
Using Random Effects to Account for High-Cardinality Categorical Features and Repeated Measures in Deep Neural Networks

Anonymous Author(s)

Affiliation

Address

email

1 Simulations

1.1 $g(Z) = Z$

- This simulation implements equation (14) from the paper
- All runs were made on a Nvidia Quadro P620 GPU on a Windows machine, implemented in Python 3.8 Numpy + Pandas suite, Keras and Tensorflow 2.2
- Code is fully available in the `lmmnn` package in attached zip and later on Github
- Running code: see details in package README file
- $n = 100,000, \sigma_e^2 = 1$
- $\epsilon \sim \mathbb{N}(0, \sigma_e^2)$
- $p = 10$, i.e. 10 fixed predictors in X where each $X_k \sim \mathbb{U}(-1, 1), k \in \{1, \dots, p\}$
- Single categorical random variable with q levels where $q \in \{100, 1000, 10000\}$
- b is a q -length vector of i.i.d random effects, sampled from a $\mathbb{N}(0, \sigma_b^2)$ distribution where $\sigma_b^2 \in \{0.1, 1, 10\}$
- $n_{iter} = 5$ replications for each (q, σ_b^2) combination
- In total $3 \times 3 \times 5 = 45$ runs, for each of 6 types: {Ignore, OHE, Embedding, lme4, MeNets, LMMNN}
- At each run 80% (80,000) of the simulated data is used as training set, of which 10% (8,000) is used as validation set which the network only uses to check for early stopping. That leaves 20% of the data (20,000) as testing set and we record y_{te} vs. \hat{y}_{te} mean squared error and standard error over 10 replication per condition.
- Max no. of epochs: 500
- batch size: 30
- We use a Keras `EarlyStopping` callback to stop training if no improvement has been observed for 10 epochs.
- Baseline DNN architecture: 4 hidden fully connected layers with $\{100, 50, 25, 12\}$ number of neurons, a ReLU activation and a Dropout of 25% in each, and a final output layer with a single neuron with no activation
- Loss used is mean squared error in all networks except for LMMNN
- Optimizer: Adam, Keras default params
- Specific architecture:

- 31 – Ignore: input is X of dimension p
- 32 – OHE: input is X and Z of dimension $p + q$
- 33 – Embedding: input is X of dimension p , in addition Z goes through an Embedding
- 34 layer which maps q levels to a $d = 0.1 \cdot q$ vector, so input dimension is $p + d$
- 35 – LMMNN: input is X of dimension p , in addition Z , y_{tr} and final single neuron output
- 36 are input to the custom NLL loss layer. We initialize $(\hat{\sigma}_e^2, \hat{\sigma}_b^2)$ to be (1.0, 1.0). Note that
- 37 Z is input to the NLL loss layer via a sparse vector, we do not actually keep a $n \times q$
- 38 matrix.
- 39 – MeNets: input is X of dimension p and the layer before last of 12 neurons is the feature
- 40 mapping used in the MeNets V-EM algorithm.

41 1.2 $g(Z) = ZW$

- 42 • This simulation implements equation (15) from the paper
- 43 • Changes from simulation in 1.1:
 - 44 – $W_{q \times d}$ is a linear transformation of Z
 - 45 – $d = 0.1 \cdot q$ and we sample W from a $\mathbb{U}(-1, 1)$ distribution
 - 46 – b is now a d -length vector of i.i.d random effects, sampled from a $\mathbb{N}(0, \sigma_b^2)$ as before
 - 47 – Max no. of epochs when $q = 10,000$: 250
 - 48 – DNNs architecture remains the same except in LMMNN where Z goes through a
 - 49 Embedding layer which maps q levels to a $d = 0.1 \cdot q$ vector before it is input to the
 - 50 NLL loss layer
 - 51 – For LMMNN we use our own custom `EarlyStoppingWithSigmasConvergence`
 - 52 callback which also makes sure the estimated variance components $(\hat{\sigma}_e^2, \hat{\sigma}_b^2)$ have
 - 53 converged for 10 epochs.

54 2 Real Data

55 2.1 UKB-PA: Estimating physical activity from self-reported behaviors from the UK 56 Biobank

- 57 • Relevant notebook: `ukb_pa.ipynb`
- 58 • Data availability: upon request from the UK Biobank only
- 59 • Physical activity (PA) definition: Subjects wore an accelerometer on their wrist for 7 days.
- 60 Physical activity is measured in ENMO units (euclidean norm minus one), calculated on the
- 61 acceleration vector in three axes, and negative values were truncated to zero. Mean wrist
- 62 ENMO in m-g was summarised across valid wear-time.
- 63 • ETL: We follow instructions by Pearce et al. (2020), implemented in R. At high level, we
- 64 filter out subjects wearing the accelerometer for less than 72 hours or having ENMO of over
- 65 80. Each of the categorical behavioral variables e.g. "frequency of walking for pleasure"
- 66 is converted to numerical with a simple mapping e.g. "once a week" is converted to 1 and
- 67 "every day" is converted to 7. Finally the PA dependent variable is standardized to have a
- 68 mean of 0 and standard deviation of 1 for each fold, for the training set.
- 69 • $n = 96,629$
- 70 • Categorical feature: Job ($q = 353$)
- 71 • Fixed features:
 - 72 – Gender
 - 73 – Heavy work time
 - 74 – Walking during work time
 - 75 – Sedentary during work time
 - 76 – Moderate-to-vigorous physical activity time (MVPA)
 - 77 – Walking for pleasure time
 - 78 – Strenuous sport time

- 79 – Other activities time
- 80 – Light DIY time
- 81 – Heavy DIY time
- 82 – TV time
- 83 – Computer time
- 84 – Sleep time
- 85 – Getting about method (OHE): walk, cycle, transport, other
- 86 – Commute method (OHE): walk, cycle, transport, other
- 87 • Baseline DNN architecture: Pearce et al. did not use DNNs, but two separate linear
- 88 regressions, for men and women. We use a simple MLP of 2 fully connected layers with
- 89 ReLU activation of 10 and 5 neurons, followed by a single output neuron with no activation.
- 90 • Split policy: 5-fold CV, from each training fold 10% of data is kept as validation data which
- 91 the network uses for early stopping
- 92 • Max no. of epochs: 2000
- 93 • Batch size: 30
- 94 • Callbacks: EarlyStopping with patience = 10
- 95 • Runtime: Google Colab, Nvidia Tesla P100 GPU

96 2.2 Drugs: Estimating Drugs 1-10 rating from textual reviews from Drugs.com

- 97 • Relevant notebook: `drugs.ipynb`
- 98 • Data availability: freely available in the UCI Machine Learning Repository [https://](https://archive.ics.uci.edu/ml/datasets/Drug+Review+Dataset+%28Drugs.com%29)
- 99 archive.ics.uci.edu/ml/datasets/Drug+Review+Dataset+%28Drugs.com%29
- 100 • ETL: We bind Gräber et al. (2018) training and testing set into a single set, on which we
- 101 perform regular 5-fold cross validation. We use only the drugs reviews, perform standard
- 102 tokenization to words using Keras Tokenizer with maximum 10,000 most common words
- 103 (which is in a sense the fixed feature dimension p) and each review is cut at maximum length
- 104 100 words.
- 105 • $n = 215,063$
- 106 • Categorical feature: Drug ($q = 3,671$)
- 107 • Baseline DNN architecture: Gräber et al. did not use DNNs, we therefore use a standard
- 108 text input architecture. After tokenization reviews are fed into a standard Embedding layer
- 109 of dimension 100. We then use a LSTM layer of 64 kernels, followed by a single output
- 110 neuron with no activation.
- 111 • Split policy: 5-fold CV, from each training fold 10% of data is kept as validation data which
- 112 the network uses for early stopping
- 113 • Max no. of epochs: 100
- 114 • Batch size: 30
- 115 • Callbacks: EarlyStopping with patience = 5
- 116 • Runtime: Google Colab, Nvidia Tesla P100 GPU

117 2.3 Airbnb: Predicting prices of Airbnb rentals

- 118 • Relevant notebook: `airbnb.ipynb`
- 119 • Data availability: freely available following Kalehbasti et al. (2019) instructions at [https:](https://github.com/PouyaREZ/AirBnbPricePrediction)
- 120 [//github.com/PouyaREZ/AirBnbPricePrediction](https://github.com/PouyaREZ/AirBnbPricePrediction)
- 121 • ETL: We run Kalehbasti et al. code as is including their variable selection procedure and
- 122 log transformation for the price. We then bind their training, validation and testing into a
- 123 single set, on which we perform regular 5-fold cross validation.
- 124 • $n = 49,976$
- 125 • Categorical feature: Host ($q = 39,393$)

- Fixed features: $p = 196$ features after variable selection, e.g.: does rental has 24-hour check-in, does it have a dryer, is it kids friendly etc.
- Baseline DNN architecture: We use Kalehbasti et al. architecture: a simple MLP of 2 fully connected layers with ReLU activation of 20 and 5 neurons, followed by a single output neuron with no activation. We also use their parameters for the Adam optimizer.
- Split policy: 5-fold CV, from each training fold 10% of data is kept as validation data which the network uses for early stopping
- Max no. of epochs: 500
- Batch size: 30
- Callbacks: EarlyStopping with patience = 10
- Runtime: Google Colab, Nvidia Tesla P100 GPU

2.4 UKB-Age: Predicting age from retinal fundus photographs from the UK Biobank

- Relevant notebook: `ukb_age.ipynb`
- Data availability: upon request from the UK Biobank only
- ETL: Fundus retinal PNG images of the UK Biobank are originally of size $2,048 \times 1,536$ pixels. We crop the images to include only the retinal circle blocked by a $1,440 \times 1,440$ square which is then resized to size 299×299 , with 3 color channels
- $n = 175,817$
- Categorical feature: Subject ($q = 85,721$)
- Baseline DNN architecture: Poplin et al. (2018) did not specify their CNN architecture exactly, only that it uses InceptionV3 at the bottom. We found that a standard CNN architecture works better:
 - Conv. 2D layer, 32 kernels, (5, 5) strides, padding valid, ReLU activation
 - Max 2D pooling, (2, 2) pool size
 - Conv. 2D layer, 64 kernels, (5, 5) strides, padding valid, ReLU activation
 - Max 2D pooling, (2, 2) pool size
 - Conv. 2D layer, 32 kernels, (5, 5) strides, padding valid, ReLU activation
 - Max 2D pooling, (2, 2) pool size
 - Conv. 2D layer, 16 kernels, (5, 5) strides, padding valid, ReLU activation
 - Max 2D pooling, (2, 2) pool size
 - Flatten, Dropout of 50
 - Fully connected layer of 100 neurons and ReLU activation
 - Single output neuron with no activation
- Split policy: 5-fold CV, from each training fold 10% of data is kept as validation data which the network uses for early stopping. For this dataset the split is made on *subjects*, i.e. the images of 80% of subjects are use in each fold as training set, so at testing time all of the categorical feature's levels are new.
- Max no. of epochs: 100
- Batch size: 20
- Callbacks: EarlyStopping with patience = 10
- Runtime: Google Colab, Nvidia Tesla P100 GPU