#### 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 <u>Kaggle competition</u> (<a href="https://www.kaggle.com/c/house-prices-advanced-regression-techniques">https://www.kaggle.com/c/house-prices-advanced-regression-techniques</a>). 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:

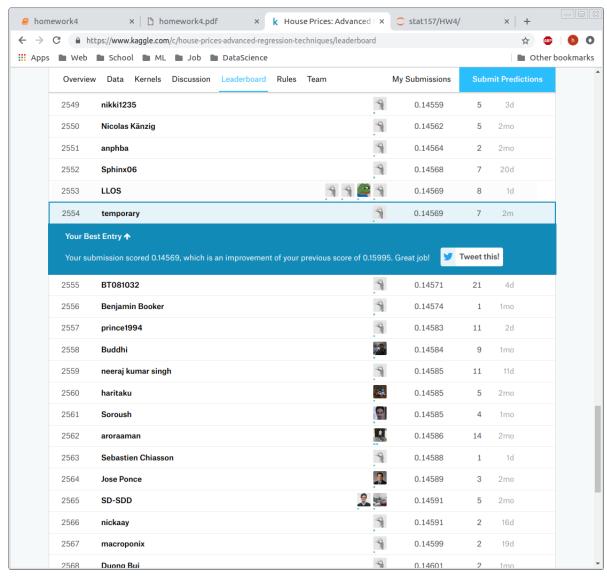
- 1. Start as earlier as possible. Though we will cover this notebook on Thursday's lecture, tuning hyper-parameters takes time, and Kaggle limits #submissions per day.
- 2. Work with your project teammates. It's a good opportunity to get familiar with each other.

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 <u>efficient data analysis toolkit</u> (<a href="http://pandas.pydata.org/pandas-docs/stable/">http://pandas.pydata.org/pandas-docs/stable/</a>). Make sure you have pandas installed for the experiments in this section.



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

executed in 813ms, finished 22:31:35 2019-02-18
```

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/kag
# utils.download('https://github.com/d2l-ai/d2l-en/raw/master/data/kag
train_data = pd.read_csv('kaggle_house_pred_train.csv')
test_data = pd.read_csv('kaggle_house_pred_test.csv')

executed in 25ms, finished 22:31:35 2019-02-18
```

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

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]] executed in 66ms, finished 22:31:35 2019-02-18
```

#### Out[4]:

	ld	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.

## 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:

\begin{equation}x \leftarrow \frac{x - \mu}{\sigma}\end{equation}

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'
    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 miss all_features = all_features.fillna(0)
    executed in 88ms, finished 22:31:35 2019-02-18
```

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, ar
    all_features = pd.get_dummies(all_features, dummy_na=True)
    all_features.shape
    executed in 62ms, finished 22:31:35 2019-02-18
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.

## 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()

v 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

executed in 38ms, finished 22:31:35 2019-02-18
```

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} \leg \frac{\hat{y}}{y} \leg e^\delta\$. This leads to the following loss function:

```
In [10]: 
    def log_rmse(net, features, labels):
        # To further stabilize the value when the logarithm is taken, set
        clipped_preds = nd.clip(net(features), 1, float('inf'))
        rmse = nd.sqrt(2 * loss(clipped_preds.log(), labels.log()).mean())
        return rmse.asscalar()
        executed in 49ms, finished 22:31:35 2019-02-18
```

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 <a href="Optimization Algorithms">Optimization Algorithms</a> (../chapter optimization/index.md) in a separate chapter.

```
In [11]: v def train(net, train features, train labels, test features, test label
                      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:
                        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
         executed in 67ms, finished 22:31:35 2019-02-18
```

#### 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" (underfit-overfit.md). 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
         executed in 48ms, finished 22:31:35 2019-02-18
```

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 ra
                                                 weight decay, batch size)
                    train l sum += train ls[-1]
                    valid l sum += valid ls[-1]
                      if i == 0:
                          d2l.semilogy(range(1, num epochs + 1), train ls, 'epochs
                                       range(1, num\_epochs + 1), valid \overline{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
          executed in 48ms, finished 22:31:35 2019-02-18
```

#### ▼ 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
    executed in 98ms, finished 22:31:35 2019-02-18
```

```
In [15]: ▼ def check(X train, y train):
               xgb = xgboost.XGBRegressor(n estimators=100, learning rate=0.08, d
                                            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 = xqb.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</pre>
                   data = data[data['1stFlrSF'] < data['1stFlrSF'].quantile(.95)]</pre>
                    data = data[data['LotArea'] < data['LotArea'].guantile(.95)]</pre>
               data = data[data.dtypes[data.dtypes != 'object'].index]
               if obj features is not None:
                    if not test:
                        data = pd.concat([data, train data.iloc[data.index][obj fe
                    else:
                        data = pd.concat([data, test data.iloc[data.index][obj feat
               data['Bath'] = data['FullBath'] + (0.5 * data['HalfBath'])
               data['BsmtBath'] = data['BsmtFullBath'] + (0.5 * data['BsmtHalfBat
               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', 'BsmtHal
               return data, ids
         executed in 25ms, finished 22:31:35 2019-02-18
```

```
train_data = pd.read_csv('kaggle house pred train.csv')
In [16]:
            test data = pd.read csv('kaggle house pred test.csv')
           obj_features = ['SaleCondition', 'MSZoning', 'BldgType', 'Neighborhood
num_features = ['GarageCars', 'Fireplaces', 'MSSubClass', 'LotFrontage
                    '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)
            # Normalize
            X_{train} = X_{train.apply}(lambda x: (x - x.mean()) / x.std()).fillna(0)
            X \text{ test} = X \text{ 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':
            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':wei
           net, train l, valid l = k fold(k, nd.array(X train), nd.array(y train)
                                         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
          executed in 59.8s, finished 22:32:35 2019-02-18
```

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

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

#### Submission

/home/hsongl101/miniconda3/envs/gluon/lib/python3.6/site-packages/ipyk ernel\_launcher.py:1: FutureWarning: Method .as\_matrix will be removed in a future version. Use .values instead.

"""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 likely a consequence of overfitting. Therefore, when we reduce the amount of training errors, we need to check whether the amount of errors in the k-fold cross-validation have also been reduced accordingly.

#### 1.6 Predict and Submit

Now that we know what a good choice of hyperparameters should be, we might as well use all the data to train on it (rather than just \$1-1/k\$ of the data that is used in the crossvalidation slices). The model that we obtain in this way can then be applied to the test set. Saving the estimates in a CSV file will simplify uploading the results to Kaggle.

Let's invoke the model. A good sanity check is to see whether the predictions on the test set resemble those of the k-fold crossvalication process. If they do, it's time to upload them to Kaggle.

A file, submission.csv will be generated by the code above (CSV is one of the file formats accepted by Kaggle). Next, we can submit our predictions on Kaggle and compare them to the actual house price (label) on the testing data set, checking for errors. The steps are quite simple:

- Log in to the Kaggle website and visit the House Price Prediction Competition page.
- · Click the "Submit Predictions" or "Late Submission" button on the right.
- Click the "Upload Submission File" button in the dashed box at the bottom of the page and select the prediction file you wish to upload.
- Click the "Make Submission" button at the bottom of the page to view your results.

#### 1.7 Hints

- 1. Can you improve your model by minimizing the log-price directly? What happens if you try to predict the log price rather than the price?
- 2. Is it always a good idea to replace missing values by their mean? Hint can you construct a situation where the values are not missing at random?
- 3. Find a better representation to deal with missing values. Hint What happens if you add an indicator variable?
- 4. Improve the score on Kaggle by tuning the hyperparameters through k-fold crossvalidation.
- 5. Improve the score by improving the model (layers, regularization, dropout).
- 6. What happens if we do not standardize the continuous numerical features like we have done in this section?

Note for converting this notebook into PDF. If you use 'File -> Download as -> PDF', you may get the error that svg cannot converted because inkscape is not installed and cannot find PNG images. The easiest way is printing this notebook as a PDF in your browser. Or, you can install inkscape to convert SVG (On macOS, you may brew cask install xquartz inkscape, on Ubuntu, you may sudo apt-get install inkscape) and change the image URL to local filenames.