

## 实验一：逻辑回归实现糖尿病预测

朱云沁 PB202061372

### Introduction

Suppose that we are going to solve a binary classification problem given a set of data  $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , where the features  $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$  and the class label  $y \in \{0, 1\}$  are generated from probability distributions. Following the minimum error rate principle, we can predict the class label as follows:

$$\hat{y} = \begin{cases} 1 & \text{if } P(y=1|\mathbf{x}) > P(y=0|\mathbf{x}) \\ 0 & \text{if } P(y=1|\mathbf{x}) < P(y=0|\mathbf{x}) \end{cases}$$

As a commonly-used linear classifier, the logistic regression model simply assumes that the decision boundary  $\{\mathbf{x} \in \mathbb{R}^d : P(y=1|\mathbf{x}) = P(y=0|\mathbf{x})\}$  is a hyperplane  $(\mathbf{w}, b)$  in the feature space by modeling the following class posterior probability:

$$P(y=1|\mathbf{x}) = \frac{1}{1 + \exp(-(\mathbf{w}^T \mathbf{x} + b))} = \sigma(\mathbf{w}^T \mathbf{x} + b) \quad (1)$$

where  $\sigma(\cdot)$  is the sigmoid function. Then, our prediction rule is equivalent to

$$\hat{y} = \begin{cases} 1 & \text{if } \mathbf{w}^T \mathbf{x} + b > 0 \\ 0 & \text{if } \mathbf{w}^T \mathbf{x} + b < 0 \end{cases}$$

### A Probabilistic Derivation

To see the rationality behind this model, we first rewrite the posterior probability in Eq. (1) by the Bayes rule

$$P(y=1|\mathbf{x}) = \frac{p(\mathbf{x}|y=1)P(y=1)}{p(\mathbf{x}|y=1)P(y=1) + p(\mathbf{x}|y=0)P(y=0)} = \frac{1}{1 + \frac{p(\mathbf{x}|y=0)P(y=0)}{p(\mathbf{x}|y=1)P(y=1)}} \quad (2)$$

We model the class probabilities by a suite of assumptions as follows.

- $y \sim \text{Berni}(p)$ , i.e. the class prior follows a Bernoulli distribution with  $P(y=1) = p$  and  $P(y=0) = 1 - p$ .
- $x_j | y \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu_j | y, \sigma_j^2)$  and  $x_j | y \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu_{j0}, \sigma_j^2)$ ,  $j = 1, 2, \dots, d$ , i.e. the features are conditionally independent given the class label, and the class-conditional distributions are both Gaussian. Note that the variance  $\sigma_j^2$  is shared by both classes.
- We simplify the notation by writing  $\mathbf{x}_{j=1} \sim \mathcal{N}(\mu_1, \Sigma)$  and  $\mathbf{x}_{j=0} \sim \mathcal{N}(\mu_0, \Sigma)$ .

We can now continue our derivation from where we left in Eq. (2).

$$P(y=1|\mathbf{x}) = \frac{1}{1 + \exp\left(\log \frac{p(\mathbf{x}|y=0)P(y=0)}{p(\mathbf{x}|y=1)P(y=1)}\right)} = \frac{1}{1 + \exp\left(\log \frac{p(\mu_{j0}|\mathbf{x})}{p(\mu_{j1}|\mathbf{x})} + \log \frac{1-p}{p}\right)} \quad (3)$$

By substituting the Gaussian class-conditional distributions, we have the log likelihood ratio

$$\log \frac{p(\mathbf{x}|y=1)}{p(\mathbf{x}|y=0)} = \log \frac{\mathcal{N}(\mathbf{x}; \mu_1, \Sigma)}{\mathcal{N}(\mathbf{x}; \mu_0, \Sigma)} = \Sigma^{-1}(\mu_1 - \mu_0)^T \mathbf{x} + \frac{1}{2}(\mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1) \quad (4)$$

Plugging Eq. (4) into Eq. (3), we obtain the posterior in the form of Eq. (1) with

$$\mathbf{w} = \Sigma^{-1}(\mu_1 - \mu_0) \\ b = \frac{1}{2}(\mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1) + \log \frac{1-p}{1-p} \quad (5)$$

In a nutshell, the logistic regression model makes predictions with the minimum error rate principle and a linear decision boundary by assuming that the log likelihood ratio is linear in  $\mathbf{x}$ .

### Learning the Parameters

As the parameter set of the involved distributions are usually unknown and our assumptions of the Gaussian likelihoods may not hold in many real applications, we cannot directly apply Eq. (5) and Eq. (1) to predict the label of a new data instance. Instead, we need to learn the parameters  $\mathbf{w}$  and  $b$  from the training data, via the maximum likelihood estimation (MLE).

$$\begin{aligned} \hat{\mathbf{w}}, \hat{b} &= \underset{\mathbf{w}, b}{\operatorname{argmax}} \prod_{i=1}^n P(y_i | \mathbf{x}_i; \mathbf{w}, b) = \underset{\mathbf{w}, b}{\operatorname{argmax}} \sum_{i=1}^n \log P(y_i | \mathbf{x}_i; \mathbf{w}, b) \\ &= \underset{\mathbf{w}, b}{\operatorname{argmax}} \sum_{i=1}^n (y_i \log P(y_i=1 | \mathbf{x}_i; \mathbf{w}, b) + (1 - y_i) \log P(y_i=0 | \mathbf{x}_i; \mathbf{w}, b)) \\ &= \underset{\mathbf{w}, b}{\operatorname{argmax}} \sum_{i=1}^n (y_i \log \sigma(\mathbf{w}^T \mathbf{x}_i + b) + (1 - y_i) \log(1 - \sigma(\mathbf{w}^T \mathbf{x}_i + b))) \\ &= \underset{\mathbf{w}, b}{\operatorname{argmin}} L(\mathbf{w}, b) \end{aligned} \quad (6)$$

where we let  $L(\mathbf{w}, b) = -\sum_{i=1}^n (y_i \log \sigma(\mathbf{w}^T \mathbf{x}_i + b) + (1 - y_i) \log(1 - \sigma(\mathbf{w}^T \mathbf{x}_i + b)))$  for notational simplicity. This convex optimization problem can be viewed equivalently as minimizing the cross entropy or the Kullback-Leibler divergence between the empirical distribution of the labels and the predicted distribution by our model.

It can be proved that Eq. (6) admits a solution when the data is NOT linearly separable. Alternatively, we can add a regularization term to ensure the existence of a solution. To keep things simple, we will not discuss this issue in this experiment.

### Gradient Descent

To apply gradient descent to find the optimal parameters, we first compute the gradient of the loss  $L(\mathbf{w}, b)$  with respect to  $\mathbf{w}$  and  $b$ . Note that  $\sigma'(x) = \sigma(x)(1 - \sigma(x))$ . Therefore,

$$\begin{aligned} \nabla_{\mathbf{w}} L(\mathbf{w}, b) &= -\sum_{i=1}^n (y_i(1 - \sigma(\mathbf{w}^T \mathbf{x}_i + b)) - (1 - y_i)\sigma(\mathbf{w}^T \mathbf{x}_i + b)) \mathbf{x}_i \\ &= -\sum_{i=1}^n (y_i - \sigma(\mathbf{w}^T \mathbf{x}_i + b)) \mathbf{x}_i \\ \nabla_b L(\mathbf{w}, b) &= -\sum_{i=1}^n (y_i(1 - \sigma(\mathbf{w}^T \mathbf{x}_i + b)) - (1 - y_i)\sigma(\mathbf{w}^T \mathbf{x}_i + b)) \\ &= -\sum_{i=1}^n (y_i - \sigma(\mathbf{w}^T \mathbf{x}_i + b)) \end{aligned}$$

The update rule of gradient descent is then given by

$$\begin{aligned} \mathbf{w}^{t+1} &\leftarrow \mathbf{w}^t - \alpha \nabla_{\mathbf{w}} L(\mathbf{w}^t, b^t) \\ b^{t+1} &\leftarrow b^t - \alpha \nabla_b L(\mathbf{w}^t, b^t) \end{aligned} \quad (7)$$

where  $\alpha$  is a hyperparameter denoting the learning rate.

### Newton's Method

The Newton's method is an iterative algorithm for optimization problems that uses the second-order derivative to update the parameters. The Hessian matrix of  $L(\mathbf{w}, b)$  can be computed as

$$\begin{aligned} \nabla_{\mathbf{w}}^2 L(\mathbf{w}, b) &= \sum_{i=1}^n \sigma(\mathbf{w}^T \mathbf{x}_i + b)(1 - \sigma(\mathbf{w}^T \mathbf{x}_i + b)) \mathbf{x}_i \mathbf{x}_i^T \\ \nabla_b^2 L(\mathbf{w}, b) &= \sum_{i=1}^n \sigma(\mathbf{w}^T \mathbf{x}_i + b)(1 - \sigma(\mathbf{w}^T \mathbf{x}_i + b)) \end{aligned}$$

The update rule of Newton's method is then given by

$$\begin{aligned} \mathbf{w}^{t+1} &\leftarrow \mathbf{w}^t - \alpha (\nabla_{\mathbf{w}}^2 L(\mathbf{w}^t, b^t))^{-1} \nabla_{\mathbf{w}} L(\mathbf{w}^t, b^t) \\ b^{t+1} &\leftarrow b^t - \alpha (\nabla_b^2 L(\mathbf{w}^t, b^t))^{-1} \nabla_b L(\mathbf{w}^t, b^t) \end{aligned} \quad (8)$$

It can be shown that the Newton's method converges quadratically to the optimal solution if the Hessian matrix is positive definite (which can be achieved by regularization) and the initial point is sufficiently close to the optimal solution.

However, the computation of the Hessian matrix is expensive for large-scale problems. In real-world applications, some approximations are often used to reduce the computational cost, leading to quasi-Newton methods such as the limited-memory BFGS (L-BFGS) method. In this experiment, we will use the standard Newton's method for simplicity.

## Experiment

### Initialization

```
In [ ]: import os
import numpy as np
import pandas as pd
import tensorflow as tf
import matplotlib.pyplot as plt
from tqdm import tqdm
```

```
os.environ['TF_CPP_MIN_LOG_LEVEL'] = '2'
```

### Dataset Description

The dataset we use in this experiment is the [Pima Indians Diabetes Database](#) from Kaggle. The dataset consists of eight medical predictor variables and one target variable, i.e., whether the patient has diabetes. The detailed description of each column is listed below.

- Feature columns (0-7):
  - Pregnancies: Number of times pregnant
  - Glucose: Plasma glucose concentration a 2 hours in an oral glucose tolerance test
  - BloodPressure: Diastolic blood pressure (mmHg)
  - SkinThickness: Triceps skin fold thickness (mm)
  - Insulin: 2-Hour serum insulin (uU/ml)
  - BMI: Body mass index (kg/m<sup>2</sup>)
  - DiabetesPedigreeFunction: Diabetes pedigree function
  - Age: Patient's age in years
- Label column (8):
  - Outcome: 0 for healthy or 1 for diabetes

```
In [ ]: df = pd.read_csv('diabetes.csv')
df.describe()
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000
mean	3.845562	120.894531	69.105469	20.536458	79.799479	31.992578	0.4716
std	3.365578	31.972618	19.355687	15.952218	115.244002	7.884160	0.3313
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0780
25%	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	0.2437
50%	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	0.3725
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	0.6262
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	2.4200

### Imputing Missing Values

Observing that the minimum values of six features (i.e. Glucose, BloodPressure, SkinThickness, Insulin, BMI, and Age) are zeros, we can infer that the missing values are encoded as zeros. Common strategies for handling missing values include dropping the corresponding samples or imputing them with the mean or median of the corresponding feature values. While the former strategy may lead to insufficient training data, the latter strategy may introduce bias to the resulting model. Here, we choose to replace the missing values with the corresponding median values.

```
In [ ]: impute_cols = [
    'Pregnancies',
    'Glucose',
    'BloodPressure',
    'SkinThickness',
    'Insulin',
    'BMI',
    'Age']
```

```
df[impute_cols] = df[impute_cols].replace(0, np.NaN)
df.fillna(df.median(), inplace=True)
df.describe()
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000
mean	4.423177	121.656250	72.386719	29.108073	140.671875	32.455208	0.4716
std	2.980481	30.438286	12.096642	8.791221	86.383060	6.875177	0.3313
min	0.000000	44.000000	24.000000	7.000000	14.000000	18.200000	0.0780
25%	2.000000	99.750000	64.000000	25.000000	121.500000	27.500000	0.2437
50%	4.000000	117.000000	72.000000	29.000000	125.000000	32.300000	0.3725
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	0.6262
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	2.4200

### Standardization

To avoid the features with larger values dominating the training process, we standardize each feature by subtracting its mean and dividing by its standard deviation.

```
In [ ]: from sklearn.preprocessing import StandardScaler

feature_cols = [
    'Pregnancies',
    'Glucose',
    'BloodPressure',
    'SkinThickness',
    'Insulin',
    'BMI',
    'DiabetesPedigreeFunction',
    'Age',
]
```

```
scaler = StandardScaler()
df[feature_cols] = scaler.fit_transform(df[feature_cols])
df.mean(), df.std()
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000
mean	4.423177	121.656250	72.386719	29.108073	140.671875	32.455208	0.4716
std	2.980481	30.438286	12.096642	8.791221	86.383060	6.875177	0.3313
min	0.000000	44.000000	24.000000	7.000000	14.000000	18.200000	0.0780
25%	2.000000	99.750000	64.000000	25.000000	121.500000	27.500000	0.2437
50%	4.000000	117.000000	72.000000	29.000000	125.000000	32.300000	0.3725
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	0.6262
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	2.4200

### Splitting into K Folds

To evaluate the performance of the trained model, we split the dataset into 5 folds and use one fold for validation and the remaining folds for training in each iteration. The final performance is reported as the average of the performance on the 5 validation sets.

```
In [ ]: from sklearn.model_selection import StratifiedKFold

kf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
```

```
df['fold'] = -1
for fold, (train_idx, val_idx) in enumerate(kf.split(df[feature_cols], df['Outcome'])):
    df.loc[train_idx, 'fold'] = fold
    df.loc[val_idx, 'fold'] = value
```

```
df.groupby('Outcome')['fold'].value_counts()
```

	Outcome	fold
0	0	100
		100
		2
		3
		100
1	0	54
		54
		2
		54
		3
	1	53
		53
		4
		53
		53

Name: count, dtype: int64

### Handling Imbalanced Data

Since the number of healthy samples is much larger than the number of diabetes samples in the dataset, the model may be biased towards the healthy class. To address this issue, we assign a larger weight to the diabetes samples in the loss function to make them more important during training. In this experiment, we treat the weight as a hyperparameter.

Alternative approaches include oversampling, generating synthetic samples or using a different loss function such as the Focal Loss.

### Using the Scikit-Learn Library

For convenience of checking the correctness of our implementation, we compare the results with those obtained by the LogisticRegression class in the Scikit-Learn library.

```
In [ ]: from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, f1_score, recall_score, precision_score
```

```
accuracy_scores = []
f1_scores = []
recall_scores = []
precision_scores = []

def get_fold(fold):
    X_train = df[df['fold'] == fold][feature_cols].values
    y_train = df[df['fold'] == fold]['Outcome'].values
    X_val = df[df['fold'] == fold][feature_cols].values
    y_val = df[df['fold'] == fold]['Outcome'].values
    return X_train, y_train, X_val, y_val
```

```
def evaluate_fold(y_true, y_pred, fold=0, print_results=True):
    accuracy = accuracy_score(y_true, y_pred)
    f1 = f1_score(y_true, y_pred)
    recall = recall_score(y_true, y_pred)
    precision = precision_score(y_true, y_pred)
    if print_results:
        print(f'Fold {fold} Accuracy: {accuracy}')
        print(f'Fold {fold} F1 Score: {f1}')
        print(f'Fold {fold} Recall: {recall}')
        print(f'Fold {fold} Precision: {precision}')
    return accuracy, f1, recall, precision
```

```
for fold in range(5):
    X_train, y_train, X_val, y_val = get_fold(fold)
    # model = LogisticRegression(penalty=None, class_weight='balanced')
    model = LogisticRegression(penalty=None)
    model.fit(X_train, y_train)
    y_pred = model.predict(X_val)
    accuracy, f1, recall, precision = evaluate_fold(y_val, y_pred, fold)
    accuracy_scores.append(accuracy)
    f1_scores.append(f1)
    recall_scores.append(recall)
    precision_scores.append(precision)
```

```
print(f'Accuracy: {np.mean(accuracy_scores)}')
print(f'F1 Score: {np.mean(f1_scores)}')
print(f'Recall: {np.mean(recall_scores)}')
print(f'Precision: {np.mean(precision_scores)}')
print()

# Weights of last fold
print(f"Coefficients: {model.coef_.flatten()}")
print(f"Intercept: {model.intercept_[0]}")
```

```
Fold 0 Accuracy: 0.7662337662337663
Fold 0 F1 Score: 0.6399999999999999
Fold 0 Recall: 0.6929292929292929
Fold 0 Precision: 0.6956521739130435

Fold 1 Accuracy: 0.8851948051948052
Fold 1 F1 Score: 0.6666666666666667
Fold 1 Recall: 0.5555555555555556
Fold 1 Precision: 0.8333333333333333
```

```
Fold 2 Accuracy: 0.7792877928779288
Fold 2 F1 Score: 0.6399999999999999
Fold 2 Recall: 0.5378378378378371
Fold 2 Precision: 0.7631578947368421

Fold 3 Accuracy: 0.7643758823529411
Fold 3 F1 Score: 0.617821275957447
Fold 3 Recall: 0.5471698113207547
Fold 3 Precision: 0.7673176731767317
```

```
Fold 4 Accuracy: 0.745988392156863
Fold 4 F1 Score: 0.6422818348623852
Fold 4 Recall: 0.668377358490566
Fold 4 Precision: 0.625

Accuracy: 0.772986544435957
F1 Score: 0.6729649121466984
Recall: 0.5785464789993912
Precision: 0.7248928936393482
```

```
Coefficients: [ 0.35864773  1.18884361 -0.03436672  0.02858375 -0.13285964  0.66320468
 0.36212322  0.23940550]
Intercept: -0.8898873751296423
```

### Implementation of Gradient Descent with TensorFlow

This implementation is deprecated due to performance issues.

```
In [ ]: class LRGradientDescent(tf.keras.Model):
    ...
    def __init__(self, n_features):
        super().__init__()
        self.w = tf.Variable(tf.zeros(n_features, 1))
        self.b = tf.Variable(tf.zeros(1, 1))
        self.alpha = 0.01 # Learning Rate
        self.class_weight = 1.0 # Weight of positive class
        self.compile()
```

```
def call(self, X):
    z = tf.matmul(X, self.w) + self.b
    return tf.sigmoid(z)

def predict(self, X):
    return self(X) > 0.5

def train_step(self, data):
    X, y = data
    y = tf.cast(tf.reshape(y, (-1, 1)), tf.float32)
    # w: GradientTape
    # with tf.GradientTape() as tape:
        y_pred = self(X)
        loss = tf.reduce_mean(
            -self.class_weight * y * tf.math.log(y_pred)
            - (1 - y) * tf.math.log(1 - y_pred)
        )
        # grads = tape.gradient(loss, [self.w, self.b])
        self.w.assign_sub(self.alpha * grads[0])
        self.b.assign_sub(self.alpha * grads[1])
    # w/o GradientTape
    loss = tf.reduce_mean(
        -self.class_weight * y * tf.math.log(y_pred)
        - (1 - y) * tf.math.log(1 - y_pred)
    )
    grad_w = tf.reduce_mean(
        -self.class_weight * y * (1 - y_pred) + (1 - y) * y_pred * X, axis=0
    )
    grad_w = tf.reshape(grad_w, (-1, 1))
    grad_b = tf.reduce_mean(
        -self.class_weight * y * (1 - y_pred) + (1 - y) * y_pred, axis=0
    )
    self.w.assign_sub(self.alpha * grad_w)
    self.b.assign_sub(self.alpha * grad_b)
    return {'loss': loss}
```

```
def test_step(self, data):
    X, y = data
    y = tf.cast(tf.reshape(y, (-1, 1)), tf.float32)
    y_true = tf.cast(y, tf.bool)
    pred = self(X)
    cross_entropy = tf.reduce_mean(
        -y * tf.math.log(y_pred) - (1 - y) * tf.math.log(1 - y_pred)
    )
    accuracy = tf.reduce_mean(tf.cast(tf.equal(y_true, y_pred > 0.5), tf.float32))
    return {'cross_entropy': cross_entropy, 'accuracy': accuracy}
```

```
X_train, y_train, X_val, y_val = get_fold(4)
model = LRGradientDescent(len(feature_cols))
model.fit(
    X_train,
    y_train,
    validation_data=(
        X_val,
        y_val,
    ),
    epochs=500,
    verbose=0,
)
```

```
y_true = y_val
y_pred = model.predict(X_val)
```

```
evaluate(y_true, y_pred)
print(f"Coefficients: {model.w.numpy().flatten()}")
print(f"Intercept: {model.b.numpy()[0]}")
```

```
logs = model.history.history
plt.figure(figsize=(6, 4))
plt.plot(logs['val_loss'], color='r', linestyle='-')
plt.xlabel('Epochs')
plt.ylabel('Loss')
plt.tight_layout()
```

```
plt.figure(figsize=(6, 4))
plt.plot(logs['val_loss'], color='r', linestyle='-')
plt.xlabel('Epochs')
plt.ylabel('Accuracy')
plt.tight_layout()
```

```
Fold 4 Accuracy: 0.745988392156863
Fold 4 F1 Score: 0.6422818348623852
Fold 4 Recall: 0.668377358490566
Fold 4 Precision: 0.625

Coefficients: [ 0.35847787  1.16684332 -0.0335982  0.02766889 -0.13026531  0.65434305
 0.36173701  0.23641544 -0.08048319]
Intercept: -0.8804831947899048
```

