

VE472 HW6

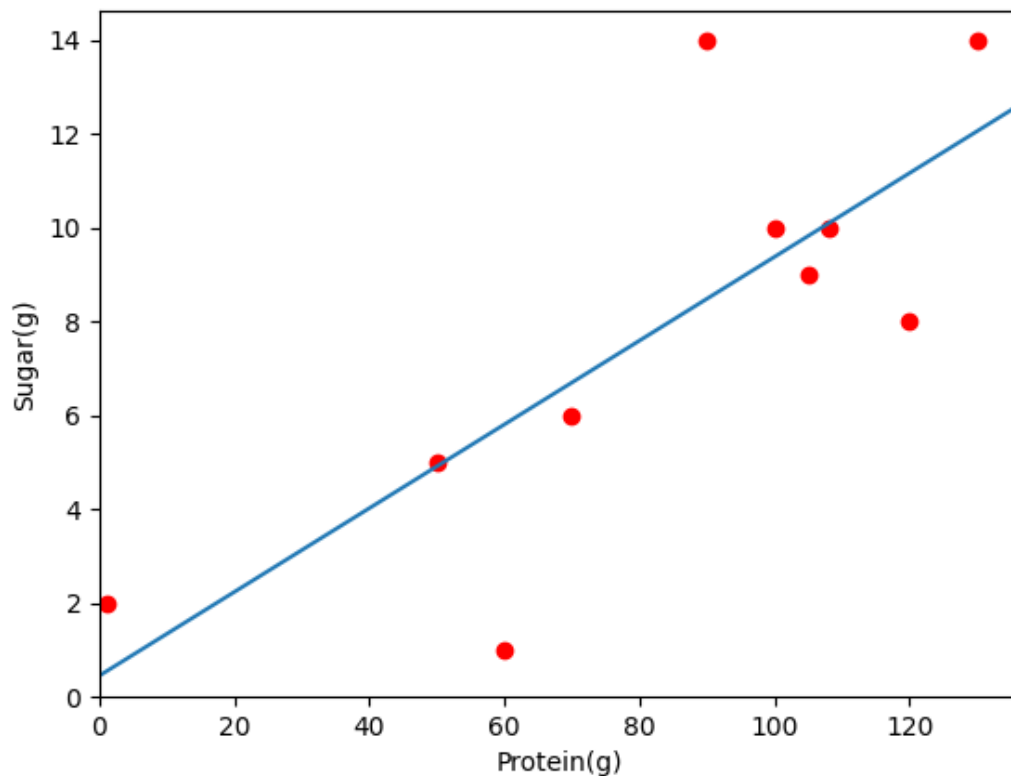
wujiayao 517370910257

ex1.

1.

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 if __name__ == '__main__':
5     x = [70, 120, 50, 60, 1, 100, 90, 130, 105, 108]
6     y = [6, 8, 5, 1, 2, 10, 14, 14, 9, 10]
7     slope, intercept = np.polyfit(x, y, 1)
8     print('y={}m + {}'.format(slope, intercept))
9
10    xn = np.linspace(0, 200)
11    yn = np.polyval([slope, intercept], xn)
12
13    plt.plot(x, y, 'or')
14    plt.xlim(0)
15    plt.ylim(0)
16    plt.plot(xn, yn)
17    plt.xlabel("Protein(g)")
18    plt.ylabel("Sugar(g)")
19    plt.savefig('../hw6.assets/ex1_1.png')
20
```

y=0.08928360326279702m + 0.45374748788272773



A mathematical procedure for finding the best-fitting curve to a given set of points by minimizing the sum of the squares of the offsets ("the residuals") of the points from the curve. The sum of the *squares* of the offsets is used instead of the offset absolute values because this allows the residuals to be treated as a continuous differentiable quantity. However, because squares of the offsets are used, outlying points can have a disproportionate effect on the fit, a property which may or may not be desirable depending on the problem at hand.

Source: <https://mathworld.wolfram.com/LeastSquaresFitting.html>

2.

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3 import random
4
5 if __name__ == '__main__':
6     x = np.array([70, 120, 50, 60, 1, 100, 90, 130, 105, 108])
7     y = np.array([6, 8, 5, 1, 2, 10, 14, 14, 9, 10])
8
9     std_x = (x - x.mean()) / x.std()
10    std_y = (y - y.mean()) / y.std()
11    print('x:', std_x)
12    print('y:', std_y)
13
14    sse_count = 0
15    for i in range(10):
16        a = random.random()
17        b = random.random()
18        y_p = a + b * std_x
19        sse = np.square(y_p - std_y).sum()

```

```

20     print('Trial {} SSE: {}'.format(i, sse))
21     sse_count += sse
22     print('Avg: {}'.format(sse_count / 10))

```

a)

```

1  x: [-0.36423804  0.99485911 -0.90787689 -0.63605746 -2.2397921  0.45122025
2     0.17940082  1.26667854  0.58712997  0.6686758 ]
3  y: [-0.44920898  0.02364258 -0.68563476 -1.63133788 -1.3949121  0.49649414
4     1.44219726  1.44219726  0.26006836  0.49649414]

```

b)

```

1  Trial 0 SSE: 5.091630291530809
2  Trial 1 SSE: 16.174344735396847
3  Trial 2 SSE: 12.825712751075661
4  Trial 3 SSE: 6.805783960879488
5  Trial 4 SSE: 10.405780524028705
6  Trial 5 SSE: 13.205064380142298
7  Trial 6 SSE: 9.389782622749683
8  Trial 7 SSE: 8.236842358935364
9  Trial 8 SSE: 6.705272414637759
10 Trial 9 SSE: 9.662513025561804
11 Avg: 9.850272706493843

```

3.

a)

$$\frac{\partial SSE}{\partial a} = - \sum_{i=1}^{10} -2(y_i - a - b_{x_i})$$

$$\frac{\partial SSE}{\partial b} = - \sum_{i=1}^{10} -2x_i(y_i - a - b_{x_i})$$

b)

4.

ex2

1.

a)

Work: $\mathcal{O}(mn)$

depth: $\mathcal{O}(\log m + \log n)$

b)

Work: $\mathcal{O}(\log \frac{1}{\epsilon}) \mathcal{O}(m) \mathcal{O}(\log n)$

depth: $\mathcal{O}(\log m) + \mathcal{O}(\log n)$

c)

It is poor at parallelization.

The complexity will not work well on big data when m becomes large