

Part 1: Build a model to predict the total delivery duration seconds

Introduction

The target problem can be defined as a regression problem for prediction based on supervised learning. This work uses Python and PyTorch to build linear/polynomial regression and DNN regression models to predict total delivery duration seconds. This design doc includes details in Data Pre-processing, Feature Extraction & Normalization, Model design, Experiments, Conclusion and Future work, which can be used together with submitted code package.

- `dataManager.py`: For Data Pre-processing, Feature Extraction & Normalization
- `model.py`: For model design
- `buildModel.py`: Run experiments

Data Pre-processing

- Define target output

According to *historical_data.csv*, the total delivery duration seconds is between *actual_delivery_time* (a) and *created_at* (b), which includes *estimated_store_to_consumer_driving_duration* (c) and *estimated_order_place_duration* (d). Considering both of them come from other prediction models, in this work, define the target output as:

- $$prediction = (a - b - c - d)$$

Then the final total delivery duration seconds will be $prediction + c + d$.

- When *estimated_store_to_consumer_driving_duration* is "NA"

Set 0 seconds instead (more than 500 records)

- When *actual_delivery_time* or *created_at* are not available

Remove records (only several records)

- When $total_onshift_dashers < total_busy_dashers$

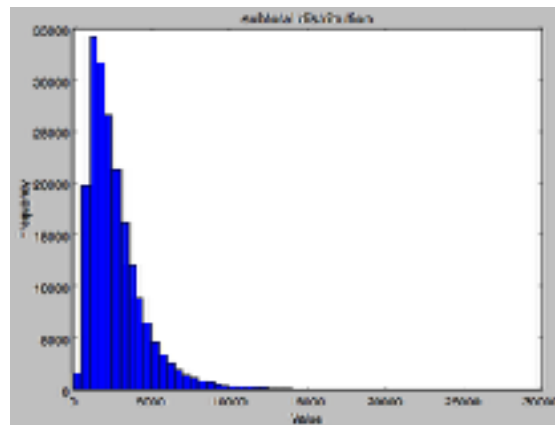
According to data description, *total_busy_dashers* is the subset of *total_onshift_dashers*, which should be always less or equal to *total_onshift_dashers*. Set them with the same value when $total_onshift_dashers < total_busy_dashers$. (more than 40000 records as noise information added)

Feature Extraction & Normalization

- **Features with Gamma distribution (5 in total)**

- subtotal
- min_item_price
- max_item_price
- total available dashers (total_onshift_dashers - total_busy_dashers)
- total_outstanding_orders

By analysis and data visualization, all of the above five features can be defined as Gamma probability distribution function. Below is the histogram of subtotal as example:



Based on the theory of Gamma probability distribution, use Cumulative Distribution Function (CDF) to normalize features into the range (0~1). Implementation is as follows:

```
import scipy.stats as st

def fitGamma(data):

    fit_alpha, fit_loc, fit_beta=st.gamma.fit(data)

    data = (data - fit_loc)/fit_beta

    data=st.gamma.cdf(data, fit_alpha)

    return data
```

- **Features by embedding selection (4 in total)**

- store_primary_category
- order_protocol

- `total_items`
- `num_distinct_items`

It is easy to understand that *store_primary_category* and *order_protocol* are category information, which can be converted to 1-dimension vector with total categories.

(*store_primary_category*: 75 categories in total, *order_protocol*: 8 categories in total)

But for *total_items* and *num_distinct_items*, although they are digital values with Gamma probability distribution, because of small value range, after fitting them by CDF, small values are converted to 0, which is useless for supervised learning. Based on this consideration, also apply embedding selection on both of them. (*total_items*: set 30 in total, *num_distinct_items*: set 20 in total).

In PyTorch, implemented by `torch.nn.Embedding(EmbCategory, EmbDimension)`.

- **Target Normalization**

Also use Gamma CDF to normalize target value to the range (0~1) for speeding up learning process.

Model Design

All model design and implementation are based on PyTorch `nn.Module` package. Model and normalization parameters will be saved after training.

- **Linear model**

Straight-forward, design by `torch.nn.Linear(inputSize, outputSize)`

- **Polynomial (Quadratic) model**

Design quadratic terms by `torch.nn.Bilinear(inputSize, inputSize, outputSize)`, then add linear terms together as output

- **DNN model**

The designed model supports multiple layers, that can be as many layers as you want, and also supports to use Sigmoid, Relu, Tanh and SELU activation functions for each layer.

- **MSE as Cost Function**

All models use MSE (from `nn.Module`) as cost function to monitor model performance, and also calculate RMSE on validation set.

In order to calculate RMSE between prediction and real target values, need to do reverse CDF of prediction, then calculate RMSE with real target values. Implementation of reverse CDF is as follows:

```
import scipy.stats as st

def reverseCDF(data, gammaPara):

    fit_alpha, fit_loc, fit_beta=gammaPara[0], gammaPara[1], gammaPara[2]

    data=st.gamma.ppf(data, fit_alpha)* fit_beta +fit_loc

    return data
```

- **Cross-validation to evaluate model performance**

After shuffling the data from historical_data.csv, the whole data is split into train data (train_data.csv) and validation data (validation_data.csv) by the ratio of 90% and 10%.

The data in train_data.csv is used for training of model parameters, and the data in validation_data.csv is used for validation to get better hyper-parameters and models.

The most important thing of using train data and validation data is it will be much easier to monitor the issue of under-fitting or overfitting, then users can find good ways to fix.

- **Use SGD and mini-batch for optimization**

Experiments

When run buildModel.py in command line, the combination of inputs of *modelDefine* and *useBilinear* can control which kind of model will be built.

- **Linear model**

Run: python buildModel.py --modelDefine 'None' --learningRate 0.05

```
Namespace(batchSize=64, dropout=0.0, embDim4Category=10, embDim4DistinctItems=10, embDim4Prot
ocol=5, embDim4TotalItems=10, epochNumber4Valid=2, learningRate=0.05, maxEpoch=30, modelDefin
e='None', savedModel='savedModel.net', trainData='train_data.csv', useBilinear=False, validDa
ta='validation_data.csv')
Epoch[2/30], Train MSE: 0.072042, Validate MSE: 0.070099, RMSE with real value: 1047.66
Epoch[4/30], Train MSE: 0.071416, Validate MSE: 0.070412, RMSE with real value: 1027.01
Epoch[6/30], Train MSE: 0.069918, Validate MSE: 0.069017, RMSE with real value: 1026.52
Epoch[8/30], Train MSE: 0.069885, Validate MSE: 0.069403, RMSE with real value: 1025.10
Epoch[10/30], Train MSE: 0.069340, Validate MSE: 0.068812, RMSE with real value: 1041.61
Epoch[12/30], Train MSE: 0.069361, Validate MSE: 0.068829, RMSE with real value: 1040.33
Epoch[14/30], Train MSE: 0.069087, Validate MSE: 0.068682, RMSE with real value: 1036.75
Epoch[16/30], Train MSE: 0.069106, Validate MSE: 0.069058, RMSE with real value: 1033.71
Epoch[18/30], Train MSE: 0.068977, Validate MSE: 0.068720, RMSE with real value: 1035.87
Epoch[20/30], Train MSE: 0.068976, Validate MSE: 0.068650, RMSE with real value: 1035.70
Epoch[22/30], Train MSE: 0.068971, Validate MSE: 0.068716, RMSE with real value: 1037.96
Epoch[24/30], Train MSE: 0.068926, Validate MSE: 0.068625, RMSE with real value: 1035.92
Epoch[26/30], Train MSE: 0.068915, Validate MSE: 0.068673, RMSE with real value: 1035.23
Epoch[28/30], Train MSE: 0.068892, Validate MSE: 0.068560, RMSE with real value: 1037.08
Epoch[30/30], Train MSE: 0.068886, Validate MSE: 0.068589, RMSE with real value: 1036.26
```

- **Polynomial (Quadratic) model**

Run: `python buildModel.py --modelDefine 'None' --learningRate 0.001 --useBilinear`

```
Namespace(batchSize=64, dropout=0.0, embDim4Category=10, embDim4DistinctItems=10, embDim4Prot
ecol=5, embDim4TotalItems=10, epochNumber4Valid=2, learningRate=0.001, maxEpoch=30, modelDefi
ne='None', savedModel='savedModel.net', trainData='train_data.csv', useBilinear=True, validDa
ta='validation_data.csv')
Epoch[2/30], Train MSE: 0.254673, Validate MSE: 0.107518, RMSE with real value: 1321.50
Epoch[4/30], Train MSE: 0.095983, Validate MSE: 0.087659, RMSE with real value: 1187.04
Epoch[6/30], Train MSE: 0.033105, Validate MSE: 0.080053, RMSE with real value: 1137.46
Epoch[8/30], Train MSE: 0.077741, Validate MSE: 0.077555, RMSE with real value: 1105.44
Epoch[10/30], Train MSE: 0.071852, Validate MSE: 0.074517, RMSE with real value: 1101.15
Epoch[12/30], Train MSE: 0.073117, Validate MSE: 0.071085, RMSE with real value: 1085.49
Epoch[14/30], Train MSE: 0.071935, Validate MSE: 0.071543, RMSE with real value: 1080.74
Epoch[16/30], Train MSE: 0.071096, Validate MSE: 0.072297, RMSE with real value: 1070.11
Epoch[18/30], Train MSE: 0.079437, Validate MSE: 0.070541, RMSE with real value: 1067.95
Epoch[20/30], Train MSE: 0.070045, Validate MSE: 0.069415, RMSE with real value: 1062.02
Epoch[22/30], Train MSE: 0.069613, Validate MSE: 0.069979, RMSE with real value: 1060.72
Epoch[24/30], Train MSE: 0.069334, Validate MSE: 0.068905, RMSE with real value: 1059.66
Epoch[26/30], Train MSE: 0.069050, Validate MSE: 0.069323, RMSE with real value: 1060.94
Epoch[28/30], Train MSE: 0.068850, Validate MSE: 0.068757, RMSE with real value: 1057.70
Epoch[30/30], Train MSE: 0.068683, Validate MSE: 0.068554, RMSE with real value: 1050.92
```

- **DNN model with linear features**

Run: `python buildModel.py --modelDefine '300|S|30|S' --learningRate 0.05`

'300|S|30|S' means there is two hidden layers, 300 and 30 represent the node number of the first and second hidden layer, S represents Sigmoid activation function.

```
Namespace(batchSize=64, dropout=0.0, embDim4Category=10, embDim4DistinctItems=10, embDim4Prot
ecol=5, embDim4TotalItems=10, epochNumber4Valid=2, learningRate=0.05, maxEpoch=30, modelDefin
e='300|S|30|S', savedModel='savedModel.net', trainData='train_data.csv', useBilinear=False, v
alidData='validation_data.csv')
Epoch[2/30], Train MSE: 0.070447, Validate MSE: 0.071029, RMSE with real value: 1050.54
Epoch[4/30], Train MSE: 0.071985, Validate MSE: 0.069620, RMSE with real value: 1033.15
Epoch[6/30], Train MSE: 0.070511, Validate MSE: 0.068970, RMSE with real value: 1039.76
Epoch[8/30], Train MSE: 0.069995, Validate MSE: 0.068868, RMSE with real value: 1042.28
Epoch[10/30], Train MSE: 0.069799, Validate MSE: 0.072016, RMSE with real value: 1050.65
Epoch[12/30], Train MSE: 0.069324, Validate MSE: 0.068761, RMSE with real value: 1038.53
Epoch[14/30], Train MSE: 0.069285, Validate MSE: 0.068961, RMSE with real value: 1033.89
Epoch[16/30], Train MSE: 0.069148, Validate MSE: 0.069195, RMSE with real value: 1032.12
Epoch[18/30], Train MSE: 0.069072, Validate MSE: 0.068743, RMSE with real value: 1037.05
Epoch[20/30], Train MSE: 0.069060, Validate MSE: 0.068805, RMSE with real value: 1035.31
Epoch[22/30], Train MSE: 0.069037, Validate MSE: 0.068719, RMSE with real value: 1036.35
Epoch[24/30], Train MSE: 0.069032, Validate MSE: 0.068721, RMSE with real value: 1037.53
Epoch[26/30], Train MSE: 0.069012, Validate MSE: 0.068768, RMSE with real value: 1035.83
Epoch[28/30], Train MSE: 0.069009, Validate MSE: 0.068714, RMSE with real value: 1036.15
Epoch[30/30], Train MSE: 0.069007, Validate MSE: 0.068721, RMSE with real value: 1037.31
```

- **DNN model with linear and quadratic features**

Run: `python buildModel.py --modelDefine '300|S|30|S' --learningRate 0.001 --useBilinear`

```

NameSpace(batchSize=64, dropout=0.0, embDim4Category=10, embDim4DistinctItems=10, embDim4Protocol=5, embDim4TotalItems=10, epochNumber4Valid=3, learningRate=0.001, maxEpoch=30, modelDefine='306|S|30|S', savedModel='savedModel.net', trainData='train_data.csv', useBilinear=True, validationData='validation_data.csv')
Epoch[2/30], Train MSE: 0.077595, Validate MSE: 0.075660, RMSE with real value: 1076.82
Epoch[4/30], Train MSE: 0.075548, Validate MSE: 0.074523, RMSE with real value: 1067.61
Epoch[6/30], Train MSE: 0.074654, Validate MSE: 0.073700, RMSE with real value: 1065.63
Epoch[8/30], Train MSE: 0.073944, Validate MSE: 0.073145, RMSE with real value: 1059.60
Epoch[10/30], Train MSE: 0.073327, Validate MSE: 0.072539, RMSE with real value: 1064.42
Epoch[12/30], Train MSE: 0.072756, Validate MSE: 0.071874, RMSE with real value: 1057.73
Epoch[14/30], Train MSE: 0.072220, Validate MSE: 0.071400, RMSE with real value: 1053.29
Epoch[16/30], Train MSE: 0.071757, Validate MSE: 0.071394, RMSE with real value: 1043.39
Epoch[18/30], Train MSE: 0.071328, Validate MSE: 0.070351, RMSE with real value: 1048.77
Epoch[20/30], Train MSE: 0.070970, Validate MSE: 0.070250, RMSE with real value: 1046.04
Epoch[22/30], Train MSE: 0.070649, Validate MSE: 0.069891, RMSE with real value: 1046.19
Epoch[24/30], Train MSE: 0.070390, Validate MSE: 0.069695, RMSE with real value: 1040.67
Epoch[26/30], Train MSE: 0.070103, Validate MSE: 0.069407, RMSE with real value: 1047.08
Epoch[28/30], Train MSE: 0.069906, Validate MSE: 0.069418, RMSE with real value: 1047.66
Epoch[30/30], Train MSE: 0.069876, Validate MSE: 0.069212, RMSE with real value: 1041.63

```

• Hyper-parameters Tuning

When do tuning on buildModel.py, hyper-parameters include:

- learningRate
- batchSize
- maxEpoch and epochNumber4Valid
- embDim4Category
- embDim4Protocol
- embDim4TotalItems
- embDim4DistinctItems
- dropout

And also including modelDefine for DNN layers.

Conclusion and Future work

- From experiments, MSE of all models on train data and validation data after 30 epochs are similar, which means models are just fine without overfitting or under-fitting based on available data.
- The performance of all models are better than the baseline on validation set after 30 epochs.
- Linear model and DNN model with linear feature are a little better and faster than models with quadratic terms of features. And it needs to set a much smaller learning rate for the training with quadratic terms. The reason might be there is no normalization for quadratic terms, so all features have different scales.

- The model performance do not have significant improvement by hyper-parameter tuning. The reason might come from data with too much noisy information.
- In future, with more clean and complete data, DNN models should be more powerful to get more accurate prediction results, considering the main advantage of DNN is to get more complex relationship between features from training.
- Also in future, consider to add feature to present different days (weekday or weekend) and different time in each day, which should also be valuable information on the prediction.