# Simulations description

## Estimator

Method used in simulations is based on joint estimation described in article 'Inference for Regression with Variables Generated by AI or Machine Learning'. Instead of using linear model as authors did

$$Y_i = \gamma^T \theta_i + \alpha^T q_i + \varepsilon_i$$

I used generalised linear model with Poisson distribution

$$Y_i \sim \mathcal{P}(\mu_i), \qquad \log \mu_i = \gamma^T \theta_i + \alpha^T q_i.$$

In model I assumed that we have 10 classes so $\theta_i$'s are 10-dimensional vectors with one 1 and nine 0's and that $q_i$ (data without measurement error) are one-dimensional. We also need error matrix $\Omega$ that on $i$-th row and $j$-th column has estimated probability of model to classify true value $j$ of $\theta$ as $i$.

Having this assumed we can compute likelihood by adding over all possible true values of $\theta_i$

$$l(Y_i, \hat{\theta}_i = d; \gamma, \alpha, \omega) = \sum_{k=1}^{10} \omega_{dk} \cdot \frac{e^{-\mu_{ik}} \mu_{ik}^{Y_i}}{Y_i!}$$

with $\mu_{ik}$ being $k$-th element of $\gamma$ plus $\alpha q_i$. and $\omega_{dk} = P(\hat{\theta}_i = d, \theta_i = k)$, that is probability of model classifying value of $\theta_i$ as $d$ when true value is $k$.

The total log-likelihood is

$$L = \sum_{i=1}^{n} \log l(Y_i, \hat{\theta}_i = d; \gamma, \alpha, \omega). \tag{*}$$

After search for values $\hat{\gamma}, \hat{\alpha}$ that maximizes $L$.

## Simulation

I made two simulation: repeated 25 simulation of 1000 observations and repeated 25 simulation of 5000 observations.

In smaller data set I drawn gamma from $U(1,2)$, $\alpha$ from $U(-0.5, 0.5)$, error matrix as scaled matrix with accuracy of 70% and errors following Poisson distribution, $\theta_i$ were drawn with given probability $p$, responses of AI model were drawn using $\theta_i$'s and error matrix and $Y_i$ were drawn from Poisson distribution.

Similar case was for bigger data set but this time probability for drawing $\theta_i$ was also randomized and so was accuracy of predictions in error matrix (from $U(60, 90)$)

Then I applied R function optimx to maximize likelihoods given by (*) with added square penalty function to prevent overfitting. This method was used twice – firstly in case when probabilities of occurring true classes were known and secondly when they were estimated from joint estimation. Results of implemented methods were compared with the results of glm. For each simulation was computed sum of squares of prediction errors (difference between true and predicted parameters).

# Results

In both bigger and smaller number of observations joint estimation was better option. In smaller set it achieved mean error 0.1728 with known probabilities and 0.1835 with unknown probabilities compared to standard glm with error 0.2155. In set of 5000 observations difference was significantly bigger – joint estimation with known probabilities had mean error 0.0098, joint estimation with predicted probabilities has mean error 0.0801 when glm had mean error 0.4839.