

COMPARISON OF MACHINE LEARNING METHODS AND LINEAR REGRESSION MODELS IN IDENTIFYING IMPORTANT PREDICTOR VARIABLES FOR SERUM DIOXIN TEQ FOR A COMMUNITY IN MICHIGAN, USA

Chen Q¹, Lee S-Y¹, Hedgeman E², Ghosh D¹, Gillespie BW¹, Lepkowski J³, Garabrant D²

¹Department of Biostatistics, University of Michigan School of Public Health, 109 S Observatory, Ann Arbor, MI 48109; ²Department of Environmental Health Sciences, University of Michigan School of Public Health, 109 S Observatory, Ann Arbor, MI 48109; ³Survey Research Center, Institute for Social Research, University of Michigan, 426 Thompson, Ann Arbor, MI 48104;

Keywords: Blood, North America, Environmental samples, Dietary intake, Soil, Humans, TEQs.

Introduction and Study Goals

The primary goal of the University of Michigan Dioxin Exposure Study (UMDES) is to investigate whether known PCDD, PCDF and PCB (hereinafter collectively referred to as 'dioxins') contamination in the Tittabawassee River sediments downstream from the Dow Chemical plant are associated with elevated serum dioxin levels in the regional population. It is well known that age and body mass index (BMI) ¹ can explain a part of human's serum dioxin levels. Therefore, we are interested in identifying which pathways of exposure and methods of elimination are more important in determining the serum dioxin levels after adjusting for age and BMI. The potential pathways include residential proximity to the river, property use, recreational activities in the contaminated region, consumption of food grown or raised in the river or region, consumption of fish and game from contaminated areas, various measures of soil contact, house dust dioxin levels, and occupational contact. Other factors that may affect the serum dioxin levels include weight change and smoking status for both sexes, and pregnancies and breast-feeding history for women. This study seeks to identify important pathways of exposure and methods of elimination for dioxin in serum.

Materials and Methods

In this study, a total of 946 subjects who have complete serum dioxin measures were included. Blood serum, house dust, and soil were analyzed by Alta Analytical Laboratory, Inc. Details of the serum, house dust, and soil sampling methods and analyses are reported elsewhere ^{2,3,4}. Additionally, the process of dealing with limit of detection issues for serum samples is reported elsewhere ⁵. A 10 based logarithm transformation of the serum dioxin concentration expressed as the Toxic Equivalent (TEQ) was taken to reduce the right skewness of the distribution of the serum dioxin concentrations.

The exposure information was obtained from the UMDES questionnaire ⁶. The respondent was asked to recall possible dioxin exposure pathways over their entire lifetime. These pathways included a full residential history, occupations, property use, recreational activities, and consumption of meat, fish, game, eggs, milk, other dairy products, and vegetables. Basic demographic (age, gender, race, education, income) and health questions (height, weight, weight loss and gain, smoking status, pregnancy, childbearing and months of breastfeeding for each child) were also included. The function forms of all potential predictor variables are reported elsewhere ⁷.

The statistical analyses were performed using two different approaches. First, we used linear regression method, in which we forced age, age square and BMI into the model, and used a backward selection procedure to identify the important predictors of the serum dioxin TEQ after adjusting for the effects of age and BMI. Details of the modeling procedure is reported elsewhere ⁷. An advantage of linear regression modeling is that it provides estimates of effect magnitudes in addition to identify the significant explanatory variables. Second, we used machine learning methods, specifically, Random Forest ⁸ and tree-boosting ⁹ models. The residuals of logarithm transformed serum dioxin

concentration after adjusting for age, age square, and BMI were obtained as outcomes in the machine learning methods. The relative importance of each potential predictor variable to the overall model was calculated. The advantage of using the tree-boosting method is that there is no need to consider input variable transformations, since all tree-boosting procedures are invariant under all strictly monotone transformations of the individual input variables ⁹. All the analyses were done using statistical software SAS 9.1 ¹⁰ and R 2.2.1 ¹¹, respectively, for linear regression models and machine learning methods. The results of these two approaches for modeling the total serum dioxin TEQ were compared. The consistencies and inconsistencies of the two sets of results are reported.

Results and Discussion

Table 1 lists the important predictor variables identified in the linear regression models, and the ranks of the estimated importance of predictor variables using machine learning methods. For the machine learning methods, the higher the rank (highest=1), the more important is the predictor variable. Both methods agree that XXX are more important in predicting the blood dioxin levels. However, XXX are identified as important pathways in the linear regression models, but not by the machine learning models. XXX have high importance measures in machine learning models, but are not significant in the linear regression models.

Table 1: Important predictor variables identified using linear regression models and machine learning methods

Category	Variables	Linear Regression	Random Forest	Tree-boosting
Dust				
Soil				
Health				
Demographics				
Food				
Residence History				
Occupations				
Property Use				
Recreational Activities				

In addition, Figure 1 displays single variable partial dependence plots on those most important predictor variables identified by both methods for serum dioxin levels. The figures shows that XXX

Figure 1

In conclusion, XXX

Acknowledgements

The authors acknowledge the Dow Chemical Company for funding the study and Ms. Sharyn Vantine for her continued assistance.

References

1. Patterson Jr, DG, Patterson, D, Canady, R, Wong, L-Y, Lee, R, Turner, W, Caudill, S, Needham, L, Henderson, A. *Organohalogen Comp* 2004; 66:2878-2883.
2. Hedgeman E, Chen Q, Gillespie BW, Franzblau A, Knutson K, Zwica L, Sima C, Lee S-Y, Lepkowski J, Ward B, Ladronka K, Olson K, Sinibaldi J, Towey T, Adriaens P, Demond A, Chang S-C, Gwinn D, Swan S, Garabrant D. *Organohalogen Comp* 2006 (Forthcoming)
3. Zwica L, Knutson K, Towey T, Hedgeman E, Franzblau A, Chen Q, Lee S-Y, Sima C, Gillespie BW, Adriaens P, Demond A, Lepkowski J, Ward B, Ladronka K, Olson K, Sinibaldi J, Chang S-C, Gwinn D, Swan S, Garabrant D. *Organohalogen Comp* 2006 (Forthcoming)
4. Adriaens P, Demond A, Towey T, Chang S-C, Chen Q, Franzblau A, Gillespie BW, Gwinn D, Hedgeman E, Knutson K, Ladronka K, Lee S-Y, Lepkowski J, Olson K, Sima C, Sinibaldi J, Swan S, Ward B, Zwica L, Garabrant D. *Organohalogen Comp* 2006 (Forthcoming)
5. Gillespie BW, Chen Q, Lee SY, Hong B, Garabrant D, Hedgeman E, Sima C, Lepkowski J, Olson K, Luksemburg W. *Organohalogen Comp* 2006 (forthcoming).
6. Olson, K¹, Lepkowski, J¹, Lohr-Ward, B¹, Ladronka, K¹, Sinibaldi, J¹, Garabrant, D², Franzblau, A², Adriaens, P³, Gillespie, B⁴, Bandyk, J¹, Chang, S-C², Chen, Q², Demond, A³, Gwinn, D⁴, Hedgeman, E², Hong, B², Knutson, K², Lee, S-Y³, Sima, C², Towey, T³, Wright, D², Zwica, L². *Organohalogen Comp* 2006 (forthcoming).
7. Garabrant D¹, Franzblau A¹, Lepkowski J², Adriaens P³, Demond A³, Hedgeman E¹, Knutson K¹, Zwica L¹, Chen Q⁴, Olson K², Ward B², Towey T³, Ladronka K², Sinibaldi J², Chang S-C³, Lee S-Y⁴, Gwinn D⁵, Sima C⁴, Swan S⁵, Gillespie BW⁴. *Organohalogen Comp* 2006 (forthcoming).
8. Breiman L. "Random Forests". *Machine Learning* 2001, 45: 5-32
9. Friedman J.H. "Greedy Function Approximation: A Gradient Boosting Machine". *The Annals of Statistics* 2001, 29: 1189-1232
10. SAS Institute. SAS/STAT User's Guide Version 8. Cary, NC: SAS Institute Inc.
11. The R Foundation for Statistical Computing, Version 2.2.1, 2005