

0. Kuo Zhang

1. There was problem during class.  
Why GD Gradient descent objective value always decreasing? While under proper conditions?

First we have to make an assumption function  
 $V_t, F(w^{t+1}) - F(w^t) < 0 \Leftrightarrow F(w^{t+1}) \leq F(w^t)$   
In this function  $F(w^{t+1})$  as next rate  $F(w^t)$  is current.  
Then we have to proof this function to proof this assumption is right:

$$w^{t+1} = w^t - \eta \cdot \nabla F(w^t)$$

$$F(w^{t+1}) - F(w^t) = \langle \nabla F(w), w^{t+1} - w^t \rangle$$

Apply Mean value theorem:

$$F: \mathbb{R} \rightarrow \mathbb{R}$$

$$\forall x, y \in \mathbb{R}$$

$$F(x) - F(y) = F'(z) \cdot (x - y)$$

$$z \in [x, y]$$

$$\longrightarrow$$

$$F: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$\forall x, y \in \mathbb{R}^d$$

$$F(x) - F(y)$$

$$= \langle \nabla F(z), x - y \rangle$$

Now, we have Function:

$$F(w^{t+1}) - F(w^t) = \langle \nabla F(w), w^{t+1} - w^t \rangle$$

$$= \langle \nabla F(w), -\eta \cdot \nabla F(w^t) \rangle$$

$$= -\eta \cdot \langle \nabla F(w), \nabla F(w^t) \rangle$$

So we know that,

$$\eta \rightarrow 0 \text{ or } \approx 0$$

$$w^t \quad w \quad w^{t+1}$$

If line segment goes to zero size,  $w$  must be very close to  $w^t$

So  $\nabla F(w) \approx \nabla F(w^t)$   
It's in same vector. is positive but, overall is negative,  $-\eta$ .  
Hence, GD Gradient descent

always decrease objective function

2. Random projection and principal component analysis are two methods that can be used to reduce the dimensionality of a dataset. Both methods are useful for making data easier to work with, but they are different in effectiveness at different types of data.

1) The random projection would outperform PCA when analyzing a large-scale dataset, while PCA would outperform the random projection when the dimension and sample size are relatively small. This is because the random projection has the time complexity of  $O(npd)$  and the complexity of PCA is  $O(np^2 + p^3)$ .

2) In mixture of Gaussians, random projection is just reducing the dimension.

3) PCA has to hold data in memory for large data, but random projection doesn't need to do so.

4) Random projection will not get worse after initial training data. But, PCA will provide the best projection, and it will get worse after that.

Conclusion:

Random projection projects data at lower dimensional space via randomly generated matrix, good for reducing high dimension, but will lose accuracy and information.

PCA is accurate for dimension reduction, but it can be computationally intensive for large data sets.



### 3. Example active learning lead to exponential saving:

For active learning, we can choose to use a crowdsourcing-like data set but do not have to label the entire data set. The active learning algorithm iteratively selects data based on some metric and sends this unlabeled data to the authority which then labels it and returns it to algorithm.

For example:

the MNIST experiments ( $\text{train-num} = 55000$ ,  $\text{val-num} = 10000$ ) using the AlexNet model Pytorch framework: 1) using the full training data to directly train the model 120 times epoch,  $\text{val-acc} = 98.992\%$

2) Using Uncertainty Strategy only 2300 labeled data  $\text{val-acc}$  needed to reach  $99.04\%$  and the remaining 52700 and into trained model for prediction and  $99.70\%$  ( $52543 / 52700$ ) is obtained.

It can be seen that for the MNIST data set. Significant results can be obtained using active learning. Example when active learning does not:

Data selected using active learning may have a large number of duplicates, thus causing annotation redundancy problems. Typically, the information content of unlabeled points is evaluated separately. In active learning we use BALD a collection function method to evaluate unlabeled content. But since individual information points may be almost identical evaluating information at each point separately is wasting resources.

For example if we simply acquire the top  $K$  most useful points, we may end up with having experts label  $K$  almost identical points.

#### 4. algorithm

- 1) We can use KNN. To be specific for each observation with missing value, we should compute the distances between this observation and all others without missing value based on the features except Blood Pressure. Then find the  $K$  observations corresponding to the smallest  $K$  distances and calculate the average value of their Blood Pressure. Finally use this value to fill '?'.

Let the indices corresponding to missing values be  $i_1, \dots, i_N$ .

Algorithm:

Choose the value of  $K$ .

For  $i = i_1, \dots, i_N$

(1) Compute  $d_j = \|x_i - x_j\|_2$  for  $j \in \{1, \dots, n\} \setminus \{i_1, \dots, i_N\}$ .

where  $x_i$  is the vector of features without Blood pressure for the  $i$ th observation.

(2) Find the  $K$  indices corresponding to the smallest  $K$  distances:  $j_1, \dots, j_K$ .

(3) Fill the missing value by  $\frac{1}{K} \sum_{j=1}^K y_{j_i}$  where  $y_j$  is the value of Blood Pressure for  $j$ th observation.



4.

2) If we know which observed entries are corrupted we can just remove them and fill the table using the algorithm from question 1).

However, if we do not know which observed entries are corrupted, we should first find them,

Use the following procedure: T: corrupted data

1. Compute empirical mean and covariance  $\mu_T, \Sigma_T$   
 Compute largest eigenvalue  $\lambda^*$  of  $\Sigma_T^{-1}$  and, eigenvector  $v^*$

If  $\lambda^*$  is small, return  $\mu_T$

Otherwise, find  $t > C_1$  such that

$$\Pr_{X \in T} \left[ |v^* \cdot (X - \mu_T)| > t \right] \geq C_2 e^{-t^2/2}$$

$$+ \frac{C_3 \epsilon}{t^2 \log(n/\epsilon T)}$$

then Remove  $X$  such that  $|v^* \cdot (X - \mu_T)| > t$   
 Return to first step.

5. ① A random forest regression model is built to predict each missing value using several unmissing values. First, let us convert the categorical variables into classification codes that can be handled by the random forest model, using first Age as an example

```
df['Age - category'] = df['Age'].astype('category')  
df['Age - category'] = df['Age - category'].cat.codes
```

② Increasing the random sample size:

```
df = pd.read_csv('xxx.csv').sample(n=50, random_state=42)
```

③ Import the random forest regressor module from Scikit learn. Define the list of features used to train the model:

```
from sklearn.ensemble import RandomForestRegressor  
features = ['Age - category', 'Weight - category', 'Height - category', ...]
```

④ Train our model with a random forest 100 estimates and a maximum depth of 100. then generate the predictions, add to new list

```
for train_index, test_index in Kf.split(df - filter):
```

```
    df_test = df - filter.iloc[test_index]
```

```
    df_train = df - filter.iloc[train_index]
```

```
    X_train = np.array(df_train[features])
```

```
    y_train = np.array(df_train['missing value'])
```

```
    X_test = np.array(df_test[features])
```

```
    y_test = np.array(df_test['missing value'])
```



5-1

model = Random Forest Regressor (n-estimators = 100,  
max-depth = 100, random-state = 42)

model.fit(X-train, y-train)

y-pred-rf.append(model.predict(X-test)[0])

y-true-rf.append(y-test[0])

⑤ mean-squared-error (y-pred-rf, y-true-rf)

⑥ Predicting missing values

model\_rf = Random Forest Regressor (n-estimators =  
100, max-depth = 100, random-state = 42)

model\_rf.fit(X-train\_rf, y-train\_rf)

model\_rf.predict(X-test\_rf)[0]