

# **Introduction of linearly independent descriptor generation (LIDG) method for interpretable modeling**

Hitoshi Fujii

# **Outline**

- **Materials informatics (MI)**
- **Introduction of linear regression analysis**
- **Ordinary least square method**
- **Multicollinearity (MCL)**
- **Examples**
- **Summary**

# Materials informatics (MI)

The goals of MI are

- (i) to find new materials which have good properties (**material search**)
- (ii) to predict the physical properties of unknown materials (**property prediction**)
- (iii) to elucidate a mechanism of physical properties (**empirical law discovery**)
- (iv) to construct a useful database

from material data with machine learning technique.

Optimization problem:

Evolutional algorithm, Monte Carlo tree search, Bayes optimization,  
Simulated annealing, VS,

Non linear regression:

Deep learning, Support vector machine, Random Forest

Sparse and interpretable modeling (linear regression):

LASSO, LIDG

# Introduction of linear regression analysis

## 1. Assumption of a linear regression model

$$\vec{y} = X\vec{\beta} + \vec{\varepsilon} \quad \vec{\varepsilon}: \text{Error (Gaussian)}$$

$\vec{y}$ : Target value

$X = [\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n]$ : Descriptor matrix

$\vec{\beta}$  : Regression coefficients

	Target property	Other properties (features/attributes/ descriptors)			
	$\vec{y}$	$\vec{x}_1$	$\vec{x}_2$	...	$\vec{x}_n$
Sample 1	$y_1$	$x_{11}$	$x_{12}$	...	$x_{1n}$
Sample 2	$y_2$	$x_{21}$	$x_{22}$	...	$x_{2n}$
...	...	...	...	...	...
Sample $m$	$y_m$	$x_{m1}$	$x_{m2}$	...	$x_{mn}$

$X (m \times n)$ :  
descriptor matrix

## 2. Find $\hat{\vec{\beta}}$ under a certain condition

- Ordinary least squares (OLS) method

$$\hat{\vec{\beta}}_{OLS} = \underset{\vec{\beta}}{\operatorname{argmin}} \left\{ \|\vec{y} - X\vec{\beta}\|_{l_2}^2 \right\}$$

- LASSO

$$\hat{\vec{\beta}}_{LASSO} = \underset{\vec{\beta}}{\operatorname{argmin}} \left\{ \|\vec{y} - X\vec{\beta}\|_{l_2}^2 + \lambda \|\vec{\beta}\|_{l_1}^1 \right\}$$

Definition of  $l_p$  norm

$$\|\vec{\beta}\|_{l_p} = \left( \sum_{i=1}^n |\beta_i|^p \right)^{1/p}$$

# Descriptor space and descriptor vectors

**Design matrix**  
**Descriptor space**  
 $(m \times n + 1)$

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & \boxed{x_{1j}} & \cdots & x_{1n} \\ 1 & x_{21} & x_{22} & \cdots & \boxed{x_{2j}} & \cdots & x_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \boxed{1} & x_{i1} & x_{i2} & \cdots & \boxed{x_{ij}} & \cdots & x_{in} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 1 & x_{m1} & x_{m2} & \cdots & \boxed{x_{mj}} & \cdots & x_{mn} \end{bmatrix} \vec{x}_i,$$

**Usual definition**

$$\vec{x}_i = [1, x_{i1}, x_{i2}, \dots, x_{in}]^T$$

Vectors in  $n + 1$  dimensional feature space  
Feature pattern of  $i$ -th sample

$$X = [\vec{x}_1, \vec{x}_2, \dots, \vec{x}_m]^T$$

$$y_i = \vec{x}_i^T \vec{\beta} + \varepsilon_i \quad \vec{y} = X \vec{\beta} + \vec{\varepsilon}$$

$$(y_i = \vec{\beta}^T \vec{x}_i + \varepsilon_i)$$

**$n$  is usually fixed**  
 **$m$  can be increased**

**In our study**

$$\vec{x}_j = [x_{1j}, x_{2j}, \dots, x_{mj}]^T$$

Vectors in  $m$  dimensional sample space  
 $j$ -th descriptor vector

$$X = [\vec{1}, \vec{x}_1, \vec{x}_2, \dots, \vec{x}_n]$$

$$\vec{y} = \sum_{j=0}^n \vec{x}_j \beta_j + \vec{\varepsilon} \quad \vec{y} = X \vec{\beta} + \vec{\varepsilon}$$

**$n$  can be increased**  
 **$m$  is usually fixed**  
for the extraction of physical law

# Simple model vs. complex model

**Simple model:**  
**(linear and sparse model)**

- low variance
- **high bias (low expression capability)**
- **physically interpretable**

**Complex model:**  
**(Non-linear model, higher order model)**

- **high variance (overfit)**
- **low bias (high expression capability)**
- **almost black box model**

**How do we solve the dilemma  
of expression capability and interpretation in linear regression ?  
→ descriptor generation and descriptor selection**

**Linear sparse model**

**descriptor generation**

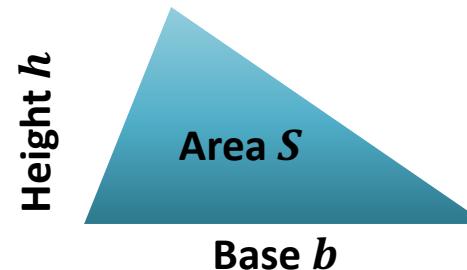
**descriptor selection**

# Importance of multiplication of descriptors

20 triangle data

Target variable	Descriptors	
Area	Base	Height
0.0001	0.2658	0.0007
0.0092	0.9789	0.0187
0.0096	0.0488	0.3931
0.0102	0.3953	0.0514
0.0195	0.0883	0.4416
0.0273	0.2153	0.2536
0.0389	0.3864	0.2012
0.0441	0.1095	0.8048
0.0562	0.8760	0.1283
0.0577	0.2917	0.3958
0.0602	0.4255	0.2832
0.0626	0.7260	0.1725
0.0788	0.1785	0.8830
0.0871	0.7943	0.2192
0.0915	0.3709	0.4936
0.0964	0.6477	0.2976
0.1051	0.5049	0.4163
0.1119	0.7204	0.3108
0.1198	0.7600	0.3154
0.1341	0.5543	0.4839

The area of triangle



Regression result by using **base** and **height** as descriptors

$$S = 0.11681 \mathbf{b} + 0.13780 \mathbf{h} - 0.03876$$

$$\text{Adj. } R^2 = 0.5515$$

Regression result by using **base \* height** as new descriptor

$$S = 0.4998 \mathbf{bh} - 2.333 \times 10^{-5}$$

$$\text{Adj. } R^2 = 1.0000$$

Although the number of descriptors is one less,  
the estimation accuracy is improved.

We can rediscover the triangle area law from this model

$$S = \frac{1}{2} bh$$

# Descriptor generation 1 (basic operations)

Basic descriptors:

$$\vec{x}_1, \vec{x}_2$$

Basic operations:

$$|\vec{x}_1|, |\vec{x}_2|,$$

$$\vec{x}_1^{-1}, \vec{x}_2^{-1},$$

$$|\vec{x}_1|^{-1}, |\vec{x}_2|^{-1},$$

$$|\vec{x}_1 + \vec{x}_2|, |\vec{x}_1 - \vec{x}_2|,$$

$$(\vec{x}_1 + \vec{x}_2)^{-1}, (\vec{x}_1 - \vec{x}_2)^{-1},$$

$$|\vec{x}_1 + \vec{x}_2|^{-1}, |\vec{x}_1 - \vec{x}_2|^{-1},$$

$$|\vec{x}_i| \equiv (|x_{1i}|, |x_{2i}|, \dots, |x_{mi}|)^T$$

$$\vec{x}_i^{-1} \equiv (x_{1i}^{-1}, x_{2i}^{-1}, \dots, x_{mi}^{-1})^T$$

$$|\vec{x}_i|^{-1} \equiv (|x_{1i}|^{-1}, |x_{2i}|^{-1}, \dots, |x_{mi}|^{-1})^T$$

$$|\vec{x}_i + \vec{x}_j| \equiv (|x_{1i} + x_{1j}|, \dots, |x_{mi} + x_{mj}|)^T$$

Of course, we can take more complex operations (such as, root, exponent, logarithm,...).  
From the view point of interpretability, however, such complex descriptors should not be used.

# Descriptor generation 2 (direct product)

Generate descriptors by multiplications (direct product) of descriptors

Ex. three base features case

( $q$  mean descriptors, like  $\vec{x}_i, \vec{x}_j, 1/\vec{x}_i, |\vec{x}_i|, \dots$ )

<b>1<sup>st</sup> order:</b>	$q_0, q_1, q_2$	$( {}_3H_1 = 3 )$
<b>2<sup>nd</sup> order:</b>	$q_0^2, q_1^2, q_2^2,$ $q_0q_1, q_0q_2, q_1q_2$	$q_iq_j \equiv (q_{1i}q_{1j}, q_{2i}q_{2j}, \dots, q_{mi}q_{mj})^T$ $( {}_3H_2 = 6 )$
<b>3<sup>rd</sup> order:</b>	$q_0^3, q_1^3, q_2^3,$ $q_0^2q_1, q_0^2q_2, q_1^2q_2,$ $q_0q_1^2, q_0q_2^2, q_1q_2^2,$ $q_0q_1q_2$	$( {}_3H_3 = 10 )$
<b>4<sup>th</sup> order:</b> ...		$( {}_3H_4 = 15 )$
<b>5<sup>th</sup> order:</b> ...		$( {}_3H_5 = 21 )$

Usually, these generated descriptors are correlated each other.  
We want to find and remove these correlation.

# Multicollinearity

**Collinearity**

$$\vec{x}_i = c_j \vec{x}_j + c_0$$

**Near collinearity**

$$\vec{x}_i \sim c_j \vec{x}_j + c_0$$

**Multicollinearity**

$$\vec{x}_i = c_j \vec{x}_j + c_k \vec{x}_k + \dots + c_0$$

**Near multicollinearity**

$$\vec{x}_i \sim c_j \vec{x}_j + c_k \vec{x}_k + \dots + c_0$$

**Why do we have to detect and remove these multicollinearity ?**

**Problems of the multicollinearity in linear regression analysis**

1. Divergence of  $\hat{\beta}$
2. Instability of  $\hat{\beta}$

# Ordinary least squares (OLS) method

$$X = [\vec{x}_1, \vec{x}_2, \vec{x}_3, \dots, \vec{x}_n]$$

$$\widehat{\vec{\beta}}_{\text{OLS}} = \underset{\vec{\beta}}{\operatorname{argmin}} \left\{ \left\| \vec{y} - X\vec{\beta} \right\|_{l_2}^2 \right\}$$

$$L(\vec{\beta}) \equiv \left\| \vec{y} - X\vec{\beta} \right\|_{l_2}^2$$

$$\frac{\partial L(\vec{\beta})}{\partial \vec{\beta}} = -[X^T \vec{y} - X^T X \vec{\beta}] = \vec{0}$$

$$\widehat{\vec{\beta}}_{\text{OLS}} = (X^T X)^{-1} X^T \vec{y} \quad (X^T X)^{-1} = \frac{\Delta}{\det|X^T X|}$$

$X^T X$  should be a regular matrix:

i.e.,  $X^T X$  has an inverse matrix

$$\det|X^T X| \neq 0$$

→  $\vec{x}_1, \vec{x}_2, \vec{x}_3, \dots, \vec{x}_n$  should be linearly independent

# If $X^T X$ is not regular matrix...?

$X = [\vec{x}_1, \vec{x}_2, \vec{x}_3, \dots, \vec{x}_n], \quad \vec{x}_1, \vec{x}_2, \vec{x}_3, \dots, \vec{x}_n$  are linearly dependent

1. Make orthogonal bases by basis transformation (PCR, PLS)
2. Solve a relaxation problem with regularizing  $X^T X$  (Ridge, LASSO)
3. Find and remove unnecessary (linearly dependent) descriptors (RREF method)

# Multicollinearity

$$\vec{x}_i = c_1 \vec{x}_1 + c_2 \vec{x}_2 + \cdots + c_n \vec{x}_n$$

$$\vec{x}_i = X^{(-i)} \vec{c}^{(-i)}$$

$$\vec{0} = X\vec{c}$$

i.e., the existence of MCL means the existence of non-trivial solutions of  $X\vec{c} = \vec{0}$ .  
→ for finding MCL, we find the non-trivial solutions.

The condition that  $X\vec{c} = \vec{0}$  has non-trivial solutions

1. Under-determined system, i.e.,  $m < n$ .
2. Rank deficient matrix, i.e.,  $\text{rank}(X) \equiv r < \min(m, n)$

$$m \begin{matrix} n \\ \begin{bmatrix} 1 & 0 & 0 & 0 & a \\ 0 & 1 & 0 & 0 & b \\ 0 & 0 & 1 & 0 & c \\ 0 & 0 & 0 & 1 & d \end{bmatrix} \end{matrix} \quad m \begin{matrix} n \\ \begin{bmatrix} 1 & 0 & 0 & a \\ 0 & 1 & 0 & b \\ 0 & 0 & 1 & c \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \end{matrix}$$

# The number of multicollinearities

Linear span (= subspace spanned by descriptor vectors):

$$W = \text{span}(\{\vec{x}_1, \vec{x}_2, \vec{x}_3, \dots, \vec{x}_n\})$$

Descriptor matrix:

$$X = [\vec{x}_1, \vec{x}_2, \vec{x}_3, \dots, \vec{x}_n]$$

$$\dim(W) = \text{rank}(X) \equiv r$$

The number of non-trivial solutions of  $X\vec{c} = \vec{0}$   
means the number of extra basis in the span  $W$

$$n - \dim(W)$$

then, the number of non-trivial solutions is given by  $n - r$   
(independent one set of solution)

# RREF method

Make row reduced echelon form (rref) of  $X$  by basic operations.

$$RXQ = X_{rref}$$

$R(m \times m)$ : row basic transformation operator (regular)

$Q(n \times n)$ : column basic transformation operator (regular and orthogonal)  
 (here suppose that it just change the order of columns)

$$X_{rref} = \begin{bmatrix} 1 & 0 & 0 & a_1 & b_1 & d_1 \\ 0 & 1 & 0 & a_2 & b_2 & d_2 \\ 0 & 0 & 1 & a_3 & b_3 & d_3 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$n - r = 6 - 3 = 3$$

$$X_{rref}[\vec{c}_1, \vec{c}_2, \vec{c}_3] = [\vec{0}, \vec{0}, \vec{0}]$$

The simplest solutions:

$$[\vec{c}_1, \vec{c}_2, \vec{c}_3] = \begin{bmatrix} -a_1 & -b_1 & -d_1 \\ -a_2 & -b_2 & -d_2 \\ -a_3 & -b_3 & -d_3 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 & a_1 & b_1 & d_1 \\ 0 & 1 & 0 & a_2 & b_2 & d_2 \\ 0 & 0 & 1 & a_3 & b_3 & d_3 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} -a_1 \\ -a_2 \\ -a_3 \\ 1 \\ 0 \\ 0 \end{bmatrix} = \vec{0}$$

# RREF method

The solutions of  $X_{rref}\vec{c} = \vec{0}$  are the solutions of  $X\vec{c} = \vec{0}$  ?

$$RXQ = X_{rref} \quad X = R^{-1}X_{rref}Q^T$$

$$X_{rref}\vec{c} = \vec{0}$$

$$R^{-1}X_{rref}\vec{c} = \vec{0}$$

$$R^{-1}X_{rref}Q^TQ\vec{c} = \vec{0}$$

$$XQ\vec{c} = \vec{0}$$

How to obtain  $Q$

$$XQ = [\vec{x}_1, \vec{x}_2, \vec{x}_3, \vec{x}_4]Q = [\vec{x}_3, \vec{x}_1, \vec{x}_2, \vec{x}_4]$$

$$IQ = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} Q = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} = Q$$

We only have to remember the column order change

# Subspace selection method

From RREF method,

$$\vec{x}_k = c_i^k \vec{x}_i + c_j^k \vec{x}_j$$

$$\vec{x}_l = c_i^l \vec{x}_i + c_j^l \vec{x}_j$$

$$\vec{x}_m = c_i^m \vec{x}_i + c_j^m \vec{x}_j$$



List representation (we call this subspace list)

$$[\vec{x}_k, \vec{x}_l, \vec{x}_m, [\vec{x}_i, \vec{x}_j]]$$

Temporary basis

The number of basis of this span  
 $\text{span}(\{\vec{x}_i, \vec{x}_j, \vec{x}_k, \vec{x}_l, \vec{x}_m\})$   
is only two.

Choose two basis in this list

Select  $\vec{x}_i, \vec{x}_j$  and remove  $\vec{x}_k, \vec{x}_l, \vec{x}_m$

or

Select  $\vec{x}_i, \vec{x}_k$  and remove  $\vec{x}_j, \vec{x}_l, \vec{x}_m$

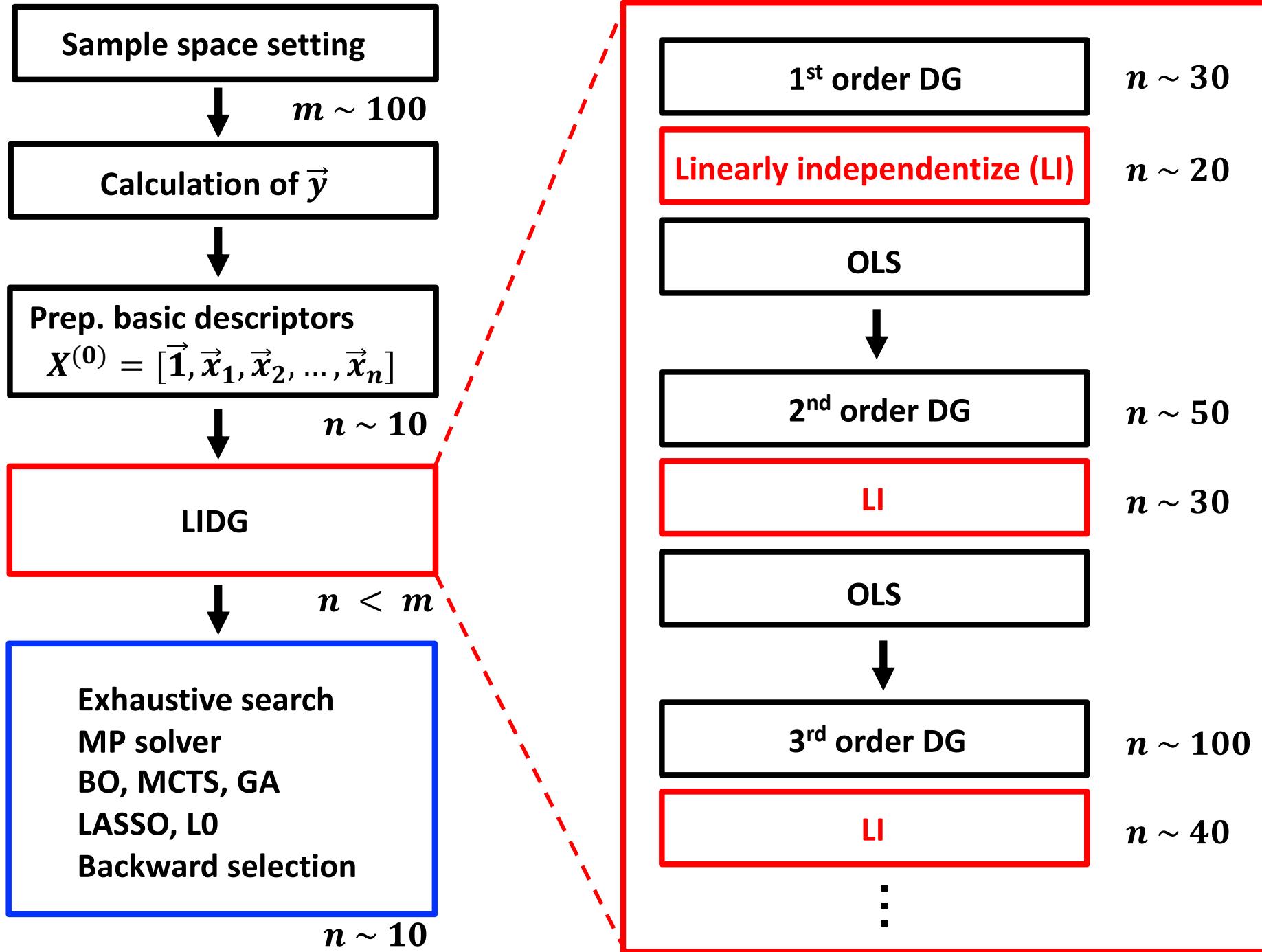
or

Select  $\vec{x}_l, \vec{x}_m$  and remove  $\vec{x}_i, \vec{x}_j, \vec{x}_k$

...

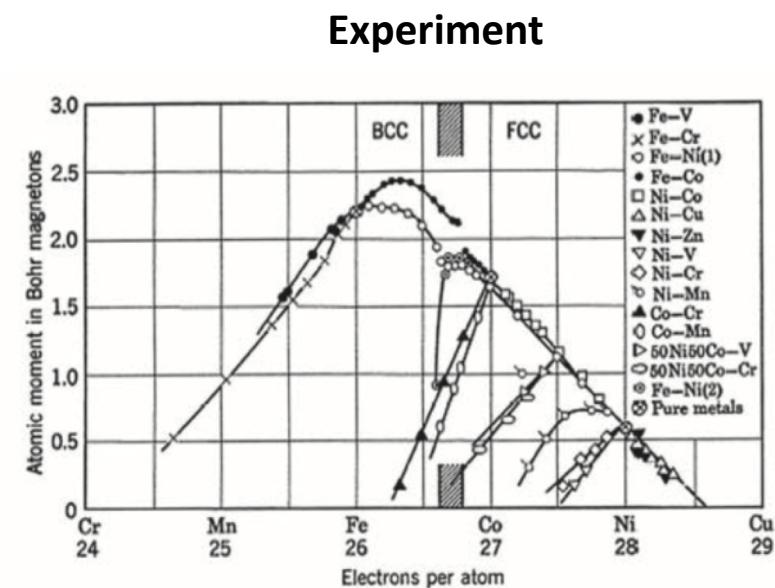
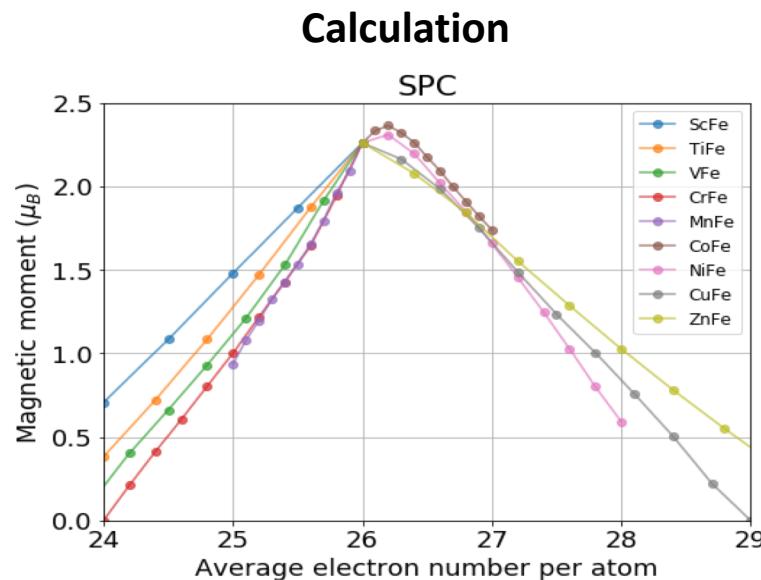
We can reduce the number of descriptors  
without any deterioration in regression accuracy

# Linearly independent descriptor generation method



# Example

AkaiKKR (KKR-CPA)  
 LDA (MJW)  
 Crystal structure: BCC  
 Lattice constant = 2.86 Å  
 Scalar relativistic  
 BZ quality = 10 (nk=256)  
 edelta = 0.0001  
 ewidth = 1.2  
 Complex energy mesh = 85  
 L max = 3



H. Saito: Physics and Applications of Invar Alloys, Maruzen, Tokyo, (1978), 18.

**Everyone may image that there is a simple model on the back of this curve  
 We try to find the simple model by using LIDG method**

**Target property ( $m = 99$ ):**

Magnetic moment of binary alloy  $A_{1-x}B_x$

$$M(A, B, x)$$

**Initial descriptors ( $n = 7$ ):**

Concentration,  $x$

Magnetic moment of pure bulk,  $M_P(A), M_P(B)$

The number of valence electrons,  $Z(A), Z(B)$

Magnetic moment of impurity atom,  $M_I(A, B)$

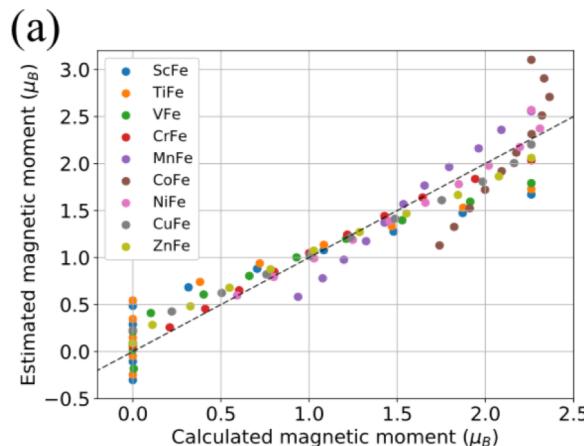
$$M_I(B, A)$$

# The results of OLS with 1st, 2nd, and 3rd order descriptors

**1st order (8)**

$$R^2 = 0.906$$

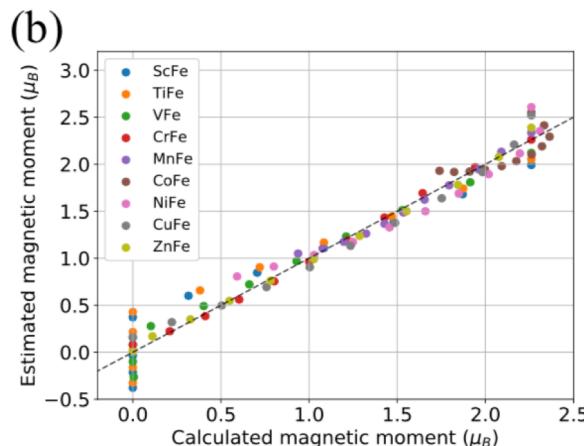
$$Q^2 = 0.887$$



**Up to 2nd order (17)**

$$R^2 = 0.969$$

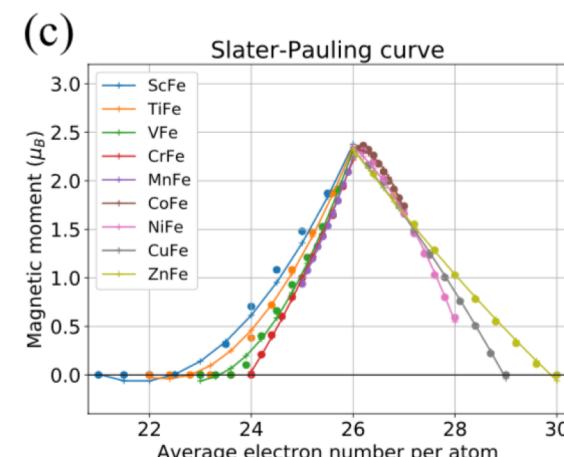
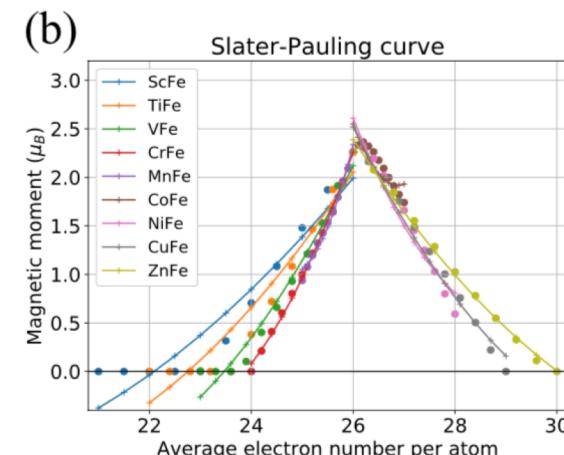
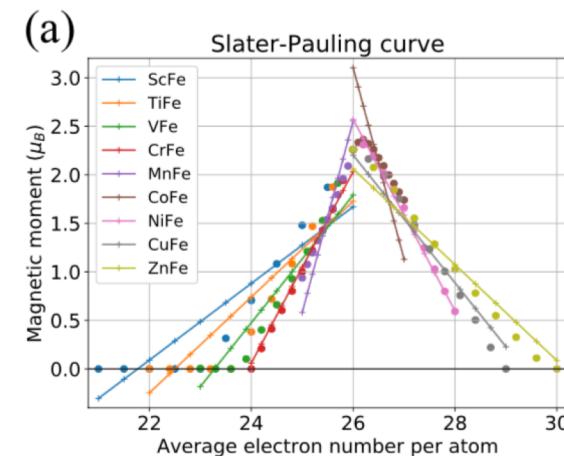
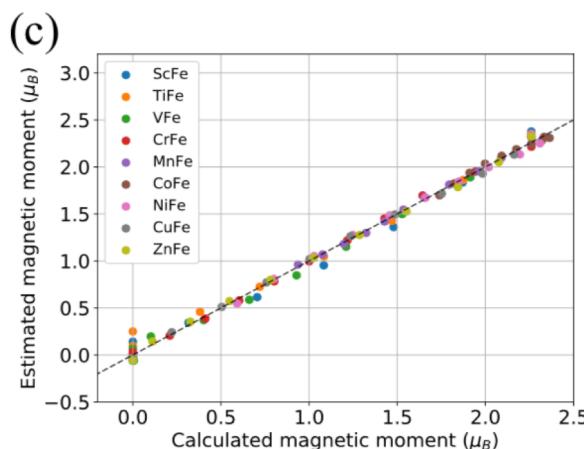
$$Q^2 = 0.950$$



**Up to 3rd order (24)**

$$R^2 = 0.996$$

$$Q^2 = 0.991$$



# Obtained model (by backward selection)

$$\begin{aligned} \mathbf{M}(A, B, x) \sim & (1 - x)\mathbf{M}_P(A) + x\mathbf{M}_P(B) \\ & + x(1 - x)[-4.03 + 0.28Z(A) + 0.62\mathbf{M}_I(B, A) + 0.20|\mathbf{M}_P(A) + \mathbf{M}_I(A, B)|] \end{aligned}$$

An interpretation  
by analogy with the regular solution approximation for binary compound system

$$\mathbf{M}(A, B, x) = \bar{\mathbf{M}}(A, B, x) + \Delta\mathbf{M}(A, B, x)$$

$$\bar{\mathbf{M}}(A, B, x) \equiv (1 - x)\mathbf{M}_P(A) + x\mathbf{M}_P(B) \quad \text{Averaged magnetic moment term (linear to } x\text{)}$$

$$\Delta\mathbf{M}(A, B, x) \equiv x(1 - x)\Omega(A, B)$$

Mixing (excessive) magnetic moment term generated by alloying

$$\Omega(A, B) = -4.03 + 0.28Z(A) + 0.62\mathbf{M}_I(B, A) + 0.20|\mathbf{M}_P(A) + \mathbf{M}_I(A, B)|$$

Interaction parameter (does not depend on  $x$ )

# Summary

- We propose RREF method for detecting MCL
- This is a new approach to solve MCL problem in linear regression analysis
- Systematic descriptor generation method is proposed
- By combining these methods, LIDG method was proposed
- LIDG method was applied to analyze SPCs and we could obtain a simple (interpretable) model with high generalization capability.

[github.com/Hitoshi-FUJII/LIDG](https://github.com/Hitoshi-FUJII/LIDG)  
Y. Kanda, H. Fujii, and T. Oguchi, STAM(2019).

## Raw data

Label	$M(A, B, x)$	$x$	$M_P(A)$	$Z(A)$	$M_P(B)$	$Z(B)$	$M_I(A, B)$	$M_I(B, A)$
Sc100Fe000	0.0000	0.0	0.0000	3	2.2604	8	-0.32667	0.00000
Sc090Fe010	0.0000	0.1	0.0000	3	2.2604	8	-0.32667	0.00000
Sc080Fe020	0.0000	0.2	0.0000	3	2.2604	8	-0.32667	0.00000
Sc070Fe030	0.0000	0.3	0.0000	3	2.2604	8	-0.32667	0.00000
Sc060Fe040	0.0000	0.4	0.0000	3	2.2604	8	-0.32667	0.00000
Sc050Fe050	0.3157	0.5	0.0000	3	2.2604	8	-0.32667	0.00000
Sc040Fe060	0.7056	0.6	0.0000	3	2.2604	8	-0.32667	0.00000
Sc030Fe070	1.0831	0.7	0.0000	3	2.2604	8	-0.32667	0.00000
Sc020Fe080	1.4793	0.8	0.0000	3	2.2604	8	-0.32667	0.00000
Sc010Fe090	1.8709	0.9	0.0000	3	2.2604	8	-0.32667	0.00000
Sc000Fe100	2.2604	1.0	0.0000	3	2.2604	8	-0.32667	0.00000
Ti100Fe000	0.0000	0.0	0.0000	4	2.2604	8	-0.69382	0.00000
Ti090Fe010	0.0000	0.1	0.0000	4	2.2604	8	-0.69382	0.00000
Ti080Fe020	0.0000	0.2	0.0000	4	2.2604	8	-0.69382	0.00000
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
Zn040Fe060	1.2862	0.6	0.0000	12	2.2604	8	0.03119	0.00000
Zn030Fe070	1.5524	0.7	0.0000	12	2.2604	8	0.03119	0.00000
Zn020Fe080	1.8440	0.8	0.0000	12	2.2604	8	0.03119	0.00000
Zn010Fe090	2.0763	0.9	0.0000	12	2.2604	8	0.03119	0.00000
Zn000Fe100	2.2604	1.0	0.0000	12	2.2604	8	0.03119	0.00000

# Descriptors

$M_p$ : Magnetic moment of pure bulk A  
 Z: # of valence electrons

$A$	$M_p(A)$	$Z(A)$
Sc	0.000	3
Ti	0.000	4
V	0.005	5
Cr	0.000	6
Mn	0.937	7
Fe	2.260	8
Co	1.740	9
Ni	0.591	10
Cu	0.000	11
Zn	0.000	12

$M_i(A,B)$ :  
 Magnetic moment of impurity A (B)  
 in bulk B (A)

$A$	$B$	$M_i(A,B)$	$M_i(B,A)$
Sc	Fe	-0.327	0.000
Ti	Fe	-0.694	0.000
V	Fe	-1.165	-0.005
Cr	Fe	-1.644	1.825
Mn	Fe	-1.769	2.387
Co	Fe	1.846	2.673
Ni	Fe	1.075	2.797
Cu	Fe	0.268	2.332
Zn	Fe	0.031	0.000

# LIDG method

