# homework4

February 18, 2019

## 0.1 Submission Info

Name : Hanmaro Song, 3032216175 Teammates : Minjune Hwang, Kyle Nguyen, Joanne Chen, Kyle Cho

# 1 Homework 4 - Berkeley STAT 157

**Your name: XX, SID YY, teammates A,B,C** (Please add your name, SID and teammates to ease Ryan and Rachel to grade.)

Handout 2/12/2019, due 2/19/2019 by 4pm in Git by committing to your repository.

In this homework, we will build a model based real house sale data from a Kaggle competition. This notebook contains codes to download the dataset, build and train a baseline model, and save the results in the submission format. Your jobs are

1. Developing a better model to reduce the prediction error. You can find some hints on the last section.

2. Submitting your results into Kaggle and take a sceenshot of your score. Then replace the following image URL with your screenshot.
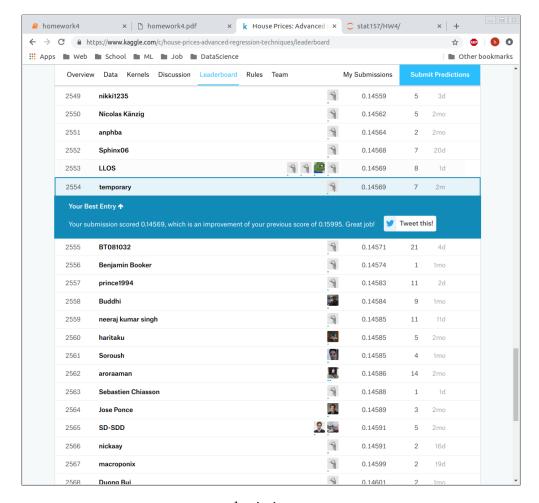
We have two suggestions for this homework:

1. Start as earlier as possible. Though we will cover this notebook on Thursday's lecture, tuning hyper-parameters takes time, and Kaggle limits #submissions per day.
2. Work with your project teammates. It's a good opportunity to get familiar with each other.

Your scores will depend your positions on Kaggle's Leaderboard. We will award the top-3 teams/individuals 500 AWS credits.

## 1.1 Accessing and Reading Data Sets

The competition data is separated into training and test sets. Each record includes the property values of the house and attributes such as street type, year of construction, roof type, basement condition. The data includes multiple datatypes, including integers (year of construction), discrete labels (roof type), floating point numbers, etc.; Some data is missing and is thus labeled 'na'. The price of each house, namely the label, is only included in the training data set (it's a competition after all). The 'Data' tab on the competition tab has links to download the data.

We will read and process the data using `pandas`, an efficient data analysis toolkit. Make sure you have `pandas` installed for the experiments in this section.

Overview    Data    Kernels    Discussion    Leaderboard    Rules    Team      My Submissions    Submit Predictions

| 2549 | nikki1235 | | 0.14559 | 5 | 3d |
| 2550 | Nicolas Känzig | | 0.14562 | 5 | 2mo |
| 2551 | anphba | | 0.14564 | 2 | 2mo |
| 2552 | Sphinx06 | | 0.14568 | 7 | 20d |
| 2553 | LLOS | | 0.14569 | 8 | 1d |
| 2554 | temporary | | 0.14569 | 7 | 2m |

**Your Best Entry ↑**
Your submission scored 0.14569, which is an improvement of your previous score of 0.15995. Great job!    Tweet this!

| 2555 | BT081032 | | 0.14571 | 21 | 4d |
| 2556 | Benjamin Booker | | 0.14574 | 1 | 1mo |
| 2557 | prince1994 | | 0.14583 | 11 | 2d |
| 2558 | Buddhi | | 0.14584 | 9 | 1mo |
| 2559 | neeraj kumar singh | | 0.14585 | 11 | 11d |
| 2560 | haritaku | | 0.14585 | 5 | 2mo |
| 2561 | Soroush | | 0.14585 | 4 | 1mo |
| 2562 | aroraaman | | 0.14586 | 14 | 2mo |
| 2563 | Sebastien Chiasson | | 0.14588 | 1 | 1d |
| 2564 | Jose Ponce | | 0.14589 | 3 | 2mo |
| 2565 | SD-SDD | | 0.14591 | 5 | 2mo |
| 2566 | nickaay | | 0.14591 | 2 | 16d |
| 2567 | macroponix | | 0.14599 | 2 | 19d |
| 2568 | Duong Bui | | 0.14601 | 2 | 1mo |

submission.png

I had to create a temporary id since the one I used before already joined this competition with my friends in a group (not from this class and that was about an year ago) so my username is temporary.

```
In [1]: # If pandas is not installed, please uncomment the following line:
        # !pip install pandas

        %matplotlib inline
        import d2l
        from mxnet import autograd, gluon, init, nd
        from mxnet.gluon import data as gdata, loss as gloss, nn, utils
        import numpy as np
        import pandas as pd
```

We downloaded the data into the current directory. To load the two CSV (Comma Separated Values) files containing training and test data respectively we use Pandas.

```
In [2]: # utils.download('https://github.com/d2l-ai/d2l-en/raw/master/data/kaggle_house_pred_t
        # utils.download('https://github.com/d2l-ai/d2l-en/raw/master/data/kaggle_house_pred_t
        train_data = pd.read_csv('kaggle_house_pred_train.csv')
        test_data = pd.read_csv('kaggle_house_pred_test.csv')
```

The training data set includes 1,460 examples, 80 features, and 1 label., the test data contains 1,459 examples and 80 features.

```
In [3]: print(train_data.shape)
        print(test_data.shape)
```

```
(1460, 81)
(1459, 80)
```

Let's take a look at the first 4 and last 2 features as well as the label (SalePrice) from the first 4 examples:

```
In [4]: train_data.iloc[0:4, [0, 1, 2, 3, -3, -2, -1]]
```

```
Out[4]:    Id  MSSubClass MSZoning  LotFrontage SaleType SaleCondition  SalePrice
        0   1          60       RL         65.0       WD        Normal     208500
        1   2          20       RL         80.0       WD        Normal     181500
        2   3          60       RL         68.0       WD        Normal     223500
        3   4          70       RL         60.0       WD        Abnorml     140000
```

We can see that in each example, the first feature is the ID. This helps the model identify each training example. While this is convenient, it doesn't carry any information for prediction purposes. Hence we remove it from the dataset before feeding the data into the network.

```
In [5]: all_features = pd.concat((train_data.iloc[:, 1:-1], test_data.iloc[:, 1:]))
```

## 1.2 Data Preprocessing

As stated above, we have a wide variety of datatypes. Before we feed it into a deep network we need to perform some amount of processing. Let's start with the numerical features. We begin by replacing missing values with the mean. This is a reasonable strategy if features are missing at random. To adjust them to a common scale we rescale them to zero mean and unit variance. This is accomplished as follows:

$$x \leftarrow \frac{x - \mu}{\sigma}$$

To check that this transforms $x$ to data with zero mean and unit variance simply calculate $\mathbf{E}[(x - \mu)/\sigma] = (\mu - \mu)/\sigma = 0$. To check the variance we use $\mathbf{E}[(x - \mu)^2] = \sigma^2$ and thus the transformed variable has unit variance. The reason for 'normalizing' the data is that it brings all features to the same order of magnitude. After all, we do not know *a priori* which features are likely to be relevant. Hence it makes sense to treat them equally.

```
In [6]: numeric_features = all_features.dtypes[all_features.dtypes != 'object'].index
        all_features[numeric_features] = all_features[numeric_features].apply(
            lambda x: (x - x.mean()) / (x.std()))
        # after standardizing the data all means vanish, hence we can set missing values to 0
        all_features = all_features.fillna(0)
```

Next we deal with discrete values. This includes variables such as 'MSZoning'. We replace them by a one-hot encoding in the same manner as how we transformed multiclass classification data into a vector of 0 and 1. For instance, 'MSZoning' assumes the values 'RL' and 'RM'. They map into vectors $(1,0)$ and $(0,1)$ respectively. Pandas does this automatically for us.

```
In [7]: # Dummy_na=True refers to a missing value being a legal eigenvalue, and creates an ind
        all_features = pd.get_dummies(all_features, dummy_na=True)
        all_features.shape
```

```
Out[7]: (2919, 354)
```

You can see that this conversion increases the number of features from 79 to 331. Finally, via the `values` attribute we can extract the NumPy format from the Pandas dataframe and convert it into MXNet's native representation - NDArray for training.

```
In [8]: n_train = train_data.shape[0]
        train_features = nd.array(all_features[:n_train].values)
        test_features = nd.array(all_features[n_train:].values)
        train_labels = nd.array(train_data.SalePrice.values).reshape((-1, 1))
```

## 1.3 Training

To get started we train a linear model with squared loss. This will obviously not lead to a competition winning submission but it provides a sanity check to see whether there's meaningful information in the data. It also amounts to a minimum baseline of how well we should expect any 'fancy' model to work.

```
In [9]: loss = gloss.L2Loss()

        def get_net():
            net = nn.Sequential()
            net.add(nn.Dense(80, activation='relu'))
            net.add(nn.Dropout(.5))
            net.add(nn.BatchNorm())
            net.add(nn.Dense(20, activation='relu'))
            net.add(nn.Dropout(.5))
            net.add(nn.BatchNorm())
            net.add(nn.Dense(1))
            net.initialize(init=init.Xavier())
            return net
```

House prices, like shares, are relative. That is, we probably care more about the relative error $\frac{y-\hat{y}}{y}$ than about the absolute error. For instance, getting a house price wrong by USD 100,000 is terrible in Rural Ohio, where the value of the house is USD 125,000. On the other hand, if we err by this amount in Los Altos Hills, California, we can be proud of the accuracy of our model (the median house price there exceeds 4 million).

One way to address this problem is to measure the discrepancy in the logarithm of the price estimates. In fact, this is also the error that is being used to measure the quality in this competition. After all, a small value $\delta$ of $\log y - \log \hat{y}$ translates into $e^{-\delta} \leq \frac{\hat{y}}{y} \leq e^{\delta}$. This leads to the following loss function:

$$L = \sqrt{\frac{1}{n} \sum_{i=1}^{n} \left( \log y_i - \log \hat{y}_i \right)^2}$$

```
In [10]: def log_rmse(net, features, labels):
             # To further stabilize the value when the logarithm is taken, set the value less
             clipped_preds = nd.clip(net(features), 1, float('inf'))
             rmse = nd.sqrt(2 * loss(clipped_preds.log(), labels.log()).mean())
             return rmse.asscalar()
```

Unlike in the previous sections, the following training functions use the Adam optimization algorithm. Compared to the previously used mini-batch stochastic gradient descent, the Adam optimization algorithm is relatively less sensitive to learning rates. This will be covered in further detail later on when we discuss the details on Optimization Algorithms in a separate chapter.

```
In [11]: def train(net, train_features, train_labels, test_features, test_labels,
                   num_epochs, learning_rate, weight_decay, batch_size):
             train_ls, test_ls = [], []
             train_iter = gdata.DataLoader(gdata.ArrayDataset(
                 train_features, train_labels), batch_size, shuffle=True)
             # The Adam optimization algorithm is used here.
             trainer = gluon.Trainer(net.collect_params(), 'adam', {
                 'learning_rate': learning_rate, 'wd': weight_decay})
             for epoch in range(num_epochs):
                 for X, y in train_iter:
```

5

```
            with autograd.record():
                l = loss(net(X), y)
            l.backward()
            trainer.step(batch_size)
        train_ls.append(log_rmse(net, train_features, train_labels))
        if test_labels is not None:
            test_ls.append(log_rmse(net, test_features, test_labels))
    return train_ls, test_ls
```

## 1.4 k-Fold Cross-Validation

The k-fold cross-validation was introduced in the section where we discussed how to deal with
"Model Selection, Underfitting and Overfitting". We will put this to good use to select the model
design and to adjust the hyperparameters. We first need a function that returns the i-th fold of the
data in a k-fold cros-validation procedure. It proceeds by slicing out the i-th segment as validation
data and returning the rest as training data. Note - this is not the most efficient way of handling
data and we would use something much smarter if the amount of data was considerably larger.
But this would obscure the function of the code considerably and we thus omit it.

```
In [12]: def get_k_fold_data(k, i, X, y):
             assert k > 1
             fold_size = X.shape[0] // k
             X_train, y_train = None, None
             for j in range(k):
                 idx = slice(j * fold_size, (j + 1) * fold_size)
                 X_part, y_part = X[idx, :], y[idx]
                 if j == i:
                     X_valid, y_valid = X_part, y_part
                 elif X_train is None:
                     X_train, y_train = X_part, y_part
                 else:
                     X_train = nd.concat(X_train, X_part, dim=0)
                     y_train = nd.concat(y_train, y_part, dim=0)
             return X_train, y_train, X_valid, y_valid
```

The training and verification error averages are returned when we train $k$ times in the k-fold
cross-validation.

```
In [13]: def k_fold(k, X_train, y_train, num_epochs,
                    learning_rate, weight_decay, batch_size):
             train_l_sum, valid_l_sum = 0, 0
             for i in range(k):
                 data = get_k_fold_data(k, i, X_train, y_train)
                 net = get_net()

                 train_ls, valid_ls = train(net, *data, num_epochs, learning_rate,
                                            weight_decay, batch_size)
                 train_l_sum += train_ls[-1]
                 valid_l_sum += valid_ls[-1]
```

6

```
#            if i == 0:
#                d2l.semilogy(range(1, num_epochs + 1), train_ls, 'epochs', 'rmse',
#                             range(1, num_epochs + 1), valid_ls,
#                             ['train', 'valid'])
#        print('fold %d, train rmse: %f, valid rmse: %f' % (
#            i, train_ls[-1], valid_ls[-1]))
    return net, train_l_sum / k, valid_l_sum / k
```

## 1.5 Model Selection

We pick a rather un-tuned set of hyperparameters and leave it up to the reader to improve the model considerably. Finding a good choice can take quite some time, depending on how many things one wants to optimize over. Within reason the k-fold crossvalidation approach is resilient against multiple testing. However, if we were to try out an unreasonably large number of options it might fail since we might just get lucky on the validation split with a particular set of hyperparameters.

### 1.5.1 Custom Preprocessing

```
In [14]: from sklearn.cross_decomposition import CCA
         import matplotlib.pyplot as plt
         from sklearn.preprocessing import normalize
         from sklearn.model_selection import train_test_split
         from sklearn.metrics import r2_score
         import xgboost
```

```
In [15]: def check(X_train, y_train):
             xgb = xgboost.XGBRegressor(n_estimators=100, learning_rate=0.08, gamma=0, subsampl
                                        colsample_bytree=1, max_depth=7)

             X_train, X_test, y_train, y_test = train_test_split(X_train, y_train)

             xgb.fit(X_train, y_train)
             pred = xgb.predict(X_test)
             print('xgb:', r2_score(y_test, pred))

             return xgb

         def preprocess(dat, test=False, obj_features=None, num_features=None):

             data = dat.copy()

             if not test:
                 data = data[data['SalePrice'] < data['SalePrice'].quantile(.95)]
                 data = data[data['1stFlrSF'] < data['1stFlrSF'].quantile(.95)]
                 data = data[data['LotArea'] < data['LotArea'].quantile(.95)]

             data = data[data.dtypes[data.dtypes != 'object'].index]
```

7

```python
        if obj_features is not None:
            if not test:
                data = pd.concat([data, train_data.iloc[data.index][obj_features]], axis=
            else:
                data = pd.concat([data, test_data.iloc[data.index][obj_features]], axis=1)

        data['Bath'] = data['FullBath'] + (0.5 * data['HalfBath'])
        data['BsmtBath'] = data['BsmtFullBath'] + (0.5 * data['BsmtHalfBath'])

        data['TotalFlrSF'] = data['1stFlrSF'] + (0.75 * data['2ndFlrSF'])
        data = data.drop(['1stFlrSF', '2ndFlrSF'], axis=1)

        if num_features is not None:
            data = data.drop(num_features, axis=1)

        # id for test data
        ids = None

        if test:
            ids = data['Id']

        data = data.drop(['FullBath', 'HalfBath', 'BsmtFullBath', 'BsmtHalfBath', 'Id', '

        return data, ids

In [16]: train_data = pd.read_csv('kaggle_house_pred_train.csv')
         test_data = pd.read_csv('kaggle_house_pred_test.csv')

         obj_features = ['SaleCondition', 'MSZoning', 'BldgType', 'Neighborhood', 'LotConfig']
         num_features = ['GarageCars', 'Fireplaces', 'MSSubClass', 'LotFrontage', 'LotArea', '(
                 'OverallCond', 'YearBuilt', 'YearRemodAdd', 'MasVnrArea',
                 'BsmtFinSF1', 'BsmtUnfSF', 'TotalBsmtSF', 'GrLivArea',
                 'BedroomAbvGr', 'TotRmsAbvGrd', 'GarageYrBlt', 'GarageArea',
                 'WoodDeckSF', 'OpenPorchSF', 'MoSold', 'YrSold', 'TotalFlrSF']


         data, _ = preprocess(train_data, obj_features=obj_features)
         test, ids = preprocess(test_data, test=True, obj_features=obj_features)

         X_train, y_train = data.drop('SalePrice', axis=1), data['SalePrice']

         # Hot-encode
         X_train = pd.get_dummies(X_train, dummy_na=True)
         X_test = pd.get_dummies(test, dummy_na=True)

         # X_train = X_train.drop(num_features, axis=1)
```

```python
# Normalize
X_train = X_train.apply(lambda x: (x - x.mean()) / x.std()).fillna(0)
X_test = X_test.apply(lambda x: (x - x.mean()) / x.std()).fillna(0)

# X_train.head(1)

# check xgb accuracy for testing improvement of preprocessing
xgb = check(X_train, y_train)

tbl = pd.DataFrame(data={'score':xgb.feature_importances_, 'features':list(X_train)})

X_train = X_train[tbl[tbl['score'] >=.01]['features'].values]
X_test = X_test[tbl[tbl['score'] >=.01]['features'].values]

k, num_epochs, lr, weight_decay, batch_size = 5, 500, 5, 0.9, 128
params = {'k': k, 'num_epochs':num_epochs, 'lr':lr, 'weight_decay':weight_decay, 'batc
net, train_l, valid_l = k_fold(k, nd.array(X_train), nd.array(y_train), num_epochs, li
                         weight_decay, batch_size)

print('%d-fold validation: avg train rmse: %f, avg valid rmse: %f'
      % (k, train_l, valid_l))

X_test.shape, X_train.shape
```

```
xgb: 0.8984429740313958
5-fold validation: avg train rmse: 0.113716, avg valid rmse: 0.129685
```

```
Out[16]: ((1459, 21), (1251, 21))
```

**Submission**

```python
In [17]: pred = net(nd.array(X_test.as_matrix())).asnumpy()
         pred = pred.reshape(-1, )

         test_submission = pd.DataFrame(data={'Id':ids, 'SalePrice':pred})
         test_submission.to_csv('submission.csv', index=False)

         test_submission.shape
```

```
/home/hsong1101/miniconda3/envs/gluon/lib/python3.6/site-packages/ipykernel_launcher.py:1: Futu
  """Entry point for launching an IPython kernel.
```

```
Out[17]: (1459, 2)
```

You will notice that sometimes the number of training errors for a set of hyper-parameters can be very low, while the number of errors for the *K*-fold cross validation may be higher. This is most