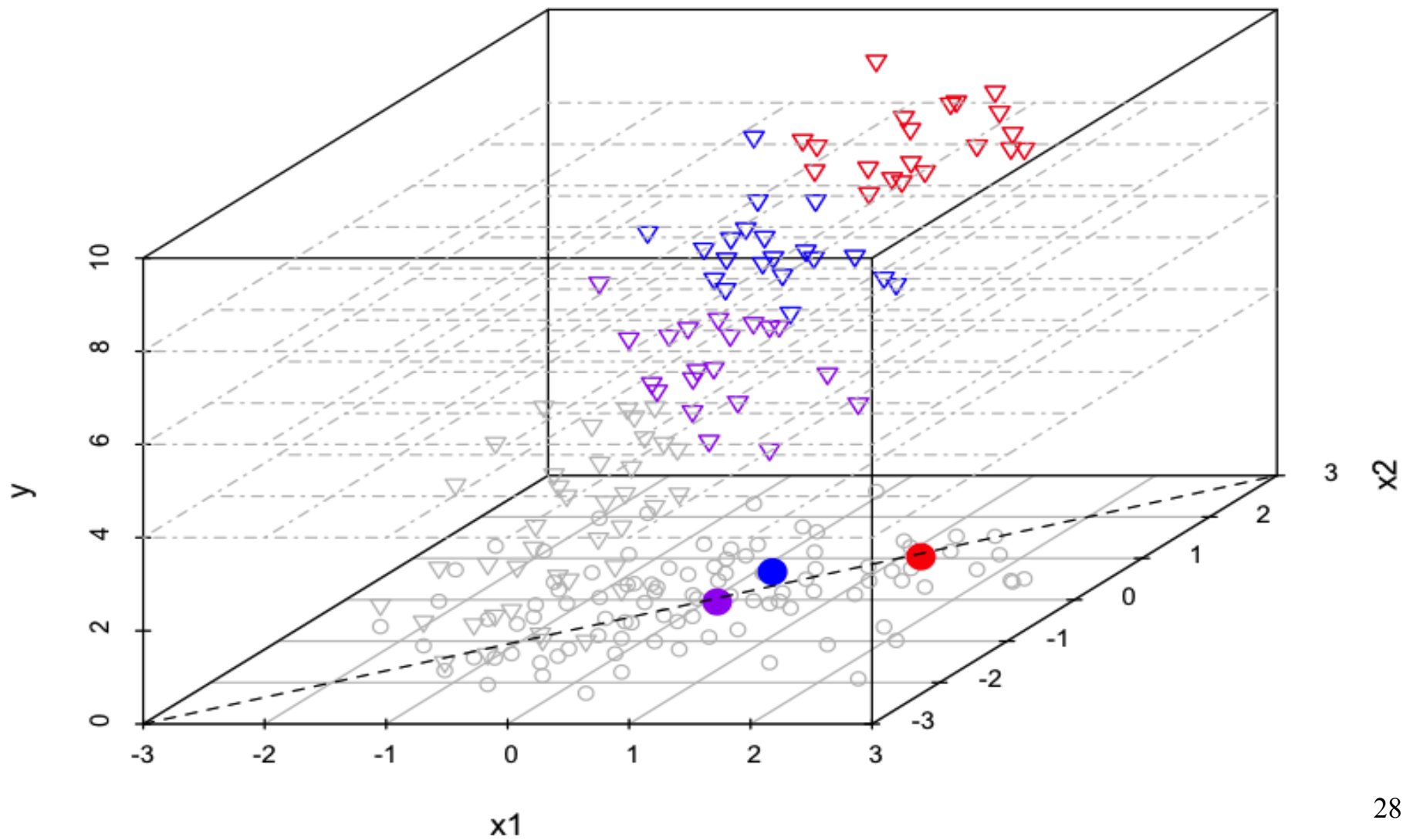


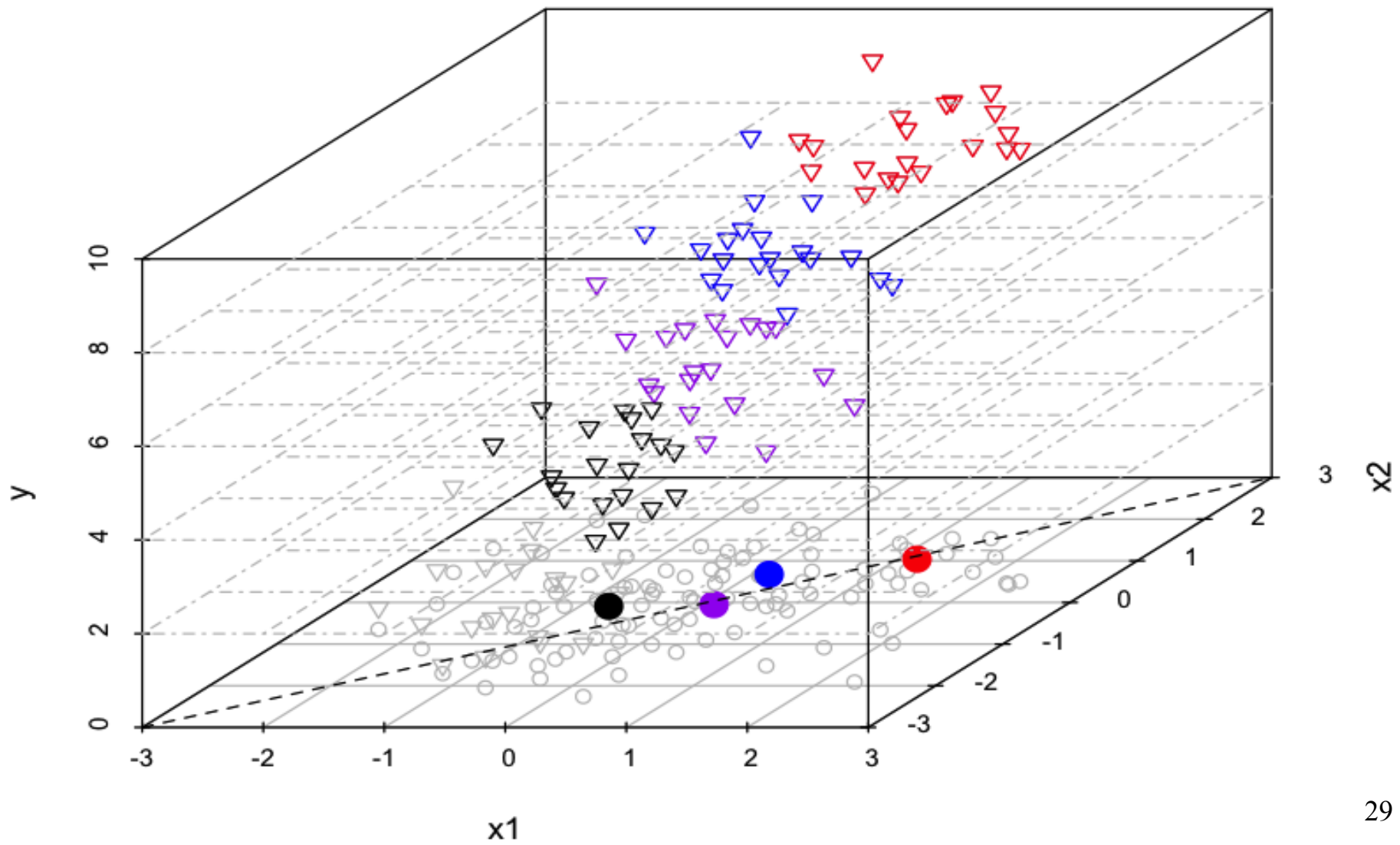


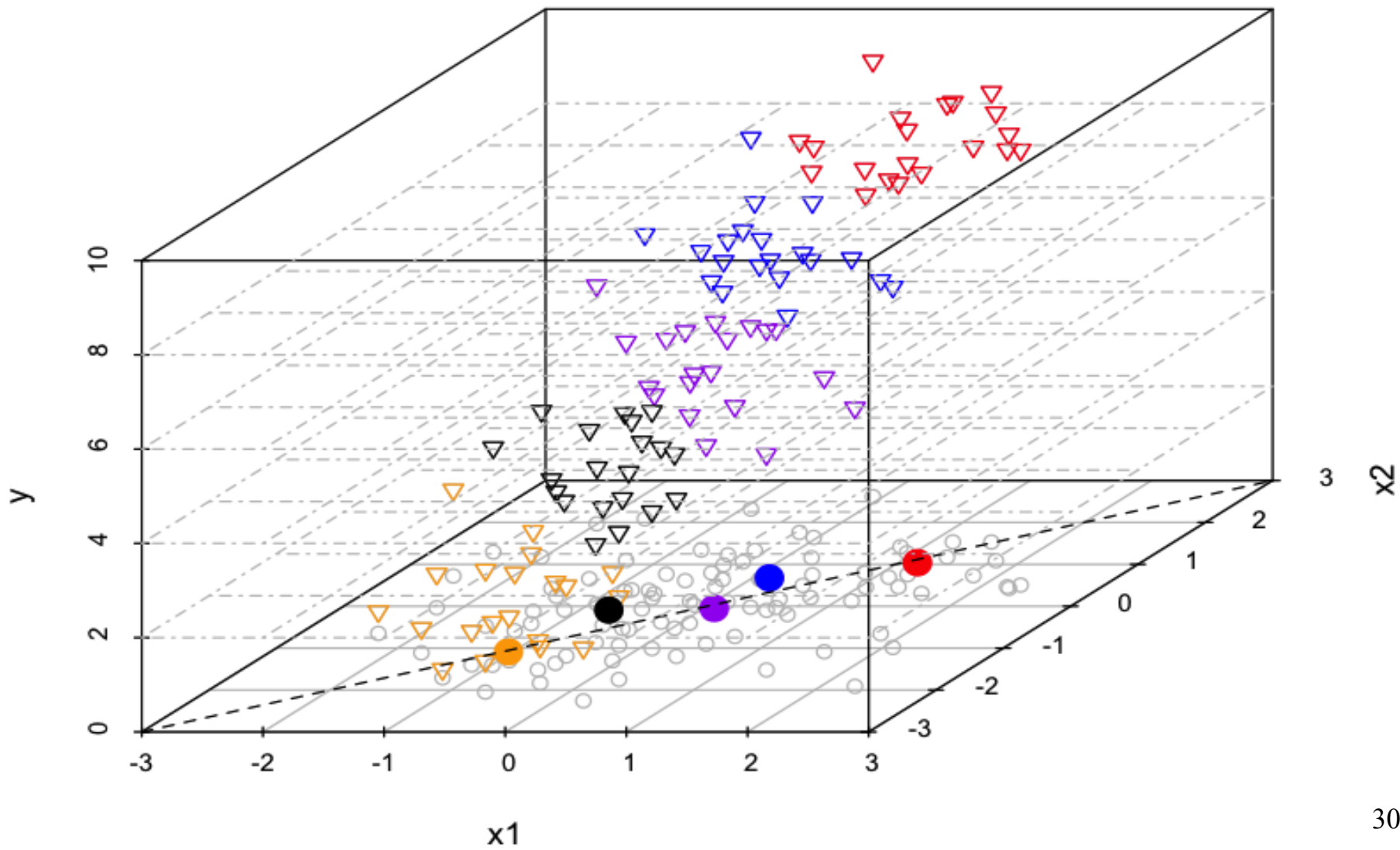












# SIR Algorithm

Let  $\Sigma_{xx}$  be the covariance matrix of  $X$ . Standardize  $X$  to :

$$Z = \Sigma_{xx}^{-1/2} \{X - E X\}$$

Divide the range of  $y_i$  into  $S$  nonoverlapping slices  $H_{s \in \{1, \dots, S\}}$   
 $n_s$  is the number of observations within each slice

Compute the mean of  $z_i$  over all slices  $\bar{z}_s = n_s^{-1} \sum_{i \in H_s} z_i$ , and  
calculate the estimate for  $Cov\{E(X|Y)\}$ :

$$\hat{M} = n^{-1} \sum_{s=1}^S n_s \bar{z}_s \bar{z}_s^T$$

Identify largest  $K$  eigenvalues of  $\hat{M}$ ,  $\hat{\lambda}_k$  and corresponding  
eigenvectors  $\hat{\eta}_k$ . Then,

$$\hat{\beta}_k = \hat{\Sigma}_{xx}^{-1/2} \hat{\eta}_k \quad (k = 1, \dots, K)$$

# SIR with Variable Selection

Only a subset of predictors are relevant:  $\beta_1, \dots, \beta_K$  are sparse

Backward subset selection (Cook 2004, Li et al. 2005)

Shrinkage estimates of  $\beta_1, \dots, \beta_K$  using  $L_1$ - or  $L_2$ -penalty :

Regularized SIR (RSIR, Zhong et al. 2005)

Sparse SIR (SSIR, Li 2007)

Correlation Pursuit (Zhong et al. 2012) : A forward selection and backward elimination procedure motivated by F-test in stepwise regression

$$F_{1, n-d-1} = (n-d-1) \frac{(\hat{R}_{d+1}^2 - \hat{R}_d^2)}{1 - \hat{R}_{d+1}^2}$$



# Correlation Pursuit (COP)

Let  $A$  be the current set of selected predictors and  $\hat{\lambda}_k^A$  the  $k$ th largest eigenvalue estimated by SIR based on predictors in  $A$ . For  $j$ th predictor ( $j \notin A$ ),  $X_j$ , define statistic

$$COP_k^{A+j} = n \frac{(\hat{\lambda}_k^{A+j} - \hat{\lambda}_k^A)}{1 - \hat{\lambda}_k^{A+j}}$$

# Correlation Pursuit (COP)

Let  $A$  be the current set of selected predictors and  $\hat{\lambda}_k^A$  the  $k$ th largest eigenvalue estimated by SIR based on predictors in  $A$ . For  $j$ th predictor ( $j \notin A$ ),  $X_j$ , define statistic

$$COP_k^{A+j} = n \frac{(\hat{\lambda}_k^{A+j} - \hat{\lambda}_k^A)}{1 - \hat{\lambda}_k^{A+j}}$$

If  $j \notin A$ ,  $COP_k^{A+j}$  ( $k=1, \dots, K$ ) are asymptotically i.i.d.  $\chi^2(1)$ , and  $COP_{1:K}^{A+j} = \sum_{k=1}^K COP_k^{A+j}$  is asymptotically  $\chi^2(K)$

# Correlation Pursuit (COP)

Let  $A$  be the current set of selected predictors and  $\hat{\lambda}_k^A$  the  $k$ th largest eigenvalue estimated by SIR based on predictors in  $A$ . For  $j$ th predictor ( $j \notin A$ ),  $X_j$ , define statistic

$$COP_k^{A+j} = n \frac{(\hat{\lambda}_k^{A+j} - \hat{\lambda}_k^A)}{1 - \hat{\lambda}_k^{A+j}}$$

If  $j \notin A$ ,  $COP_k^{A+j}$  ( $k=1, \dots, K$ ) are asymptotically i.i.d.  $\chi^2(1)$ , and  $COP_{1:K}^{A+j} = \sum_{k=1}^K COP_k^{A+j}$  is asymptotically  $\chi^2(K)$ .

The stepwise procedure is consistent if  $p = O(n^r)$ ,  $r < 1/2$ .

# Correlation Pursuit (COP)

Let  $A$  be the current set of selected predictors and  $\hat{\lambda}_k^A$  the  $k$ th largest eigenvalue estimated by SIR based on predictors in  $A$ . For  $j$ th predictor ( $j \notin A$ ),  $X_j$ , define statistic

$$COP_k^{A+j} = n \frac{(\hat{\lambda}_k^{A+j} - \hat{\lambda}_k^A)}{1 - \hat{\lambda}_k^{A+j}}$$

If  $j \notin A$ ,  $COP_k^{A+j}$  ( $k = 1, \dots, K$ ) are asymptotically i.i.d.  $\chi^2(1)$ , and  $COP_{1:K}^{A+j} = \sum_{k=1}^K COP_k^{A+j}$  is asymptotically  $\chi^2(K)$ .

The stepwise procedure is consistent if  $p = O(n^r)$ ,  $r < 1/2$ .

Dimension  $K$  and threshold in forward selection (backward elimination) are chosen by cross-validation.

# SIR via MLE

Let  $A$  be the set of relevant predictors and  $C = A^c$ ,  $d = |A|$

$$X_A | Y \in H_s \sim N(\mu_s, \Sigma)$$

# SIR via MLE

Let  $A$  be the set of relevant predictors and  $C = \neg A$ ,  $d = |A|$

$$X_A | Y \in H_s \sim N(\mu_s, \Sigma)$$

$$X_C | X_A, Y \in H_s \sim N(X_A \beta, \Sigma_0)$$

# SIR via MLE

Let  $A$  be the set of relevant predictors and  $C = A^C$ ,  $d = |A|$

$$X_A | Y \in H_s \sim N(\mu_s, \Sigma)$$

$$X_C | X_A, Y \in H_s \sim N(X_A \beta, \Sigma_0)$$

$\mu_s = \alpha + \Gamma \gamma_s$ , where  $\gamma_s \in R^K$  and  $\Gamma$  is a  $d \times K$  orthogonal matrix

# SIR via MLE

Let  $A$  be the set of relevant predictors and  $C = A^c$ ,  $d = |A|$

$$X_A | Y \in H_s \sim N(\mu_s, \Sigma)$$

$$X_C | X_A, Y \in H_s \sim N(X_A \beta, \Sigma_0)$$

$\mu_s = \alpha + V^K$ , belongs to a  $K$ -dimensional affine space ( $K < d$ )



# SIR via MLE

Let  $A$  be the set of relevant predictors and  $C = A^c$ ,  $d = |A|$

$$X_A | Y \in H_s \sim N(\mu_s, \Sigma)$$

$$X_C | X_A, Y \in H_s \sim N(X_A \beta, \Sigma_0)$$

$\mu_s = \alpha + V^K$ , belongs to a  $K$ -dimensional affine space ( $K < d$ )

MLE of the span of subspace  $V^K$  coincides with SIR directions  
(Cook 2007, Szretter and Yohai 2009)

# SIR via MLE

Let  $A$  be the set of relevant predictors and  $C = \neg A$ ,  $d = |A|$

$$X_A | Y \in H_s \sim N(\mu_s, \Sigma)$$

$$X_C | X_A, Y \in H_s \sim N(X_A \beta, \Sigma_0)$$

$\mu_s = \alpha + V^K$ , belongs to a  $K$ -dimensional affine space ( $K < d$ )

MLE of the span of subspace  $V^K$  coincides with SIR directions  
(Cook 2007, Szretter and Yohai 2009)

Given current  $A$  and predictor  $X_{j \notin A}$ , we want to test

$H_0: X_j$  is irrelevant, vs.  $H_1: X_j$  is relevant

# SIR via MLE

Let  $A$  be the set of relevant predictors and  $C = A^c$ ,  $d = |A|$

$$\mathbf{X}_A | Y \in H_s \sim N(\boldsymbol{\mu}_s, \boldsymbol{\Sigma})$$

$$\mathbf{X}_C | \mathbf{X}_A, Y \in H_s \sim N(\mathbf{X}_A \boldsymbol{\beta}, \boldsymbol{\Sigma}_0)$$

$\boldsymbol{\mu}_s = \boldsymbol{\alpha} + \mathbf{V}^K$ , belongs to a  $K$ -dimensional affine space ( $K < d$ )

MLE of the span of subspace  $\mathbf{V}^K$  coincides with SIR directions  
(Cook 2007, Szretter and Yohai 2009)

Given current  $A$  and predictor  $X_{j \notin A}$ , we want to test

$H_0: X_j$  is irrelevant, vs.  $H_1: X_j$  is relevant

$$\frac{P_{M_1}(\mathbf{X} | Y)}{P_{M_0}(\mathbf{X} | Y)} = \frac{P_{M_1}(X_j | \mathbf{X}_A, Y)}{P_{M_0}(X_j | \mathbf{X}_A, Y)}$$

# SIR via MLE

Let  $A$  be the set of relevant predictors and  $C = A^c$ ,  $d = |A|$

$$\mathbf{X}_A | Y \in H_s \sim N(\mu_s, \Sigma)$$

$$\mathbf{X}_C | \mathbf{X}_A, Y \in H_s \sim N(\mathbf{X}_A \boldsymbol{\beta}, \Sigma_0)$$

$\mu_s = \alpha + V^K$ , belongs to a  $K$ -dimensional affine space ( $K < d$ )

MLE of the span of subspace  $V^K$  coincides with SIR directions  
(Cook 2007, Szretter and Yohai 2009)

Given current  $A$  and predictor  $X_{j \notin A}$ , we want to test

$H_0: X_j$  is irrelevant, vs.  $H_1: X_j$  is relevant

$$LR_j = \frac{P_{\hat{M}_1}(X_j | \mathbf{X}_A, Y)}{P_{\hat{M}_0}(X_j | \mathbf{X}_A, Y)}$$

# LR Test vs. COP

Given current  $A$ , the likelihood ratio (LR) test statistic of  $H_0: X_j$  is irrelevant, vs.  $H_1: X_j$  is relevant

$$\begin{aligned} 2\text{LR}_j &= -n \left( \sum_{k=1}^K \log(1 - \hat{\lambda}_k^{A+j}) - \sum_{k=1}^K \log(1 - \hat{\lambda}_k^A) \right) \\ &= n \sum_{k=1}^K \log \left( 1 + \frac{\hat{\lambda}_k^{A+j} - \hat{\lambda}_k^A}{1 - \hat{\lambda}_k^{A+j}} \right) \end{aligned}$$

Under  $H_0: X_j$  is irrelevant

$$\text{COP}_k^{A+j} = n \frac{(\hat{\lambda}_k^{A+j} - \hat{\lambda}_k^A)}{1 - \hat{\lambda}_k^{A+j}} \rightarrow_p \chi^2(1), \quad \frac{(\hat{\lambda}_k^{A+j} - \hat{\lambda}_k^A)}{1 - \hat{\lambda}_k^{A+j}} \rightarrow_p 0$$

# LR Test vs. COP

Given current  $A$ , the likelihood ratio (LR) test statistic of  $H_0: X_j$  is irrelevant, vs.  $H_1: X_j$  is relevant

$$\begin{aligned} 2\text{LR}_j &= -n \left( \sum_{k=1}^K \log(1 - \hat{\lambda}_k^{A+j}) - \sum_{k=1}^K \log(1 - \hat{\lambda}_k^A) \right) \\ &= n \sum_{k=1}^K \log \left( 1 + \frac{\hat{\lambda}_k^{A+j} - \hat{\lambda}_k^A}{1 - \hat{\lambda}_k^{A+j}} \right) \end{aligned}$$

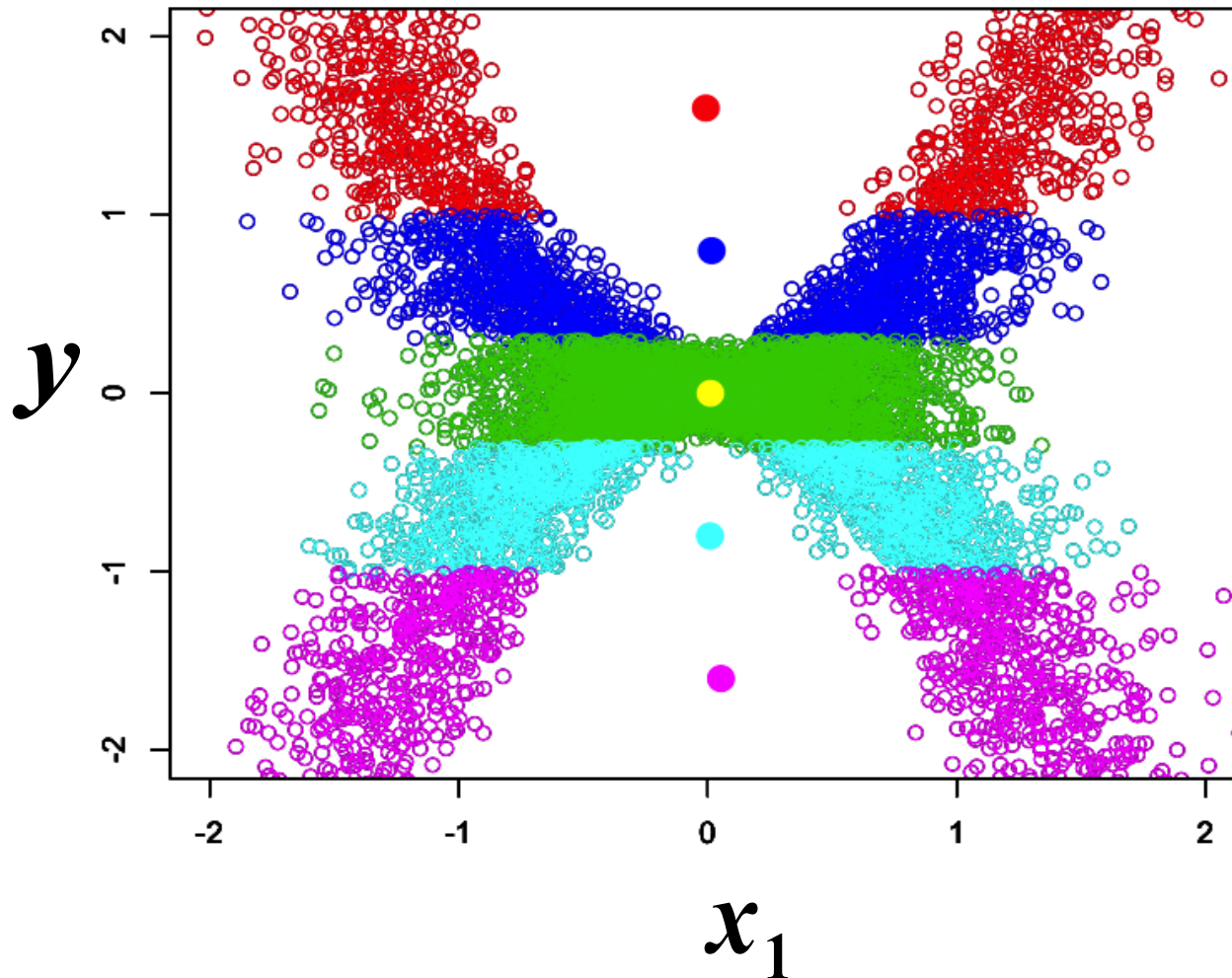
Under  $H_0: X_j$  is irrelevant

$$\text{COP}_k^{A+j} = n \frac{(\hat{\lambda}_k^{A+j} - \hat{\lambda}_k^A)}{1 - \hat{\lambda}_k^{A+j}} \rightarrow_p \chi^2(1), \quad \frac{(\hat{\lambda}_k^{A+j} - \hat{\lambda}_k^A)}{1 - \hat{\lambda}_k^{A+j}} \rightarrow_p 0$$

$$2\text{LR}_j \rightarrow_p \text{COP}_{1:K}^{A+j} = \sum_{k=1}^K \text{COP}_k^{A+j} \rightarrow_p \chi^2(K)$$

# Beyond the First-order

- $E(X_1|Y)=0$



# An Augmented Model

- For  $h = 1, \dots, H$ ,

$$\begin{aligned}\mathbf{X}_{\mathcal{A}}|Y \in S_h &\sim N(\mu_h, \Sigma_h), \\ \mathbf{X}_{\mathcal{A}^c}|\mathbf{X}_{\mathcal{A}}, Y \in S_h &\sim N(\alpha + \beta' \mathbf{X}_{\mathcal{A}}, \Sigma_0).\end{aligned}$$

- If  $\text{Cov}(\mathbf{X})$  is not degenerate, the *minimum* set  $\mathcal{A}$  satisfying the above model is unique.
- Likelihood ratio test for  $H_0 : \mathcal{A} = \mathcal{C}$  v.s.  $H_1^* : \mathcal{A} = \mathcal{C} \cup \{j\}$ ,

$$L_{j|\mathcal{C}}^* = \frac{P_{\mathcal{M}_1^*}(\mathbf{x}|y)}{P_{\mathcal{M}_0}(\mathbf{x}|y)} = \frac{P_{\mathcal{M}_1^*}(x_j|\mathbf{x}_{\mathcal{C}}, y)}{P_{\mathcal{M}_0}(x_j|\mathbf{x}_{\mathcal{C}}, y)}$$



# Likelihood Ratio Test

- Assume  $[\hat{\sigma}_j^{(h)}]^2$  is the estimated variance by regressing  $X_j$  on  $X_C$  in slice  $S_h$  and  $\hat{\sigma}_j^2$  is the overall estimated variance. Then,

$$\hat{D}_{j|C}^* = \frac{2}{n} \log(L_{j|C}^*) = \log \hat{\sigma}_j^2 - \sum_{h=1}^H \frac{n_h}{n} \log [\hat{\sigma}_j^{(h)}]^2$$

- When  $\mathcal{A} \subset \mathcal{C}$  and  $|\mathcal{C}| = d$ ,  $n\hat{D}_{j|C}^* \xrightarrow{d} \chi_{(H-1)(d+2)}^2$

$$\hat{D}_{j|C}^* \sim \log \left( 1 + \frac{Q_0}{\sum_{h=1}^H Q_h} \right) - \sum_{h=1}^H \frac{n_h}{n} \log \left( \frac{Q_h/n_h}{\sum_{h=1}^H Q_h/n} \right)$$

where  $Q_0 \sim \chi_{(H-1)(d+1)}^2$  and  $Q_h \sim \chi_{n_h-d-1}^2$  ( $h = 1, \dots, H$ ) are mutually independent according to *Cochran's theorem*.

- As  $n \rightarrow \infty$ , we have

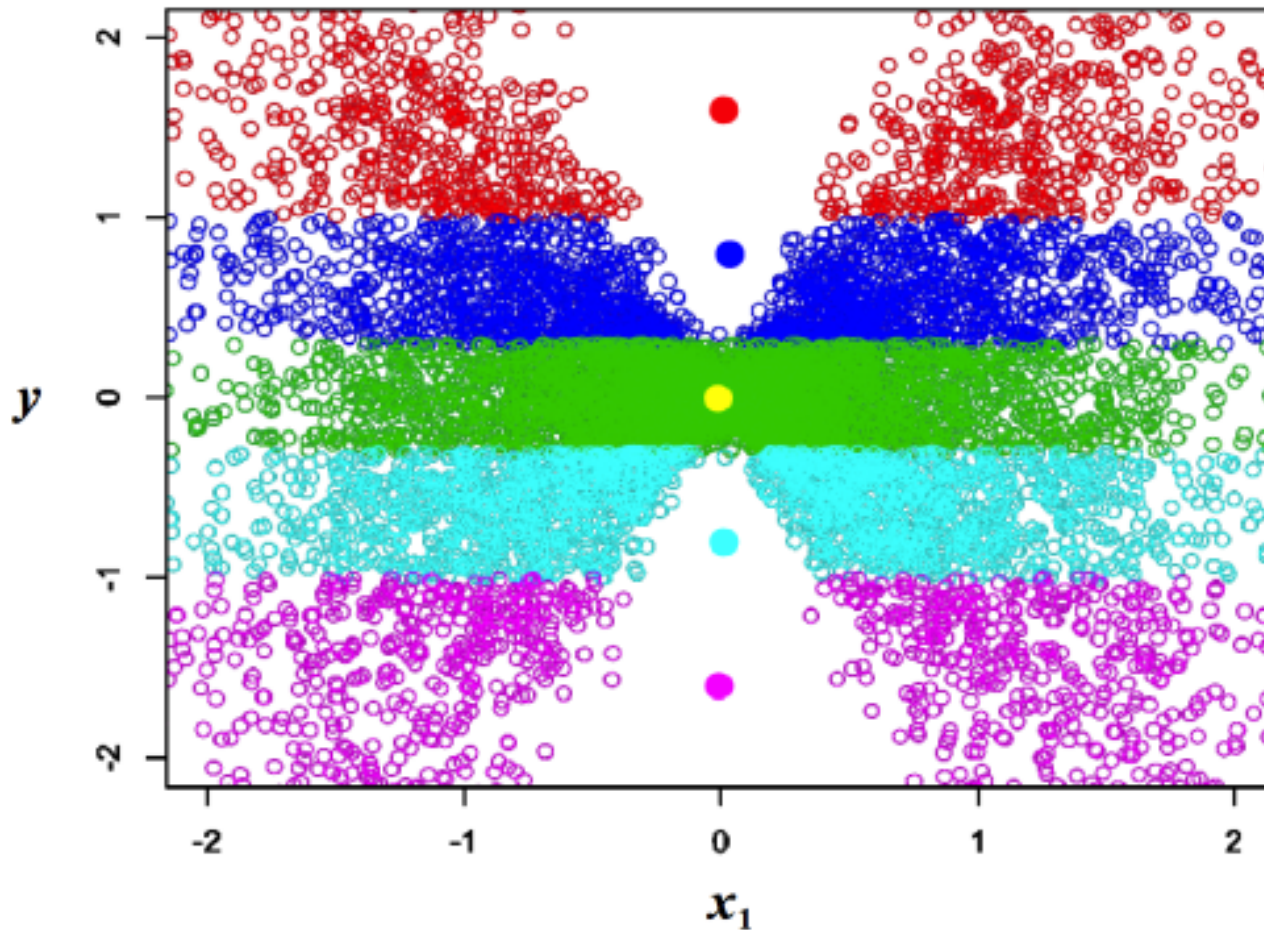
$$\begin{aligned} & \widehat{D}_{j|c}^* \xrightarrow{\text{a.s.}} D_{j|c}^* \\ &= \log \left( 1 + \frac{\text{Var}(M_j) - \text{Cov}(M_j, \mathbf{X}_c) [\text{Var}(\mathbf{X}_c)]^{-1} \text{Cov}(M_j, \mathbf{X}_c)'}{\mathbb{E}(V_j)} \right) \\ &+ \log \mathbb{E}(V_j) - \mathbb{E} \log(V_j) \end{aligned}$$

where  $M_j = \mathbb{E}(X_j | \mathbf{X}_c, \mathcal{S}(Y))$  and  $V_j = \text{Var}(X_j | \mathbf{X}_c, \mathcal{S}(Y))$ .

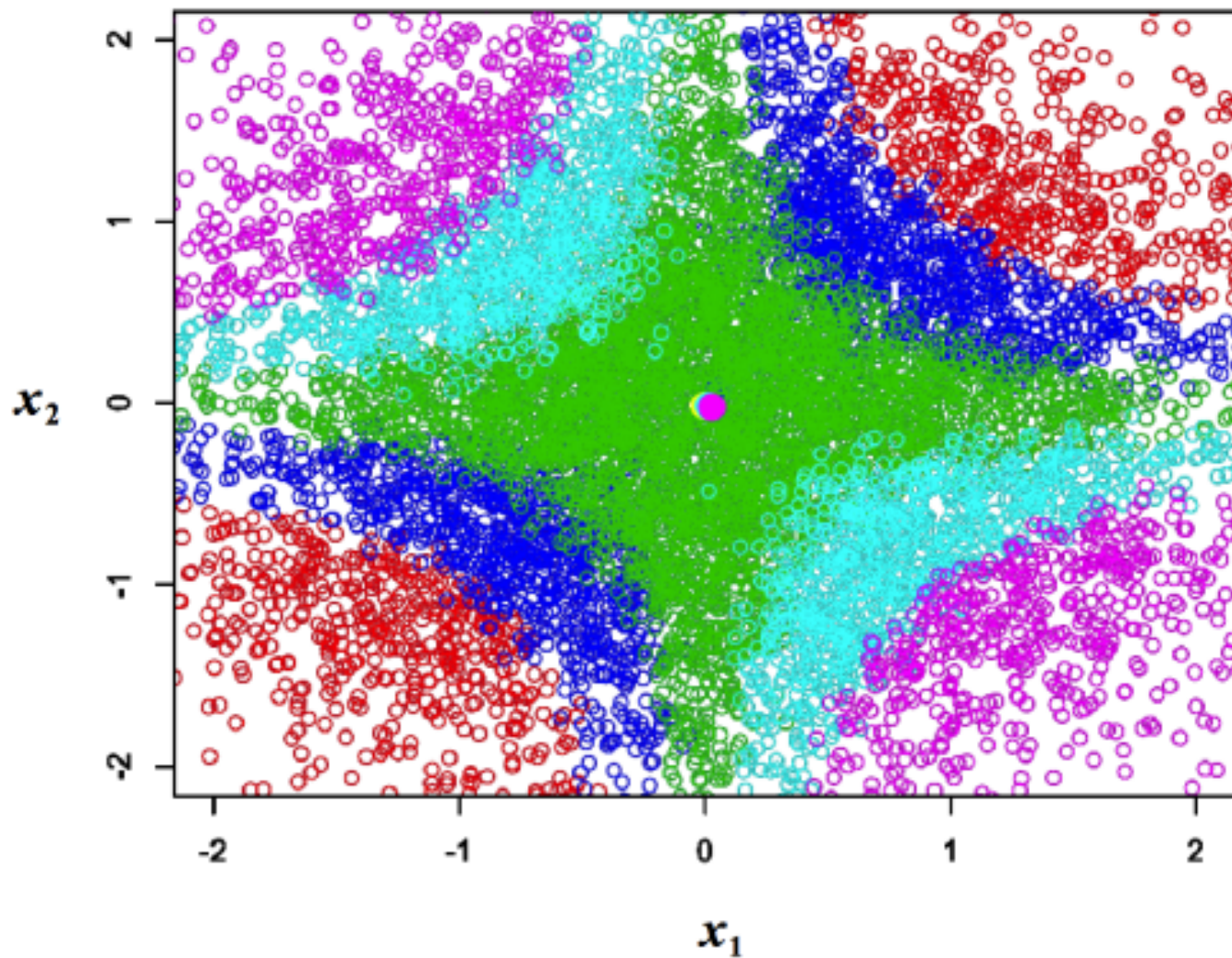
- $\widehat{D}_{j|c}^* \xrightarrow{\text{a.s.}} 0$  if and only if  $\mathbb{E}(X_j | \mathbf{X}_c, \mathcal{S}(Y)) = \mathbb{E}(X_j | \mathbf{X}_c)$  and  $\text{Var}(X_j | \mathbf{X}_c, \mathcal{S}(Y)) = \text{Var}(X_j | \mathbf{X}_c)$
- A stepwise procedure based on  $\widehat{D}_{j|c}^*$  is consistent when  $p = O(n^\gamma)$  with  $\gamma < 1/2$ . ●

# Example revisit

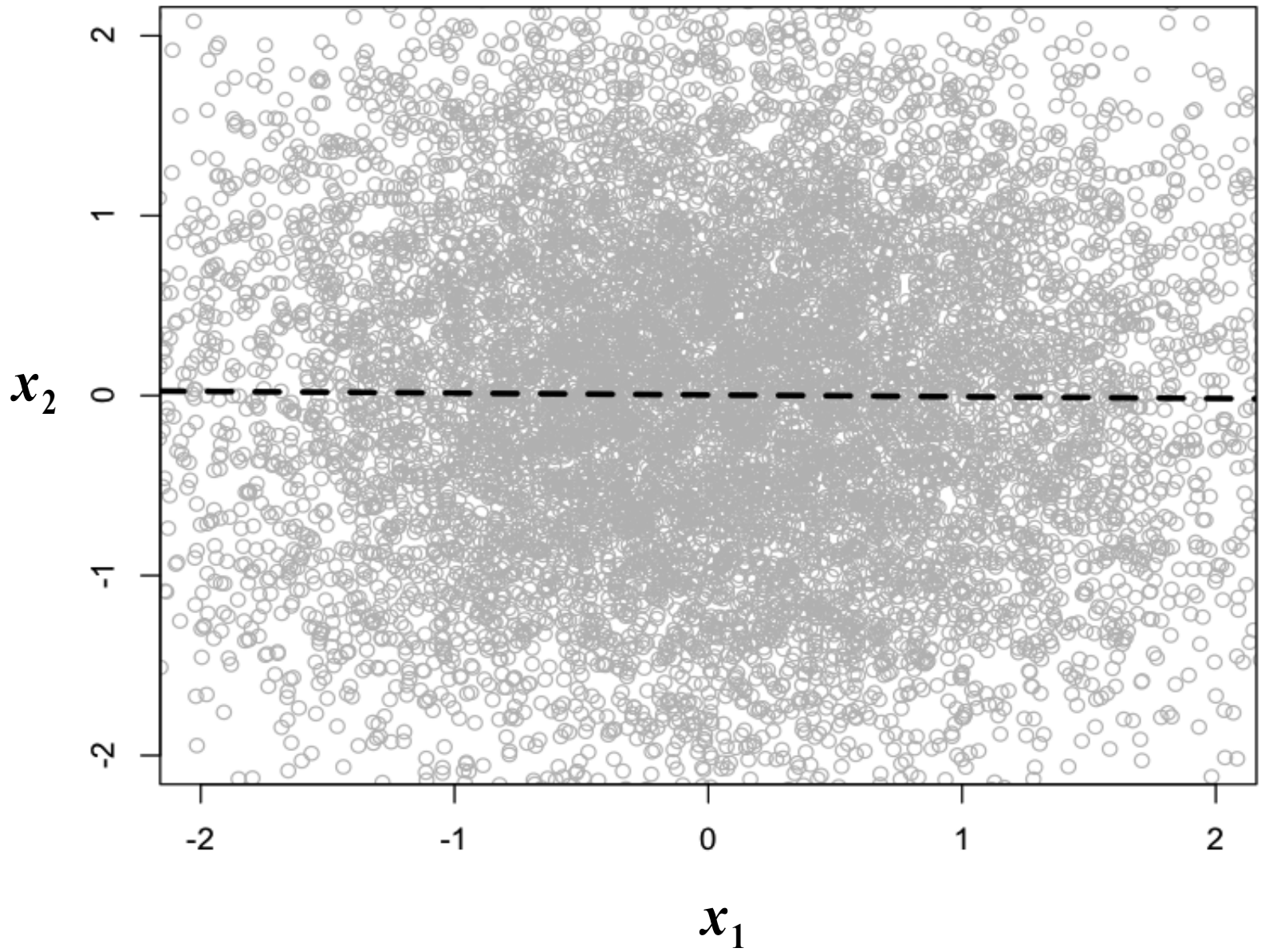
- $\text{Var}(X_1|Y)$  depends on  $Y$ .

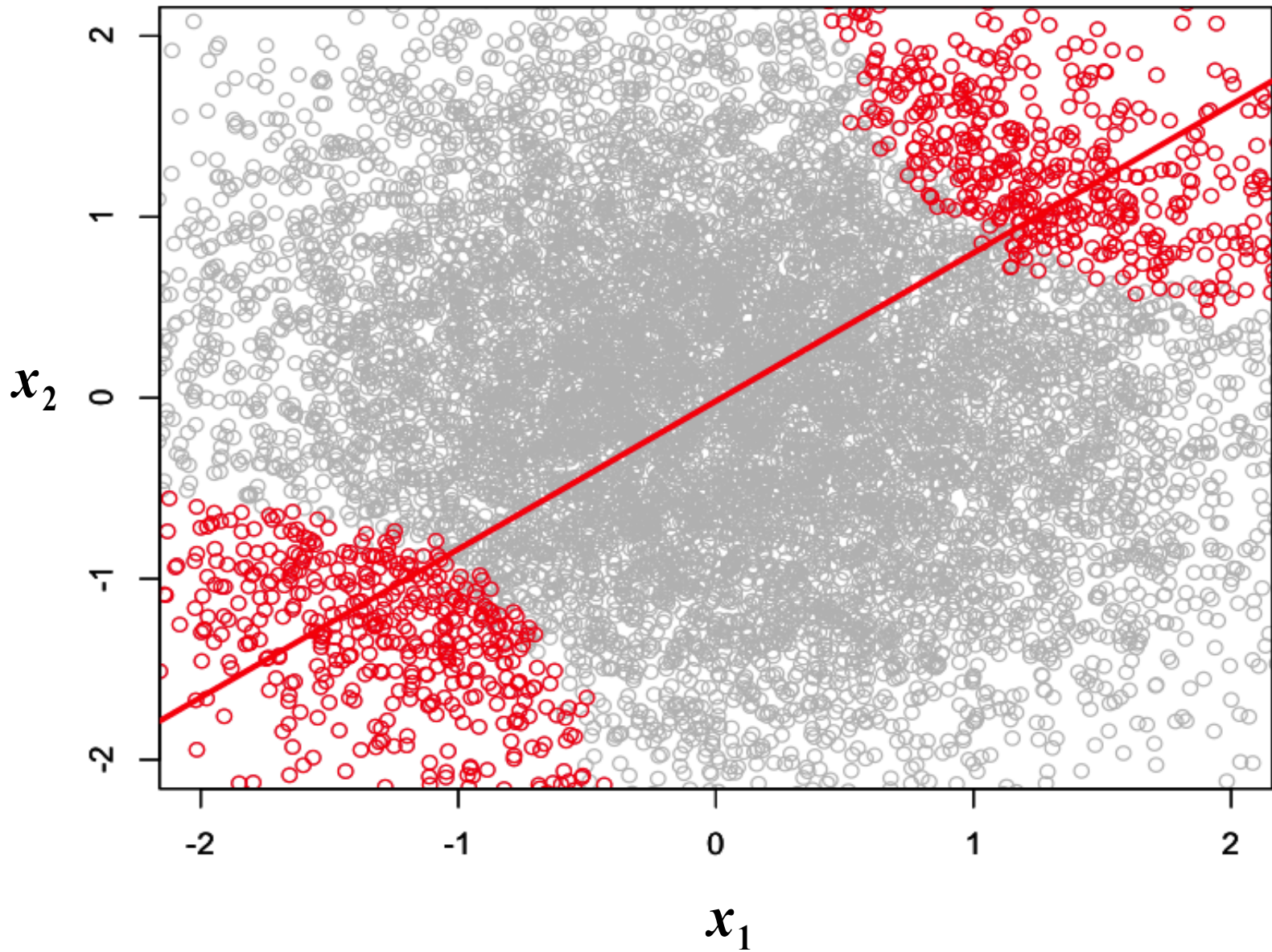


- $\mathbb{E}(X_2|X_1, Y)$  depends on  $Y$ .

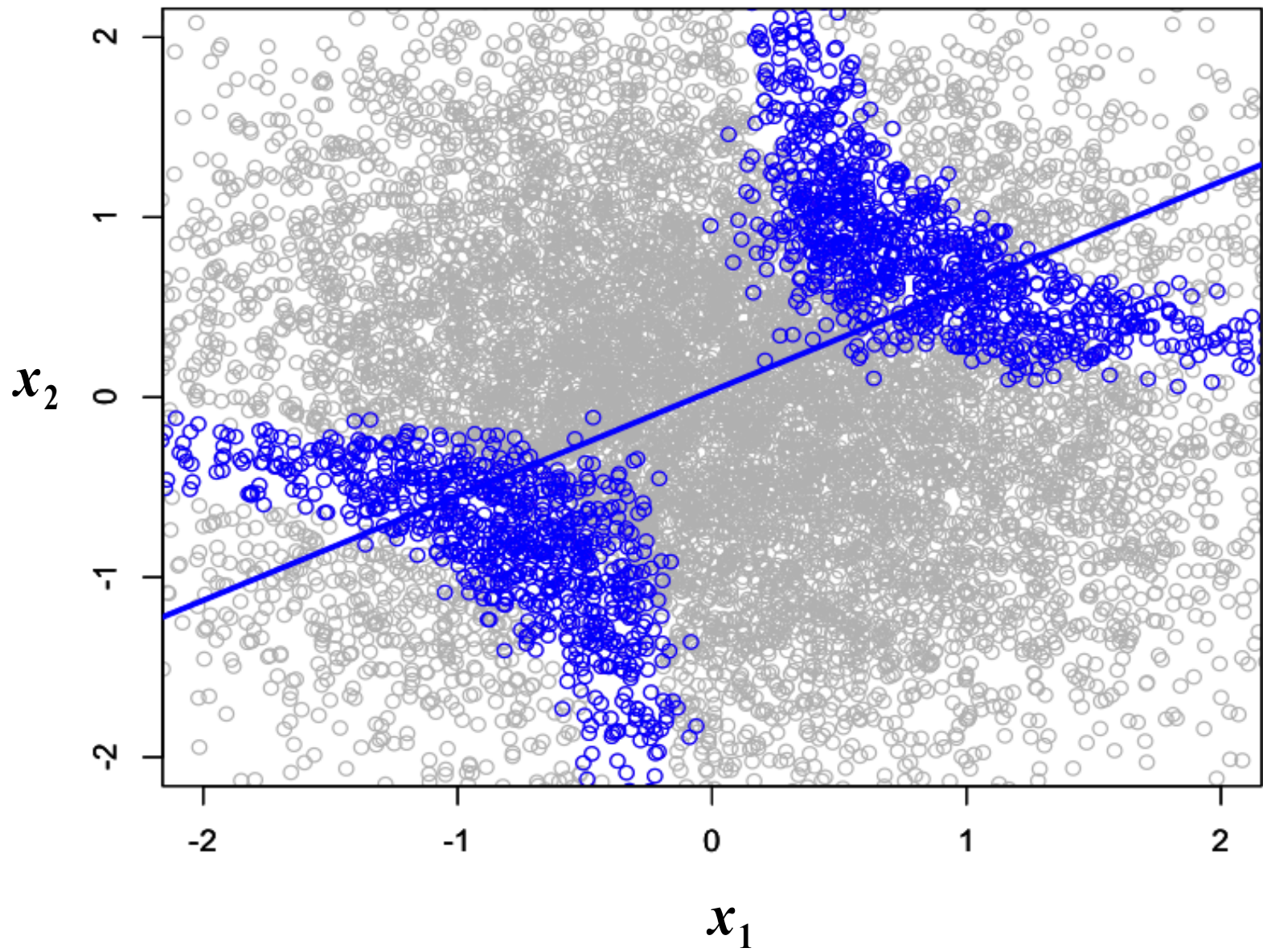


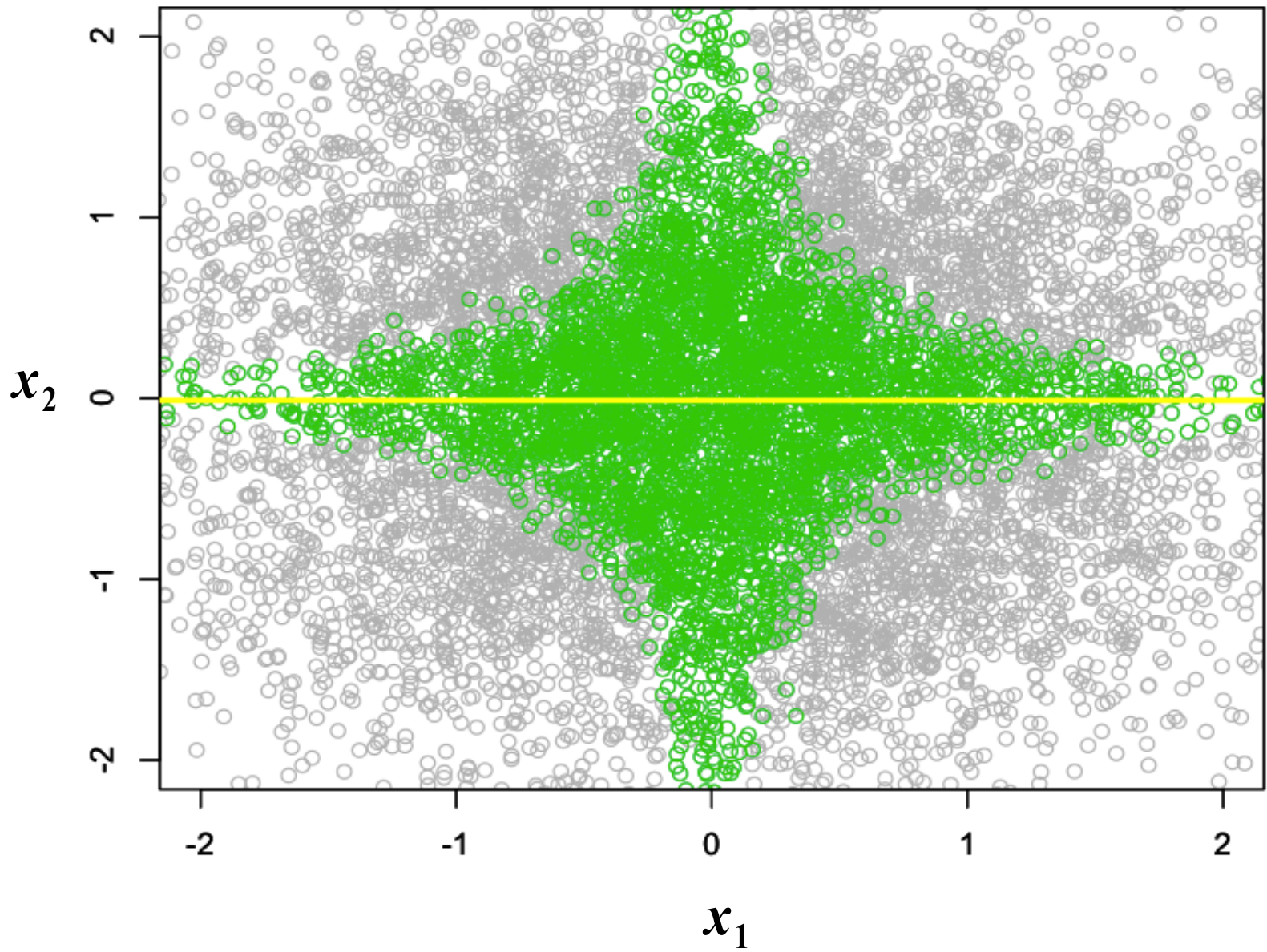




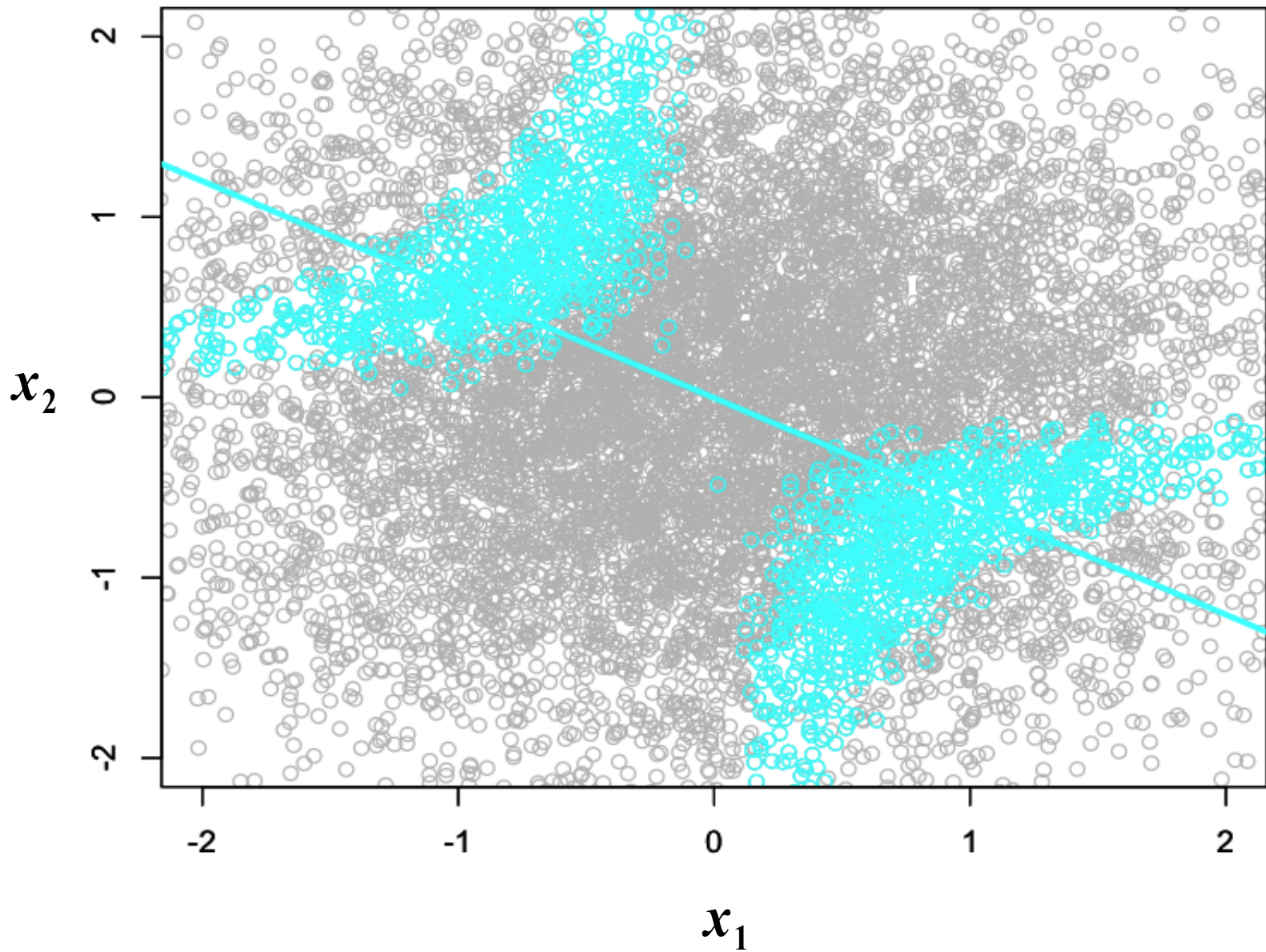


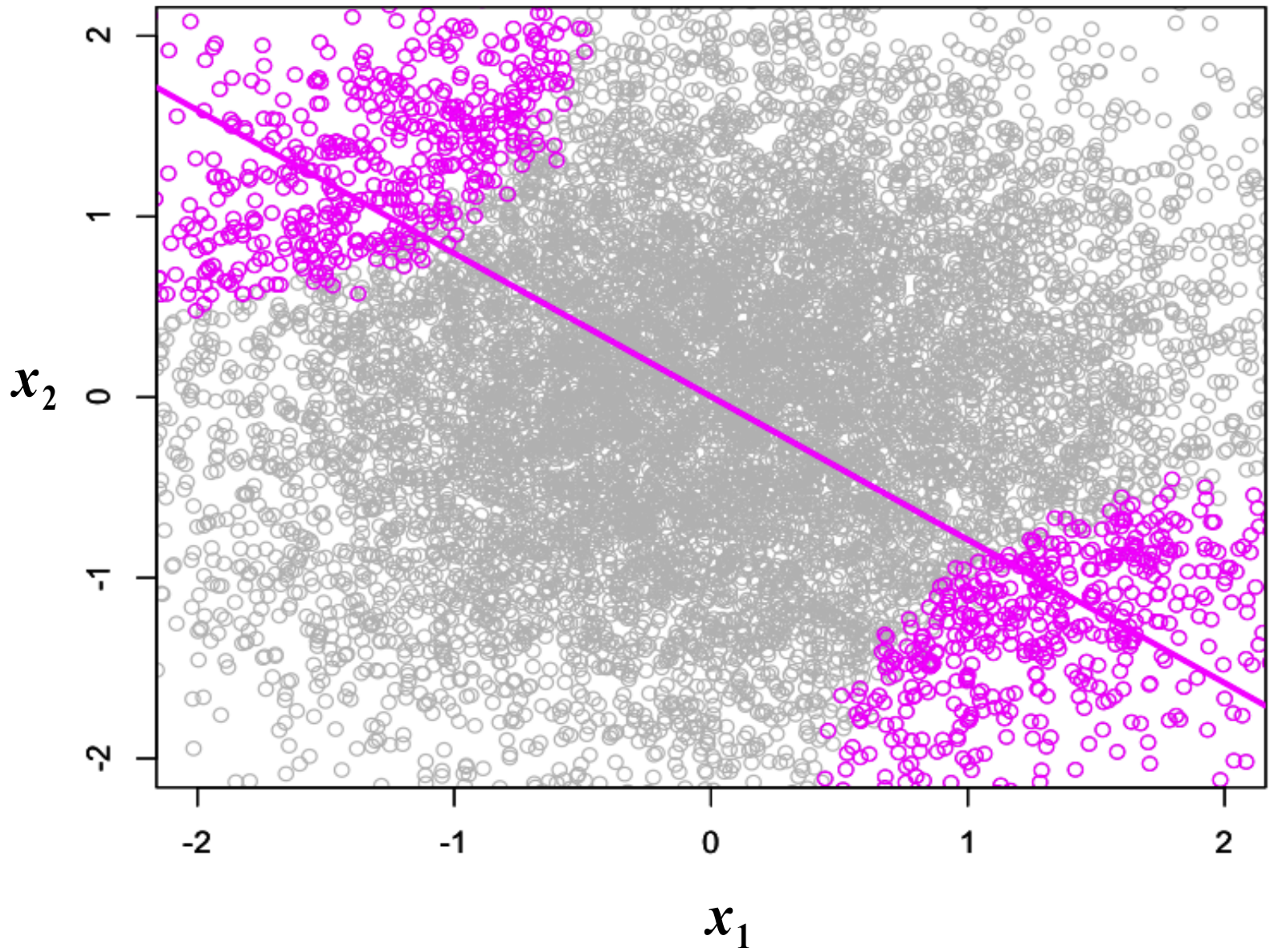




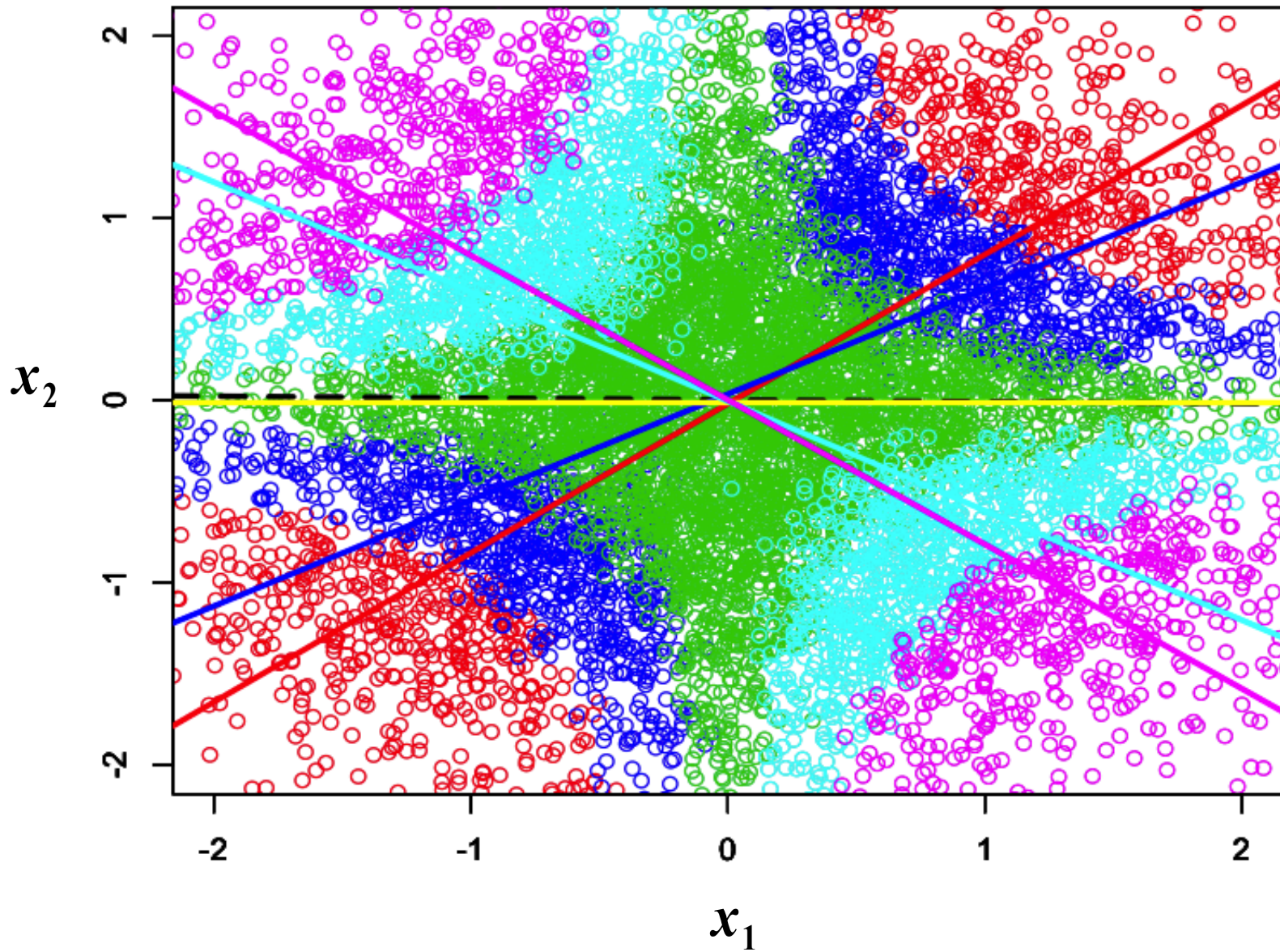












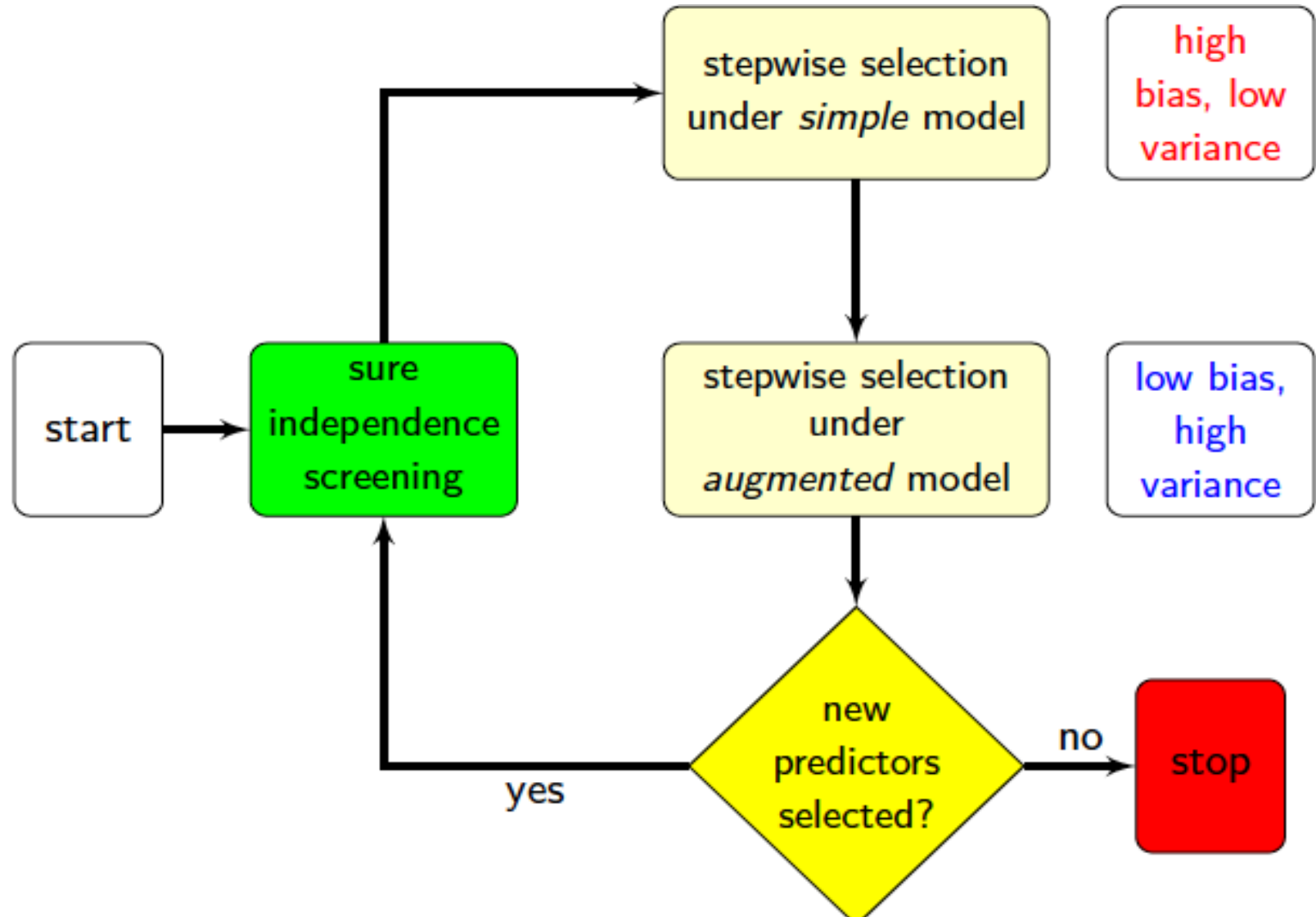
## Sure independence screening (SIS) when $p \gg n$

- Independence screening: first independently selects variables based on their marginal relationships with the response and then applies refined methods in the second step.
- Rank predictors according to  $\hat{D}_{j|C}^*$  with  $C = \emptyset$ :

$$\hat{D}_{j|C=\emptyset}^* \xrightarrow{\text{a.s.}} \log \left( 1 + \frac{\text{Var}(\mathbb{E}(X_j | \mathcal{S}(Y)))}{\mathbb{E}(\text{Var}(X_j | \mathcal{S}(Y)))} \right) + \log \mathbb{E}(\text{Var}(X_j | \mathcal{S}(Y))) - \mathbb{E} \log(\text{Var}(X_j | \mathcal{S}(Y)))$$

- The first  $n - 1$  variables have a high probability to include the true predictors  $\mathcal{A}$  (almost surely under moderate conditions) even when  $\log(p) = O(n^\gamma)$  with  $\gamma < 1$ .

# Siri: An interweaving strategy



# Theoretical Properties

- Under moderate conditions, a stepwise procedure with forward selection and backward elimination is consistent when  $p = O(n^\gamma)$  with  $\gamma < 1/2$ .
- By choosing the threshold appropriately, the addition step will not stop selecting variables until all the true predictors have been included.
- Once all the true predictors have been included, all the redundant variables will be removed from the selected variables.



## Conditions for consistency

- For  $j \in \mathcal{A}$ , we have

$$X_j | \mathbf{X}_{\mathcal{A}-\{j\}}, Y \in S_h \sim N \left( \alpha_j^{(h)} + \beta' \mathbf{X}_{\mathcal{A}-\{j\}}, \sigma_j^2 \right).$$

- Condition 1 (detectability): Let  $\alpha_j(Y) = \sum_h^H \alpha_j^{(h)} \mathbb{I}(Y \in S_h)$ . There exist  $\xi > 0$  and  $\kappa > 0$  such that

$$\text{Var}(\alpha_j(Y)) \geq \xi n^{-\kappa} \text{ for } j \in \mathcal{A}.$$

- Condition 2 (dependency): The eigenvalues of  $\text{Var}(\mathbf{X})$  and  $\text{Var}(\mathbf{X} | Y \in S_h)$  ( $h = 1, \dots, H$ ) have positive lower and upper bounds.
- Condition 3 (dimensionality):  $\lim_{n \rightarrow \infty} (p) = \infty$  and  $p = o(n^\rho)$  with  $\rho > 0$  and  $2\rho + 2\kappa < 1$ .

# Consistency of Stepwise Procedure

- Under Condition 1-3, as  $n \rightarrow \infty$ , there exist constants  $c > 0$  and  $\kappa \geq 0$  such that

$$\Pr \left( \min_{C: C^c \cap \mathcal{A} \neq \emptyset} \max_{j \in C^c} \hat{D}_{j|C} \geq cn^{-\kappa} \right) \rightarrow 1, \text{ and}$$

$$\Pr \left( \max_{C: C^c \cap \mathcal{A} = \emptyset} \max_{j \in C^c} \hat{D}_{j|C} < Cn^{-\kappa} \right) \rightarrow 1 \text{ for any } C > 0.$$

- If we choose  $t_a = cn^{-\kappa}$  and  $t_d = (c/2)n^{-\kappa}$ , then the addition step will not stop selecting variables until all the true predictors have been included.



# Implementation Issues

- Under Condition 1-3, as  $n \rightarrow \infty$ , there exist constants  $c > 0$  and  $\kappa \geq 0$  such that

$$\Pr \left( \min_{C: C^c \cap \mathcal{A} \neq \emptyset} \max_{j \in C^c} \hat{D}_{j|C} \geq cn^{-\kappa} \right) \rightarrow 1, \text{ and}$$

$$\Pr \left( \max_{C: C^c \cap \mathcal{A} = \emptyset} \max_{j \in C^c} \hat{D}_{j|C} < Cn^{-\kappa} \right) \rightarrow 1 \text{ for any } C > 0.$$

- If we choose  $t_a = cn^{-\kappa}$  and  $t_d = (c/2)n^{-\kappa}$ , then the addition step will not stop selecting variables until all the true predictors have been included.
- Once all the true predictors have been included, all the redundant variables will be removed from the selected variables.

# Simulation I (linear)

$$Y = X_p \beta + \epsilon, \epsilon \sim N(0,1), \text{Cov}(X_i, X_j) = 0.5^{|i-j|}$$

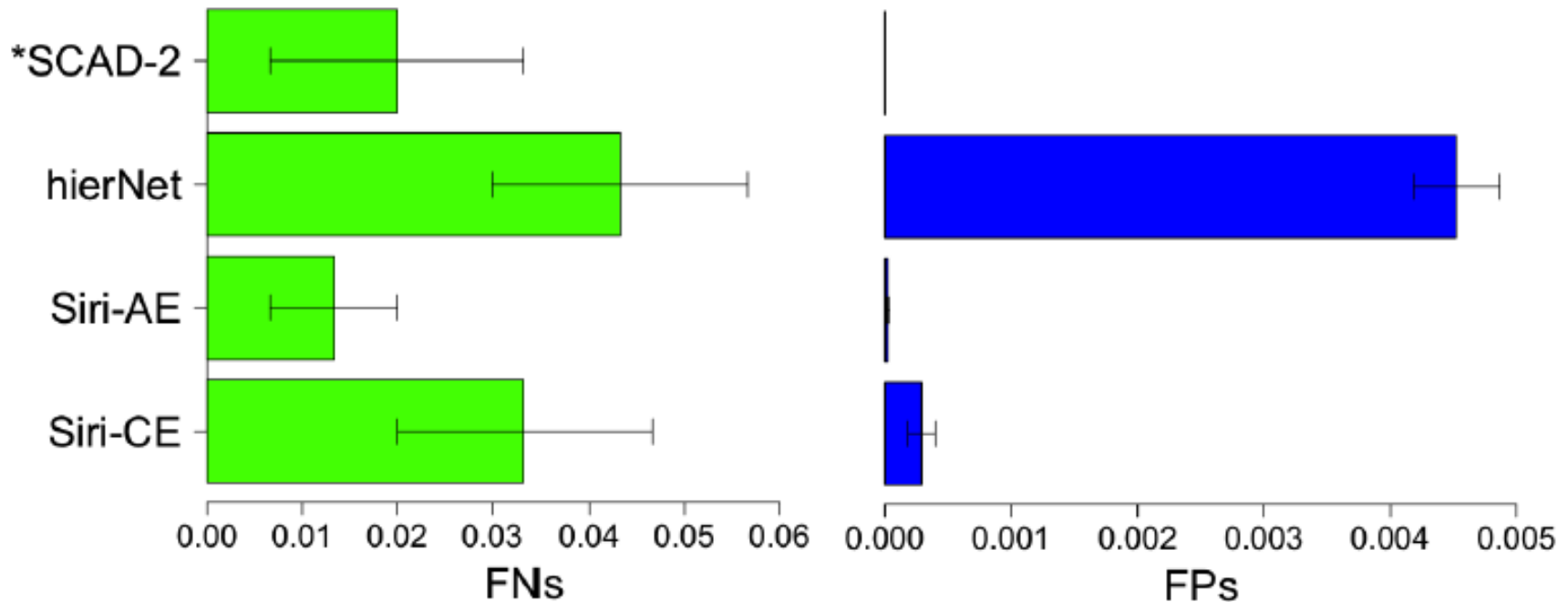
$$n = 200, p = 1000, \beta = (3, 1.5, 1, 1, 2, 1, 0.9, 1, 1, 1, 0, \dots, 0)^T$$

Method	FP(0, 990)	FN(0, 10)
<b>SIRI-C</b> [CV minimizing classification error]	<b>1.86 (0.222)</b>	<b>1.66 (0.117)</b>
<b>SIRI-M</b> [CV minimizing mean square error]	<b>0.76 (0.120)</b>	<b>1.75 (0.114)</b>
<b>COP</b>	<b>1.62 (0.165)</b>	<b>1.67 (0.118)</b>
<b>SIS-SCAD</b>	<b>0.10 (0.030)</b>	<b>0.64 (0.069)</b>
<b>LASSO</b>	<b>5.40 (0.188)</b>	<b>0.00 (0.000)</b>

# Simulation II: hierarchical interactions

- $\mathbf{X} \sim \text{MVN}(\mathbf{0}, \mathbb{I}_p)$  with  $n = 200$ ,  $p = 1000$ , and

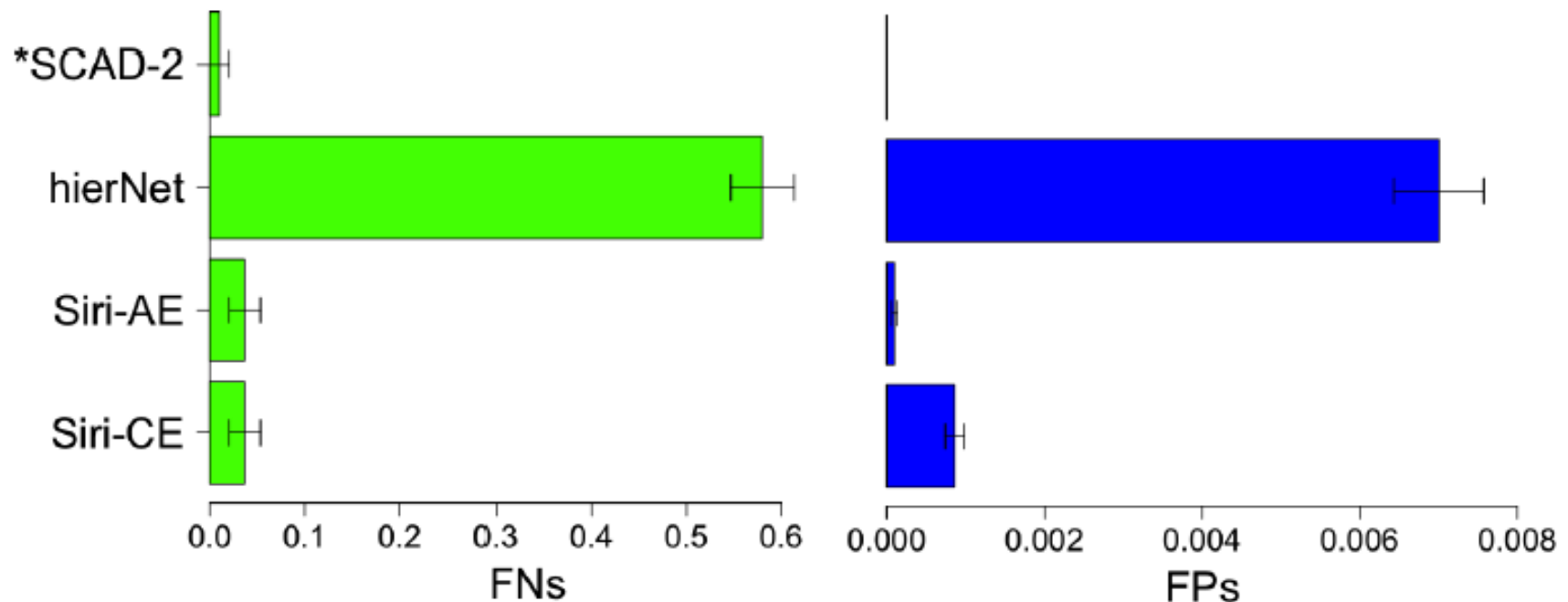
$$Y = X_1 + X_1X_2 + X_1X_3 + 0.2\epsilon$$



# Simulation III: non-hierarchical

- $\mathbf{X} \sim \text{MVN}(\mathbf{0}, \mathbb{I}_p)$  with  $n = 200$ ,  $p = 1000$ , and

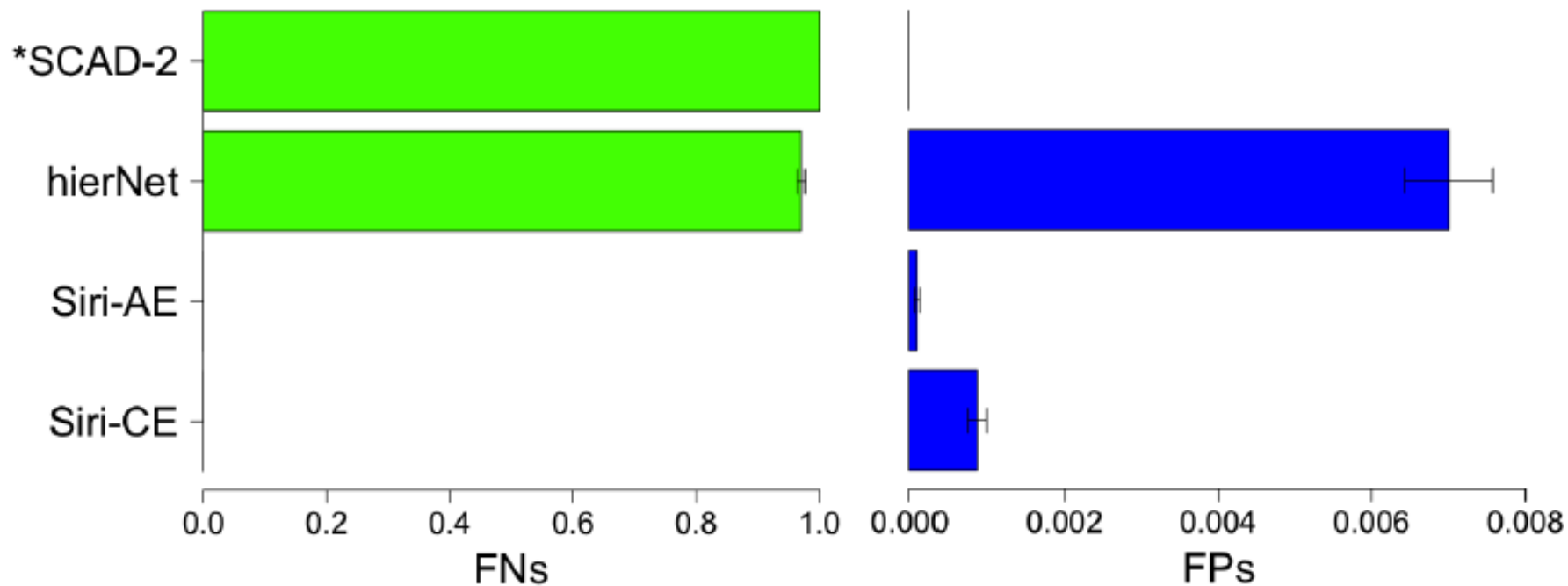
$$Y = X_1X_2 + X_1X_3 + 0.2\epsilon$$



# Simulation IV: Non-multiplicative

- $\mathbf{X} \sim \text{MVN}(\mathbf{0}, \mathbb{I}_p)$  with  $n = 200$ ,  $p = 1000$ , and

$$Y = \frac{X_1}{X_2 + X_3} + 0.2\epsilon$$



# Simulation V (heteroscedastic, single index)

$$Y = \frac{0.2\epsilon}{1.5 + \sum_{j=1}^8 X_j}, \quad X_p \sim \text{independent normal}$$

$n = 1000, p = 1000$

Method	FP(0, 992)	FN(0, 8)
SIRI-C	2.00 (0.163)	<b>0.42 (0.138)</b>
SIRI-M	<b>0.43 (0.079)</b>	4.60 (0.274)
COP	1.26 (0.128)	3.32 (0.192)
SIS-SCAD	3.23 (0.356)	8.00 (0.000)
LASSO	0.64 (0.255)	8.00 (0.000)

# Simulation VI (hub with linear effect)

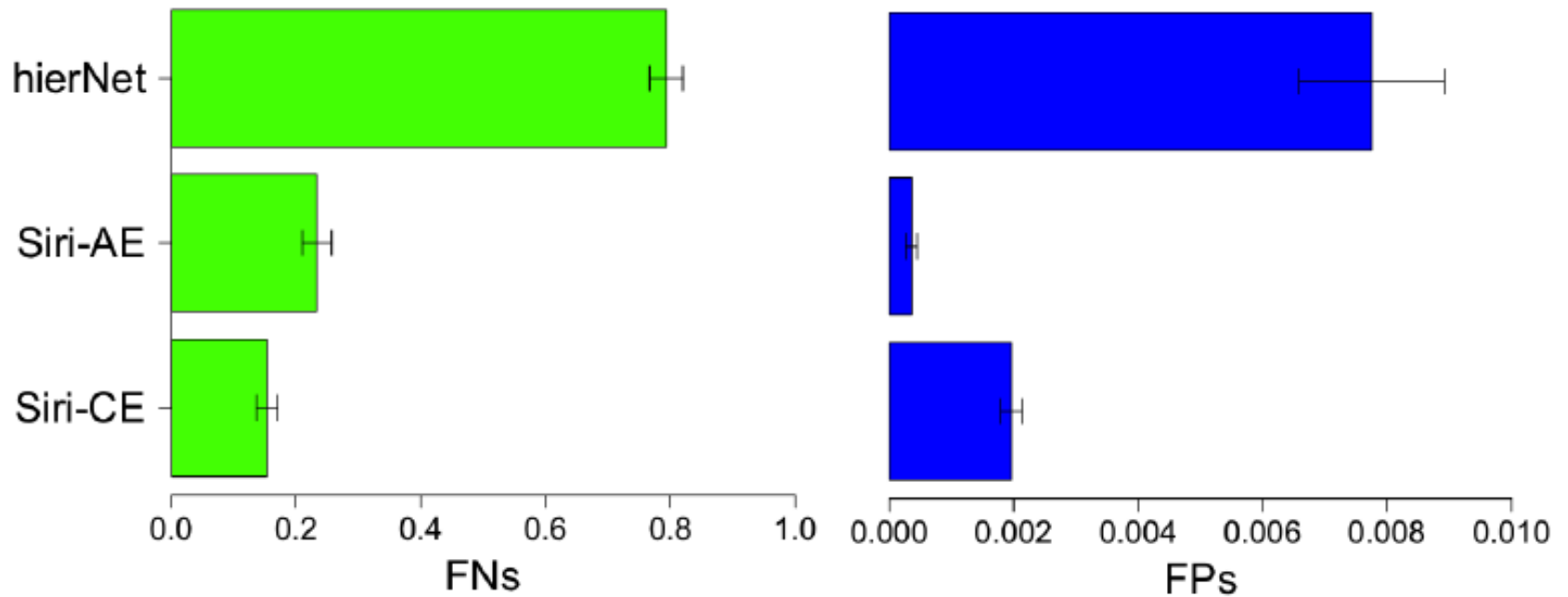
$Y = X_1 + X_1 \times (X_2 + X_3) + 0.2\epsilon$ ,  $X_p \sim$  independent normal  
 $n = 200$ ,  $p = 1000$

Method	FP(0, 997)	FN(0, 3)
SIRI-C	0.39 (0.115)	0.12 (0.046)
SIRI-M	0.03 (0.017)	<b>0.04 (0.020)</b>
SIS-SCAD-2	<b>0.00 (0.000)</b>	0.45 (0.068)

# Simulation VII (three-way interaction)

$$Y = X_1 \times X_2 \times X_3 + 0.2\epsilon, \quad X_p \sim \text{independent normal}$$

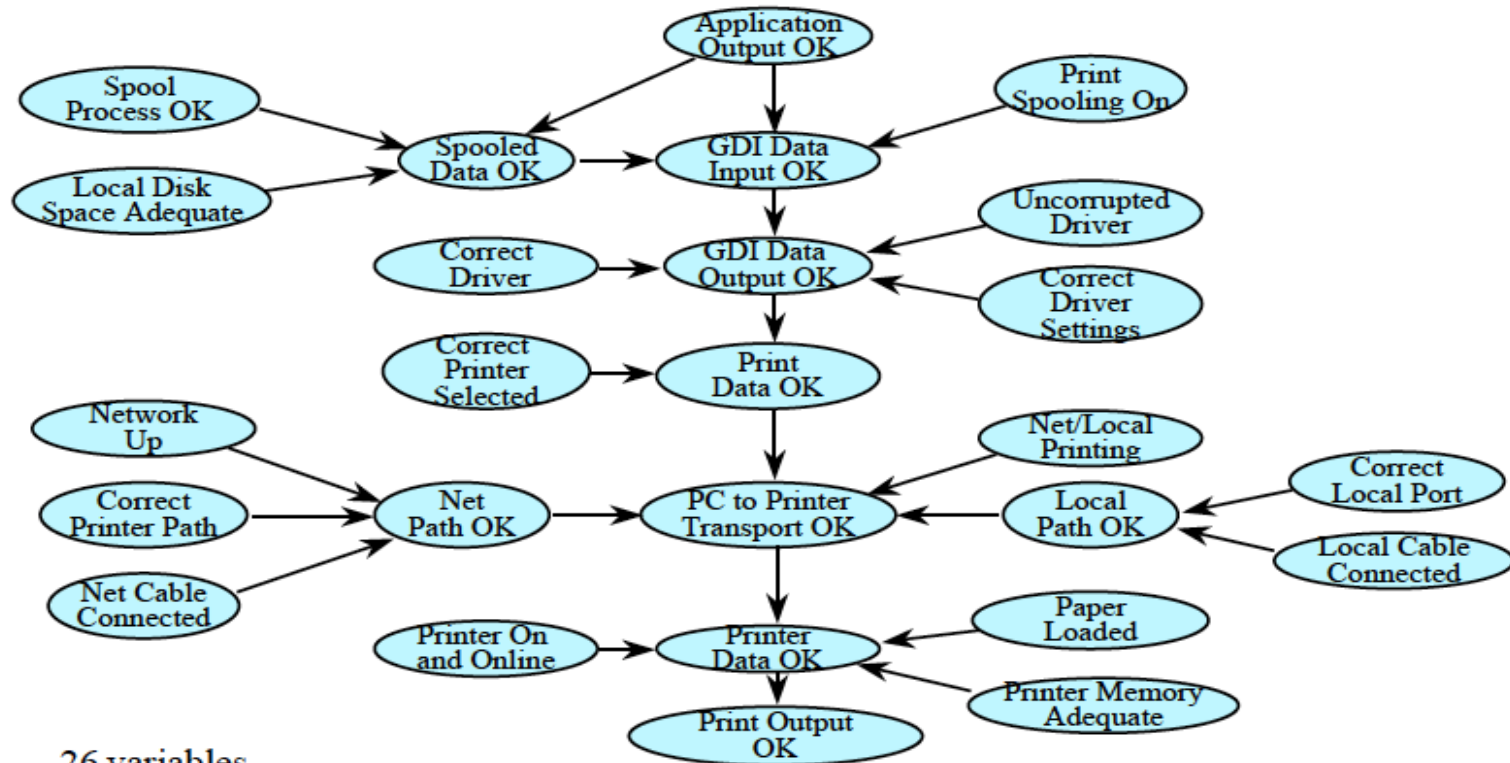
$n = 500, p = 1000$





# Bayesian Networks

## Example: Printer Troubleshooting

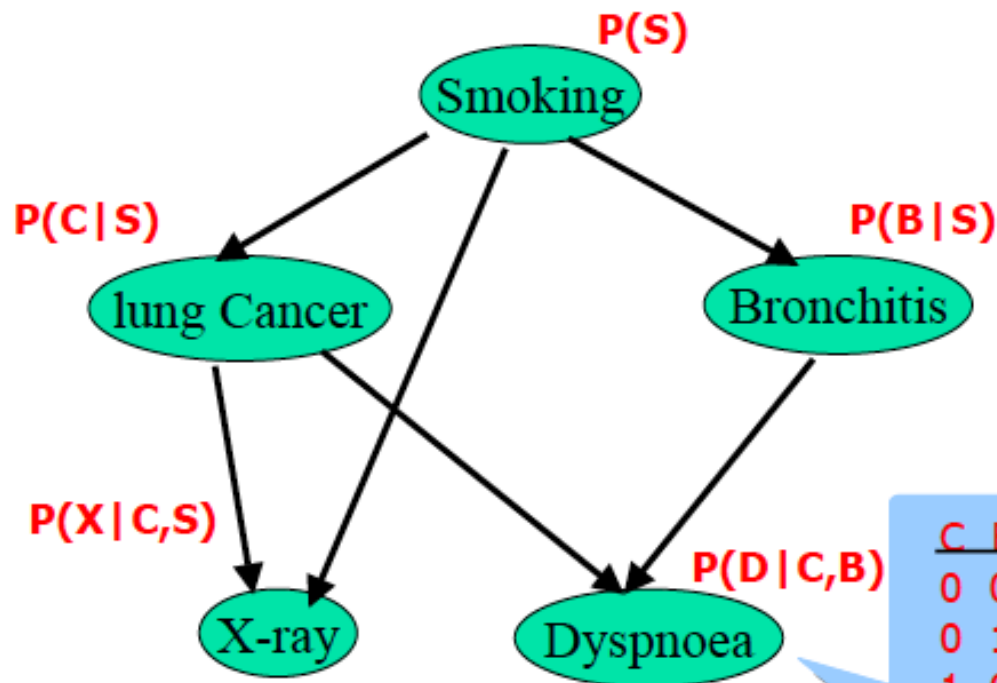


26 variables

Instead of  $2^{26}$  parameters we get

$$99 = 17 \times 1 + 1 \times 2^1 + 2 \times 2^2 + 3 \times 2^3 + 3 \times 2^4$$

# Bayesian Network: $\text{BN} = (\mathbf{G}, \Theta)$



$\mathbf{G}$  - directed acyclic graph (DAG)  
nodes – random variables  
edges – direct dependencies

$\Theta$  - set of parameters in all conditional probability distributions (CPDs)

CPD:

C	B	D=0	D=1
0	0	0.1	0.9
0	1	0.7	0.3
1	0	0.8	0.2
1	1	0.9	0.1

**CPD of node X:**  
 $P(X | \text{parents}(X))$

Compact representation of joint distribution in a **product form** (chain rule):

$$P(S, C, B, X, D) = P(S) P(C|S) P(B|S) P(X|C,S) P(D|C,B)$$

$1 + 2 + 2 + 4 + 4 = 13$  parameters instead of  $2^5 = 32$

# Learning BN structures

- Global approach (Score/likelihood based):

- Posterior inference:

$$P(G | Data) \propto \int P(Data | \theta_G, G) p(\theta_G | G) p(G) d\theta_G$$

- Or score-based criterion

- $AIC = -2 \log P(Data | \hat{\theta}_G, G) + 2 p_G$
- $BIC = -2 \log P(Data | \hat{\theta}_G, G) + \log(n) p_G$

# Learning structures

- Local approaches: using conditional independence statements as constraints.
    - Represented by “Inductive causality” (IC) algorithm due to Pearl (2000).
1. First, the skeleton of the network (undirected graph underlying the network structure) is learned by recursively testing the conditional independence between nodes.
  2. Set all direction of the arcs that are part of a v-structure, which is a triplet of nodes incident on a converging connection  $X_j \rightarrow X_i \leftarrow X_k$
  3. Set the directions of the order arcs as needed to satisfy the acyclicity constraint.

# Finding Markov blanket for each node using Growth-Shrink (GS) algorithm

- It is like a stepwise regression. For each node  $X_i$ , we treat it as the response variable and
  - (a) gradually add variables that are predictive of  $X_i$ ;
  - (b) Backward removing those “redundant”  $X_j$ 's obtained from the growth phase.

# Discussion

- Cross-validation to select the dimension and thresholds
- Back to full Bayesian model with dynamic slicing
  - We want to have flexibility in choosing slicing boundaries
  - Connection with Mutual-Information Criterion (MIC)
  - Many interesting possibilities
- Robustness to the distribution of predictors

# Acknowledgment

Bo Jiang – who did all the work

Dr. Tingting Zhang, Dr. Wenxuan Zhong

Joseph K. Blitzstein