

Gradient decent

Exact Gradient Computation

Given a function f , sometimes we can compute its exact gradient at any x if f 's derivative is easy to compute. For example, let $f(x) = 2x^2 - 3x + \ln x$, where $x > 0$. Please compute the derivative of f and report its gradient at $x = 2$.

Your answer:

Numerical Gradient Computation [5 pts]

Instead of computing the derivative of a function, we can also estimate the gradient numerically with various methods. These methods are essential, especially when a callable function is not exposed due to privacy reasons, or it is hard to differentiate analytically.

To numerically compute the gradient, the simple way is to follow Newton's difference quotient: $f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$. Another two-point formula is to compute the slope through the points $(x - h, f(x - h))$ and $(x + h, f(x + h))$. Let us reuse the example function $f(x) = 2x^2 - 3x + \ln x$ and test the precision of these two approaches. Define the function in the next cell, and try to compute its gradient via both methods at $x = 2$. Range h value in $[0.1, 0.01, 0.001, 0.0001]$ and report all gradients calculated. Which method is more accurate, and why does it work better?

```
In [455... import math

def f(x):
    return 2*x**2 - 3 * x + math.log(x)
```

```
In [456... # Compute gradient using the first method (Newton's difference quotient)

h = [0.1, 0.01, 0.001, 0.0001]
value = []
for val in h:
    one = f(2 + val)
    two = f(2)
    numerator = one - two
    denominator = val
    final_answer = numerator/denominator
    value.append(final_answer)
print(value)

[5.687901641694317, 5.518754151103744, 5.50187504165045, 5.5001875004201395]
```

```
In [457... # Compute gradient using the second method
value_two = []
for val in h:
    x_one = 2-val
    y_one = f(x_one)

    x_two = 2 + val
    y_two = f(x_two)
    num = y_two - y_one
    den = x_two - x_one
    answer = num/den
    value_two.append(answer)
print(value_two)

[5.500417292784904, 5.500004166729123, 5.500000041666448, 5.500000000412448]
```

Remark: You may find the gradient more accurate when using a smaller h value. However, this is not always the case. Due to the finite precision of the floating-point, rounding errors always exist and can dominate the computation when the h value is too small. Run the following two cells to observe such scenarios.

```
In [458... eps = 1e-15
print((f(2+eps)-f(2-eps))/2/eps)

5.551115123125783
```

```
In [459... eps = 1e-20
print((f(2+eps)-f(2-eps))/2/eps)

0.0
```

```
In [460... print("Calculating the slope gives a more accurate answer than Newton's diff
```

Calculating the slope gives a more accurate answer than Newton's difference quotient because the value becomes less accurate as the h value increases (as we go further away from the given point $x = 2$). This is because as the difference increases, it becomes more sensitive to the changes of the function between the two points, especially when it is not a linear slope. In our case, the function is not linear slope, which is why calculating the derivative through slope is more accurate.

Logistic regression

Logistic regression is a classification tool that models the probability of an event taking place by having the log odds for the event be a linear combination of one or more independent variables. Specifically, let $\vec{x} = \langle x_1, \dots, x_m \rangle$ be an m dimensional vector of independent variables (features), and y be the corresponding binary dependent variable (label). The probability of having $y = 1$ is modeled as

$$P_y = \frac{1}{1 + e^{-(b_0 + b_1 \cdot x_1 + \dots + b_m \cdot x_m)}} = \frac{1}{1 + e^{-(b_0 + \vec{b}_{1:m} \cdot \vec{x})}}$$

Given a set of data points $\langle \vec{x}_k, y_k \rangle$ with $k \in [1, n]$, how can we fit the model with these data, i.e., how to choose the best $\vec{b} = b_0, b_1, \dots, b_m$?

One way is to write out the likelihood

$$\prod_{k:y_k=1} P_{y_k} \prod_{k:y_k=0} (1 - P_{y_k})$$

and find b_0, b_1, \dots, b_m that maximize its logarithm,

$$l = \sum_{k:y_k=1} \ln(P_{y_k}) + \sum_{k:y_k=0} \ln(1 - P_{y_k})$$

In contrast to computing the closed form gradient of a Least-squares loss in a linear model (chapter 5 of MML book), doing the same for logistic regression is not possible. Gradient descent can be used to optimize such function l , and we will implement it step-by-step. First, write a function `log_likelihood` in the next cell that computes the log-likelihood given data points and \vec{b} . [5 pts]

In [384... `import numpy as np`
`import sklearn`

```
In [461]: def log_likelihood(X,y,b):
    """
    X: n*m numpy data array.
    y: one dimension numpy data array of length n
    b: one dimension numpy data array of length m+1

    Return the log likelihood.
    """

    n, m = X.shape
    ones = np.ones((n, 1))
    x = np.copy(X)
    x = np.concatenate([ones, x], axis=1)
    z = x.dot(b)
    z *= -1
    p = 1 / (1 + np.exp(z))
    logistic = 0
    for val in range(len(p)):
        if (p[val] != 1 and p[val] != 0 and not np.isnan(p[val])):
            logistic += y[val] * np.log(p[val]) + (1-y[val]) * np.log(1-p[val])
    return logistic
```

Test your log_likelihood function with a small example below.

```
In [462]: X=np.array([[0.1],[0.5],[1.]])
y=np.array([0,0,1])
b=np.array([0.,1.])
# Your answer should be around -2.03
log_likelihood(X,y,b)
```

Out[462]: -2.031735331771901

In []:

Now that we have a function to maximize, the next step is to compute the gradient of the log-likelihood with respect to parameter \vec{b} . Use the method with Newton's difference quotient, and set $h = 0.0001$. Implement the function compute_gradient in the next cell. [7 pts]

```
In [463... def compute_gradient(X,y,b):
# The inputs are the same as the ones of log_likelihood
h = 0.0001

grad = np.zeros(b.shape)
for i in range(b.shape[0]):
    b_plus = np.copy(b)
    b_plus[i] += h
    b_minus = np.copy(b)
    grad[i] = (log_likelihood(X, y, b_plus) - log_likelihood(X, y, b_minus))

return grad
```

```
In [464... # Test your function here, preserve the output
compute_gradient(X,y,b)
```

```
Out[464]: array([-0.87853115, -0.09479906])
```

Once we know how to compute the gradients, we can optimize the objective, which is log-likelihood in our case, using gradient descent. It iteratively changes the parameters in a small "step" towards the gradient direction, i.e., the direction where the objective increases at the fastest pace. Formally, denote the calculated gradients as $\Delta(\vec{b})$, we can update our parameters via $\vec{b} = \vec{b} + \gamma \cdot \Delta(\vec{b})$, where γ is the size of the "step". Repeat this process until the objective stop improving or a pre-set max number of iterations is reached. **Note in practice, the value of gradient changes over iterations and can be very large/small, so you should normalize the gradient vector every iteration, i.e., scale it to $\frac{\Delta(\vec{b})}{\|\Delta(\vec{b})\|_2}$, before using it to compute the new \vec{b} . Therefore, the update rule for parameters becomes $\vec{b} = \vec{b} + \gamma \cdot \frac{\Delta(\vec{b})}{\|\Delta(\vec{b})\|_2}$.**

Implement the gradient_descent function below. [7 pts]

```
In [465... def gradient_descent(X, y, initial_b, step_size, max_iteration):
    """
    X: n*m numpy data array.
    y: one dimension numpy data array of length n
    initial_b: one dimension numpy data array of length m+1
    step_size: scalar, the size of one step update
    max_iteration: scalar, the max number of iterations
    Return the updated coefficient vector b.
    """

    copyb = np.copy(initial_b)
    for count in range(max_iteration):
        grad = compute_gradient(X, y, copyb)
        scaleVal = 0
        for x in range(len(grad)):
            scaleVal += grad[x]**2
        copyb += grad/scaleVal * step_size

    return copyb
```

Test the function with the previous example again. Print for each sample from X, based on your model, the probability of having label=1.

```
In [478... optimized_b = gradient_descent(X, y, b, 0.1, 1000)

# compute and print the probability for each row in X below using optimized_

# Print the probability of having label=1 for each sample from X
z = []
for x in range(len(X)):
    value = 1 / (1+ np.exp(-(optimized_b[0]+optimized_b[1]*X[x])))
    z.append(value)
    print("Probability of having label = 1 for X = " + str(X[x]) + " is " +

print(z)

Probability of having label = 1 for X = [0.1] is [1.55085739e-97]
Probability of having label = 1 for X = [0.5] is [1.51222147e-102]
Probability of having label = 1 for X = [1.] is [8.23985975e-109]
[array([1.55085739e-97]), array([1.51222147e-102]), array([8.23985975e-109])
]
```

Next, we apply the implemented logistic regression model to a real dataset. The dataset is a trimmed breast-cancer-Wisconsin dataset from UCI machine learning Repository. Only 100 data points are offered in the training set to make sure the computation can be finished swiftly, no matter how you implement the optimizer. The training dataset is loaded in the next cell, and the vector \vec{b} is also randomly initialized.

Fit three models with the training set using different step size ranging in [0.01,0.05,0.1] and set the max number of iterations as 10000. How do the final log-likelihood value and the number of iterations change with different step sizes? [7 pts]

```
In [479... f = open("breast-cancer-wisconsin.data", "r")
X_train = []
y_train = []
for line in f:
    tmp = []
    for part in line.strip().split(",")[1:-1]:
        tmp.append(float(part))
    y_train.append((0 if line.strip().split(",")[-1]=="2" else 1))
    X_train.append(tmp)
X_train = np.array(X_train)
y_train = np.array(y_train)
random_b = np.random.uniform(0,1,size=(10))
```

```
In [481... # Fit three models with different step size, report the final log-likelihood
# number of iterations and the final coefficient vector b.

step_s = [0.01, 0.05, 0.1]
maxI = 10000
for val in step_s:
    final_val = gradient_descent(X_train, y_train, random_b, val, maxI)
    log_final = log_likelihood(X_train, y_train, final_val)
    print("Final log-likelihood for step size of " + str(val) + " is ")
    print(log_final)
    print()
    print("Final coefficient vector b for step size of " + str(val) + " is ")
    print(str(final_val))
    print()
```

Final log-likelihood for step size of 0.01 is
-341.38316583268397

Final coefficient vector b for step size of 0.01 is
[0.93710591 -0.02141338 0.28514312 0.87181098 0.2010544 0.87709955
0.53525439 0.0506262 0.43571095 0.15960707]

Final log-likelihood for step size of 0.05 is
-8.14105151319611

Final coefficient vector b for step size of 0.05 is
[-11.90270091 0.74883558 -0.87601827 0.63135412 0.87510495
0.29514687 0.09852812 0.6255476 0.76950141 1.00979473]

Final log-likelihood for step size of 0.1 is
-9.76248204728458

Final coefficient vector b for step size of 0.1 is
[-20.46523756 1.47893533 -2.48575658 1.53034596 1.24269018
0.75599993 -0.13183902 1.49665361 1.64403522 2.2593107]

Finally, load the test dataset, and predict for each sample in the test set what labels it should have using the model obtained. Compare your results with the ground truth labels, and report the accuracy rate. [4 pts]

```
In [482... f = open("test_data.txt", "r")
X_test = []
y_test = []
for line in f:
    tmp = []
    for part in line.strip().split(",")[1:-1]:
        tmp.append(float(part))
    y_test.append((0 if line.strip().split(",")[-1]=="2" else 1))
    X_test.append(tmp)
```


In [487... *# Predict based on your models and report the accuracy*

```

z = 0
answer = []
for val in range(len(X_test)):
    if val == 0:
        answer.append(1/(1+np.exp(-1*final_val[0])))
    else:
        lst = [i[val-1] for i in X_test]
        for num in lst:
            z += final_val[val] * num
        answer.append(1/(1+np.exp(-1*z)))
z = 0

print("Ground truth labels ")
print(y_test)
print()

print("My results ")
print(answer)
print()

rounded_answer = np.round_(answer, decimals = 6)

print("My answer rounded ")
print(rounded_answer)
print()

print("The accuracy rate is 80%")

```

Ground truth labels

[0, 1, 1, 1, 1, 1, 0, 1, 0, 1]

My results

[1.2943754169239662e-09, 1.0, 1.8244534296893125e-51, 1.0, 1.0, 0.9999999999999999, 0.0015621766020403106, 1.0, 1.0, 1.0]

My answer rounded

[0. 1. 0. 1. 1. 1. 0.001562 1.
1. 1.]

The accuracy rate is 80%

In []: