

# Stochastic & Data Analysis Methods & Applications in Statistics and Demography

*Editors*

**James R. Bozeman**

**Teresa Oliveira**

**Christos H. Skiadas**

**ISAST**



**Stochastic and Data Analysis Methods and  
Applications in Statistics and Demography**

*Edited by*

**James R. Bozeman, Teresa Oliveira and Christos H. Skiadas**



**ISAST 2016**

**James R. Bozeman, Teresa Oliveira and  
Christos H. Skiadas, *Editors***

**Stochastic and Data Analysis Methods and  
Applications in Statistics and Demography**

**Copyright © 2016 by ISAST**

All rights reserved. No part of this book may be reproduced or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, or otherwise, without prior written permission of ISAST.

ISBN (print) : 978-618-5180-18-8

ISBN (e-book) : 978-618-5180-19-5

# On ordered categorical modelling for complex skill development

Natalya Pya<sup>1</sup> and Arman Kussainov<sup>2</sup>

<sup>1</sup> Department of Mathematical Sciences, School of Science and Technology, Nazarbayev University, Astana, Kazakhstan

(E-mail: natalya.pya@nu.edu.kz)

<sup>2</sup> Physics and Technology Department, al-Farabi Kazakh National University, Almaty, Kazakhstan

(E-mail: arman.kussainov@gmail.com)

**Abstract.** In cognitive science there has been considerable interest in the understanding of expertise development. Models for exploring human complex skill development are often based on comparisons between experts and novices, and use measurements of performance at different levels of skills as predictors. In this paper we study the development of expertise by analysing video game telemetry data collected from a real-time strategy game. Data that relate to cognitive-motor abilities, attentional and perceptual processes were collected from StarCraft 2 game players from seven levels of expertise. We develop an extended generalized additive model for ordered categorical data to investigate the effects of predictors on skill development.

**Keywords:** ordered categorical, cognitive science, generalized additive models, skill learning.

## 1 Introduction

Thompson *et al.*[18] conducted a study exploring human complex skill development. Their aim was to identify potential predictors of expertise in real-time strategy (RTS) video games using the telemetric data collected from RTS StarCraft 2 game players. StarCraft 2 is a popular video game which has millions of players worldwide. Thompson *et al.*[18] examined measures that relate to cognitive-motor abilities, attentional and perceptual processes. Using random forest classifiers, they disproved the assumption that importance of variables across skill levels remains static. Moreover, they argue that telemetric data can become a standard tool for studying human cognition and learning. As different expert areas such as, e.g. chess, basketball, surgery, are expected to show sufficient consistency in development of expertise, many studies have been devoted to exploring skill development in strategy games (Chase and Simon[8], Charness[9], Ericsson and Charness[10]).

This paper proposes to investigate development of expertise using additive regression modelling. The paper develops an extended generalized additive model for ordered categorical data (Wood *et al.*[19]) to study the effects of predictors on skill

---

*Stochastic and Data Analysis Methods and Applications in Statistics and Demography*  
(pp. 677-685)

James R. Bozeman, Teresa Oliveira and Christos H. Skiadas (Eds)

© 2016 ISAST



learning. Modelling categorical responses using smooth functions of predictors allows us to confirm Thompson *et al.*[18] findings and to further investigate the effects of predictors on skill development.

## 2 Video game data

The telemetric data were collected from 3,360 RTS StarCraft 2 game players from 7 levels of expertise. The dataset is public and available at UCI Machine Learning Repository (Bache and Lichman[5]). For each player, the level of expertise is measured by the league in which they contend and serves as an ordered response  $Y_i$ .  $Y_i$  takes a value from  $r = 1, 2, \dots, 7$ , indicating Bronze, Silver, Gold, Platinum, Diamond, Master, and Professional leagues. There are eighteen predictor variables available including measures of attentional control, perceptual process and cognitive-motor speed. The examination of the data and preliminary modelling revealed 13 variables relevant to skill development. Table 1 summarizes the predictors under study. The time at which values of the predictors is recorded is in terms of timestamps in the StarCraft 2 replay file. `GapBwPACs`, `ActionLatency`, `NumberOfPACs`, and `ActionsInPAC` are four variables that refer to a certain period of time during which a player performs at a specific location. Perception action cycle (PAC) was defined by Thompson *et al.*[18] as screen fixations with one or more actions. For the complete information about the variables used in the study, see Thompson *et al.*[18].

## 3 Modelling approach

Many models have been proposed to analyze ordered categorical data which became well-known by virtue of Cox[7] and Plackett[16]. The most appealing regression models for ordered categories are cumulative logit (proportional-odds version of the cumulative logit) models expressed in terms of a latent usually unobservable continuous variable proposed by McCullagh[15], Anderson and Philips[4], Hastie and Tibshirani[13]. McCullagh[15] and Anderson and Philips[4] introduced parametric regression models with ordered categorical responses, whereas Hastie and Tibshirani[13] extended this to a non-parametric version. The parameter estimation for those models is based on maximizing likelihood assuming independent multinomial observations using Fisher scoring algorithm. The cumulative logit models were also discussed in Anderson[3], Agresti[1], Agresti[2], Goodman[12]. Fahrmeir and Lang[11], Kneib and Fahrmeir[14] developed a general class of semi-parametric additive regression models for categorical responses from a Bayesian perspective.

### Extended generalized additive model

The model proposed in this paper is within a new general framework to generalized additive modelling for non-exponential family responses introduced by Wood *et al.*[19]. The framework of Wood *et al.*[19] proposes two methods for the generalized additive models (GAM) generalization: an extended GAM fitting for the cases with a single linear predictor and a log likelihood expressed as a sum over the log likelihood for each response datum; and a general model estimation when log likelihood

**Table 1.** Telemetric data characteristics

Name	Description	Min	Max
APM	Action per minute	22.06	389.83
SelectByHotkeys	Number of unit or building selections made using hotkeys per timestamp	0	0.043
AssignToHotkeys	Number of units or buildings assigned to hotkeys per timestamp	0	$1.75 \cdot 10^{-3}$
UniqueHotkeys	Number of unique hotkeys used per timestamp	0	10
MinimapAttacks	Number of attack actions on minimap per timestamp	0	$3.02 \cdot 10^{-3}$
NumberOfPACs	Number of perception action cycles (PAC) per timestamp	$6.79 \cdot 10^{-4}$	$7.97 \cdot 10^{-3}$
GapBwPACs	Mean duration in milliseconds between PACs	6.667	237.143
ActionLatency	Mean latency from the onset of PACs to their first action in milliseconds	24.09	176.37
ActionsInPAC	Mean number of actions within each PAC	2.039	18.558
TotalMapExplored	The number of 24x24 game coordinate grids viewed by the player per timestamp	5	58
WorkersMade	Number of SCVs, drones, and probes trained per timestamp	$7.7 \cdot 10^{-5}$	$5.15 \cdot 10^{-3}$
UniqueUnitsMade	Unique unites made per timestamp	2	13
ComplexAbilUsed	Abilities requiring specific targeting instructions used per timestamp	0	$3.08 \cdot 10^{-3}$

depends non-linearly on smooth functions of predictors. The first method includes such distributions outside the exponential family as beta, zero inflated Poisson, negative binomial, Tweedie, scaled t distribution and ordered categorical data. The GAM fitting method is extended for these models. The second extension requires different approach for model fitting and general and reliable smoothing parameter estimation. It covers such models as Cox proportional hazard (Cox[6]) and Cox process models, generalized additive models for location scale and shape proposed by Rigby and Stasinopoulos[17] and multivariate additive models (Yee and Wild[21]). Below is a brief description of modelling with ordered categorical responses within a new extended GAM.

Consider independent response observations,  $y_i$ , that take values from  $r = 1, \dots, R$ , where  $r$  is ordered category label. A latent variable  $u_i = \mu_i + \epsilon_i$  is introduced with the c.d.f. of  $\epsilon_i$  being  $F$ . Then, given  $-\infty = \alpha_0 < \alpha_1 < \dots < \alpha_R = \infty$ ,  $y_i = r$  if a latent variable  $u_i$  is such that  $\alpha_{r-1} < u_i \leq \alpha_r$ ,

$$P(Y_i = r) = F(\alpha_r - \mu_i) - F(\alpha_{r-1} - \mu_i).$$

The usual choice for the c.d.f. of  $\epsilon$  is the standard logistic or normal. For identifiability reasons  $\alpha_1 = -1$ , so there are  $R - 2$  extra unknown parameters. To impose

increasing ordering on the cutting points,  $\alpha_r$  are set as

$$\alpha_r = \alpha_1 + \sum_{j=1}^{r-1} \exp(\theta_j), \quad 1 < r < R,$$

so  $\theta_j$  are parameters to be estimated. The mean value of the latent variable depends on the predictor variable in the following way,

$$\mu_i = \mathbf{A}_i \boldsymbol{\gamma} + \sum_j f_j(x_{ji}),$$

where  $\mathbf{A}$  is a model matrix for the strictly parametric terms,  $\boldsymbol{\gamma}$  is a vector of unknown parameters,  $f_j$  is an unknown smooth function of the predictor variable  $x_j$ , where  $x_j$  can be vector valued. Each smooth term is represented by reduced rank spline smoothers  $f_j(x_j) = \sum_k \beta_{kj} b_{kj}(x_j)$ , where  $b_{kj}$  are known spline basis functions,  $\beta_{kj}$  unknown coefficients. Then, the mean of the latent variable can be expressed as  $\boldsymbol{\mu} = \mathbf{X}\boldsymbol{\beta}$ , with the model matrix  $\mathbf{X}$  combining  $\mathbf{A}$  and matrix of spline basis, and  $\boldsymbol{\gamma}$  being a part of  $\boldsymbol{\beta}$ .

The log likelihood of the model can be written as

$$l = \sum_{i=1}^n l_i(y_i, \mu_i, \boldsymbol{\theta}),$$

where  $l_i$  is the log likelihood for each observation,  $\boldsymbol{\theta}$  is a  $(R - 2)$ -vector of the extra parameters,  $\theta_j$ , that control the thresholds. Then, the deviance corresponding to the observation  $y_i$  is defined in the standard way as  $D_i = 2(\tilde{l}_i - l_i)$ , where  $l_i = \max_{\mu_i} l_i(y_i, \mu_i, \boldsymbol{\theta})$  is the saturated log likelihood. Given  $\boldsymbol{\theta}$ , the parameters  $\boldsymbol{\beta}$  are estimated by minimization of the penalized deviance

$$\mathcal{D}(\boldsymbol{\beta}, \boldsymbol{\theta}) = \sum_i D_i(\boldsymbol{\beta}, \boldsymbol{\theta}) + \sum_j \lambda_j \boldsymbol{\beta}^T \mathbf{S}^j \boldsymbol{\beta},$$

where a quadratic penalty term  $\boldsymbol{\beta}^T \mathbf{S}^j \boldsymbol{\beta}$  measuring function smoothness is associated with each smooth  $f_j$  and  $\lambda_j$  being a smoothing parameter. Penalized iteratively re-weighted least squares (PIRLS) is applied for  $\boldsymbol{\beta}$  estimation. Estimation of  $\boldsymbol{\theta}$  is achieved by minimization of the negative Laplace approximate marginal likelihood (LAML),

$$\mathcal{V} = \frac{\mathcal{D}(\hat{\boldsymbol{\beta}}, \boldsymbol{\theta})}{2} - \tilde{l}(\boldsymbol{\theta}) + \frac{\log |\mathbf{X}^T \mathbf{W} \mathbf{X} + \mathbf{S}^\lambda| - \log |\mathbf{S}^\lambda|_+}{2} - \frac{M_p}{2} \log(2\pi),$$

where  $\mathbf{S}^\lambda = \sum_j \lambda_j \mathbf{S}^j$  and  $|\mathbf{S}^\lambda|_+$  is the product of the positive eigenvalues of  $\mathbf{S}^\lambda$ ,  $M_p$  is the number of zero eigenvalues of  $\mathbf{S}^\lambda$ . Newton's or a quasi-Newton's method is used for  $\mathcal{V}$  minimization. Several issues with numerical instability have to be taken into account to make the optimization procedure efficient and reliable. This is fully covered in Wood *et al.*[19]. Generalized additive modelling with ordered categorical data as well as other extensions are implemented in an R package `mgcv` available at CRAN (Wood[20]).

**Additive model for video game data**

The preliminary backward selection, first in the framework of a generalized linear model and then in the framework of an extended GAM, revealed thirteen covariates relevant to skill development (see section 2). The extended GAM for ordered categorical data with  $R = 7$  was fitted with the selected set of predictors using smooth terms. We first consider a model with all selected predictors having non-linear effects on the mean of the ordered categorical latent variable.

*Model 1:*

$$\begin{aligned} \mu_i = & f_1(\text{NumberOfPACs}_i) + f_2(\text{UniqueHotkeys}_i) + f_3(\text{WorkersMade}_i) \\ & + f_4(\text{GapBwPACs}_i) + f_5(\text{ActionLatency}_i) + f_6(\text{AssignToHotkeys}_i) \\ & + f_7(\text{MinimapAttacks}_i) + f_8(\text{APM}_i) + f_9(\text{SelectByHotkeys}_i) \\ & + f_{10}(\text{ActionsInPAC}_i) + f_{11}(\text{TotalMapExplored}_i) \\ & + f_{12}(\text{UniqueUnitsMade}_i) + f_{13}(\text{ComplexAbilUsed}_i), \end{aligned}$$

where the model terms  $f_1 - f_{13}$  are unknown smooth functions of the corresponding predictors. Thin plate regression splines are used for their representations. The predictor values were preprocessed using square root or log transformation in order to avoid gaps with very small amount of data that account for the Professional league. There was a significant linear dependence of `NumberOfPACs`, `UniqueHotkeys` and `WorkersMade` on the mean of the latent variable, so that it was sufficient to add strictly parametric structure for these three predictors. The resulted model has the following structure.

*Model 2:*

$$\begin{aligned} \mu_i = & \beta_1 \cdot \text{NumberOfPACs}_i + \beta_2 \cdot \text{UniqueHotkeys}_i + \beta_3 \cdot \text{WorkersMade}_i \\ & + f_1(\text{GapBwPACs}_i) + f_2(\text{ActionLatency}_i) + f_3(\text{AssignToHotkeys}_i) \\ & + f_4(\text{MinimapAttacks}_i) + f_5(\text{APM}_i) + f_6(\text{SelectByHotkeys}_i) \\ & + f_7(\text{ActionsInPAC}_i) + f_8(\text{TotalMapExplored}_i) \\ & + f_9(\text{UniqueUnitsMade}_i) + f_{10}(\text{ComplexAbilUsed}_i), \end{aligned}$$

where  $\beta_1$ ,  $\beta_2$  and  $\beta_3$  are unknown parameters.

Including the bivariate smooth of the APM and the second most important variable, `SelectByHotkeys`, gives better model in comparison with the model with univariate effect of the APM. Moreover, constructing a tensor product interaction of `GapBwPACs` and `ActionLatency`, with their main effects being included separately further improves the model fit. The following additive structure for the mean value of the ordered categorical latent variable was considered as the third model.

*Model 3:*

$$\begin{aligned} \mu_i = & \beta_1 \cdot \text{NumberOfPACs}_i + \beta_2 \cdot \text{UniqueHotkeys}_i + \beta_3 \cdot \text{WorkersMade}_i \\ & + f_1(\text{GapBwPACs}_i) + f_2(\text{ActionLatency}_i) \\ & + f_3(\text{GapBwPACs}_i, \text{ActionLatency}_i) + f_4(\text{AssignToHotkeys}_i) \\ & + f_5(\text{MinimapAttacks}_i) + f_6(\text{APM}_i, \text{SelectByHotkeys}_i) \\ & + f_7(\text{ActionsInPAC}_i) + f_8(\text{TotalMapExplored}_i) \\ & + f_9(\text{UniqueUnitsMade}_i) + f_{10}(\text{ComplexAbilUsed}_i), \end{aligned}$$

where all the predictors except for the first three have nonparametric smooth effects. A tensor product interaction of `GapBwPACs` and `ActionLatency`, is used for representing  $f_3$  with the main effects comprised in  $f_1$  and  $f_2$ .

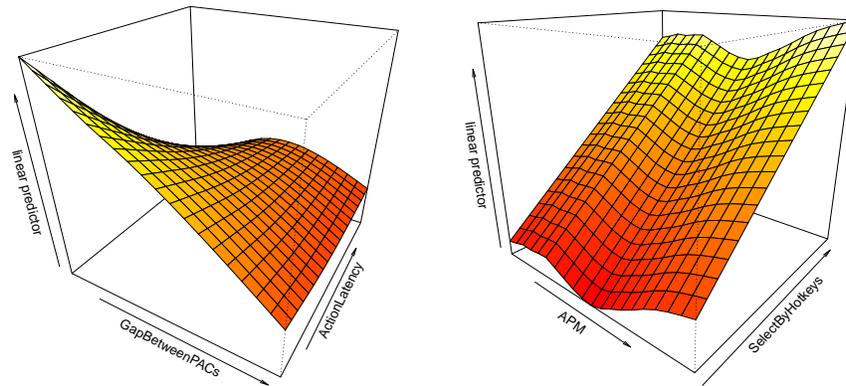
## 4 Results and discussion

In addition to the above mentioned models, we fitted submodels with certain terms omitted. The model selection procedure showed that the best model in terms of the Akaike information criterion is the full model 3. Other model performance measures such as generalized cross validation score, adjusted  $r^2$  and percentage deviance explained were also better for model 3 than for other considered models.

Figure 1 illustrates the estimated effects of the two bivariate smooths of model 3. `APM` variable is used as a measure of cognitive speed. This variable is shown to have the highest rank of the predictive importance Thompson *et al.* ([18]) in distinguishing Bronze-Professional classifier. The first decreasing trend of the `APM` (figure 1, right panel) is due to the high correlation between the `APM` and `SelectByHotkeys` variables (higher values of the covariate effect correspond to higher league level). Both predictors have increasing trends when considered separately as smooths of a single variable.

The estimates of the univariate effects are shown in figure 2. As expected, the main effects of the two characteristics of the PACs, `GapBwPACs` and `ActionLatency`, are decreasing with increase in the skill level (panels (a) and (b)), while the other two have increasing trends (panel (e) and the positive parametric effect of `NumberOfPACs`). An adaptive smoother was used to estimate the effect of the `ActionsInPAC`, that allowed the degree of smoothing to vary along the range of the predictor values.

The strong increasing effect of the `MinimapAttacks` (fig. 2, panel (d)) which is used as a measure of the attentional control, supports the hypothesis of Thompson *et al.* [18] that more skillful players act on minimap more often. On the contrary the `TotalMapExplored` has a decreasing trend (panel (f)), more skillful players view the total map less often. Usage of hotkeys allows players to build and control game units in more efficient way, so that the higher values of the `AssignToHotkeys` variable would correspond to higher level of expertise (panel (c)). The estimated effects of the last two predictors do not have such monotonic features as for the other smooths. Players in the highest leagues seem to use moderately abilities that require particular targeting instructions (panel (h)), and keeping from the production of the Unique Units (panel (g)). Whereas, the lowest league players make medium number of units while avoiding complex abilities. The estimated monotone increasing trend of



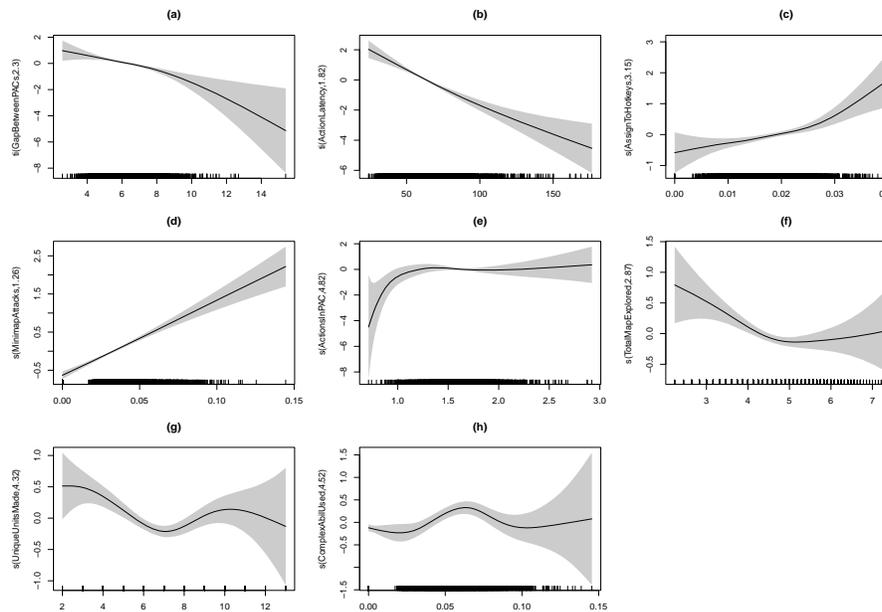
**Fig. 1.** Video game data: the estimated interactions between Gap Between PACs and Action Latency variables, and between APM and Select By Hotkeys.

the `WorkersMade` shows that to progress players must produce more workers. However, the importance of this variable diminishes for the highest league (Thompson *et al.* [18]), which is explained by possible automatization of the worker production skill. Moreover, to advance in skills, players are required to put more effort on managing their learning and increasing cognitive demand, which is reflected by the positive linear trend of the `UniqueHotKeys` predictor.

This study supports Thompson's *et al.* [18] proposition that telemetric data can be used as a standard tool for studying human cognition and learning. Moreover, the proposed model confirms the previous findings that the assumption that importance of predictors across skill levels remains static is not correct. We showed that constructing non-isotropic tensor product splines used to model smooth interactions improves prediction of skill development. Modelling categorical responses using smooth functions of predictors allows to capture skill learning in a continuous fashion.

## References

1. A. Agresti. *Categorical data analysis*, volume 359. John Wiley and Sons, 2002.
2. A. Agresti. *Analysis of ordinal categorical data*, volume 656. Wiley, 2010.
3. J.A. Anderson, J. A. Regression and ordered categorical variables. *Journal of the Royal Statistical Society. Series B*, 1–30, 1984.



**Fig. 2.** Video game data: the estimated univariate smooth effects.

4. J. Anderson and P. Philips. Regression, discrimination and measurement models for ordered categorical variables. *Applied Statistics*, 22–31, 1981.
5. K. Bache and M. Lichman. UCI machine learning repository, 2013.
6. D.D.R. Cox. Regression models and life tables. *Journal of the Royal Statistical Society. Series B*, 34(2), 187–220, 1972.
7. D.D.R. Cox. *The analysis of binary data*. volume 32, CRC Press, 1989.
8. W.G. Chase and H.A. Simon. Perception in chess. *Cognitive Psychology*, 4, 55–81, 1973.
9. N. Charness. Age, skill, and bridge bidding: A chronometric analysis, *Journal of Verbal Learning and Verbal Behavior*, 22(4), 406–416, 1983.
10. K.A. Ericsson, N. Charness. Expert performance. *American Psychologist*, 49(8), 725–747, 1994
11. L. Fahrmeir and S. Lang. Bayesian semiparametric regression analysis of multicategorical time-space data. *Annals of the Institute of Statistical Mathematics*, 53(1), 11–30, 2001.
12. L.A. Goodman. The analysis of dependence in cross-classifications having ordered categories, using log-linear models for frequencies and log-linear models for odds. *Biometrics*, 149–160, 1983.
13. T. Hastie and R. Tibshirani. Non-parametric logistic and proportional odds regression. *Applied statistics*, 260–276, 1987.
14. T. Kneib and L. Fahrmeir. Structured additive regression for categorical spacetime data: A mixed model approach. *Biometrics*, 62(1), 109–118, 2006.
15. P. McCullagh. Regression models for ordinal data. *Journal of the Royal Statistical Society. Series B*, 109–142, 1980
16. R.L. Plackett. *The Analysis of Categorical Data*. Griffin, London, 1981.
17. R. Rigby and D.M. Stasinopoulos. Generalized additive models for location, scale and shape. *Journal of the Royal Statistical Society. Series C*, 54(3), 507–554, 2005

18. J.J. Thompson, M.R. Blair, L. Chen, and A.J. Henrey. Video game telemetry as a critical tool in the study of complex skill learning. *PloS One*, 8, 9, e75129, 2013.
19. S. Wood, N. Pya, and B. Säfken. Smoothing parameter and model selection for general smooth models. *Journal of the American Statistical Association*. DOI: 10.1080/01621459.2016.1180986.
20. S. Wood. mgcv: Mixed GAM Computation Vehicle with GCV/AIC/REML smoothness estimation. R package version 1.8-6, 2015.
21. T.W. Yee and C. Wild. Vector generalized additive models. *Journal of the Royal Statistical Society. Series B*, 481–493, 1996