

DSC 40A

Theoretical Foundations of Data Science I

In This Video

Can we use linear regression to fit nonlinear functions to data?

Recommended Reading

Course Notes: Chapter 2, Section 1

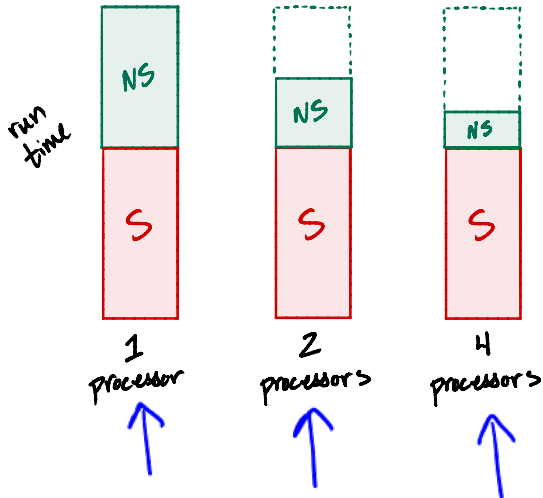
Example: Parallel Processing



Problem

- ▶ Some parts of a program are necessarily **sequential**.
- ▶ E.g., downloading the data must happen before analysis.
- ▶ More processors do not speed up **sequential** code.
- ▶ But they do speed up **non-sequential** code.

Speedup



Amdahl's Law

The time T it takes to run a program on p processors is:

$$T(p) = t_S + \frac{t_{NS}}{p}$$

where t_S and t_{NS} are the time it takes the sequential and non-sequential parts to run on one processor, respectively.

Amdahl's Law

The time T it takes to run a program on p processors is:

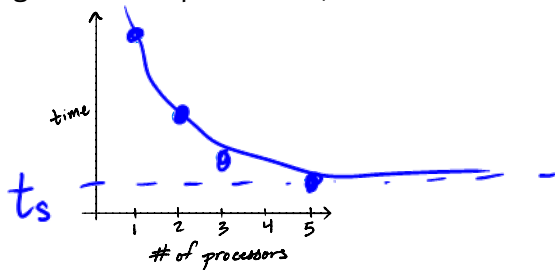
$$T(p) = t_S + \frac{t_{NS}}{p}$$

where t_S and t_{NS} are the time it takes the sequential and non-sequential parts to run on one processor, respectively.

Problem: we don't know t_S and t_{NS} .

Fitting Amdahl's Law

- **Solution:** we will learn t_S and t_{NS} from data.
- Run with varying number of processors, record total time:



- Find prediction rule $H(p) = \frac{t_{NS}}{p} + t_S$ by minimizing MSE.
$$= t_{NS} \left(\frac{1}{p} \right) + t_S$$

General Problem

- ▶ Given data $(x_1, y_1), \dots, (x_n, y_n)$.
- ▶ Fit a **non-linear** rule $H(x) = w_1 \cdot \frac{1}{x} + w_0$ by minimizing MSE:

$$R_{sq}(H) = \frac{1}{n} \sum_{i=1}^n (H(x_i) - y_i)^2$$

Using definition of H :

$$R_{sq}(H) = \frac{1}{n} \sum_{i=1}^n (\underline{w_1} \cdot \frac{1}{x_i} + \underline{w_0} - y_i)^2$$

Minimizing MSE

- Take partial derivatives, set to zero, solve. You'll find:

$$w_1 = \frac{\sum_{i=1}^n \left(\frac{1}{x_i} - \frac{1}{n} \sum_{i=1}^n \frac{1}{x_i} \right) (y_i - \bar{y})}{\sum_{i=1}^n \left(\frac{1}{x_i} - \frac{1}{n} \sum_{i=1}^n \frac{1}{x_i} \right)^2}$$

$$w_0 = \bar{y} - w_1 \cdot \frac{1}{n} \sum_{i=1}^n \frac{1}{x_i}$$

Minimizing MSE

- Take partial derivatives, set to zero, solve. You'll find:

$$w_1 = \frac{\sum_{i=1}^n \left(\frac{1}{x_i} - \frac{1}{n} \sum_{i=1}^n \frac{1}{x_i} \right) (y_i - \bar{y})}{\sum_{i=1}^n \left(\frac{1}{x_i} - \frac{1}{n} \sum_{i=1}^n \frac{1}{x_i} \right)^2}$$

$w_0 = \bar{y} - w_1 \cdot \frac{1}{n} \sum_{i=1}^n \frac{1}{x_i}$

- Define

$$z_i = \frac{1}{x_i}$$

$$\bar{z} = \frac{1}{n} \sum_{i=1}^n z_i = \frac{1}{n} \sum_{i=1}^n \frac{1}{x_i}$$

$$w_1 = \frac{\sum_{i=1}^n (z_i - \bar{z})(y_i - \bar{y})}{\sum_{i=1}^n (z_i - \bar{z})^2} \quad w_0 = \bar{y} - w_1 \bar{z}$$

Fitting Non-Linear Trends z

To fit a prediction rule of the form $H(x) = w_1 \left(\frac{1}{x} \right) + w_0$:

1. Create a new data set $(z_1, y_1), \dots, (z_n, y_n)$, where $z_i = \frac{1}{x_i}$.
2. Fit $H(z) = w_1 z + w_0$ using familiar least squares solutions:

$$w_1 = \frac{\sum_{i=1}^n (z_i - \bar{z})(y_i - \bar{y})}{\sum_{i=1}^n (z_i - \bar{z})^2}$$

$$w_0 = \bar{y} - w_1 \cdot \bar{z}$$

3. Use w_1 and w_0 in original prediction rule, $H(x)$.

Example: Amdahl's Law

- We have timed our program:

Processors	Time (Hours)
1	8
2	4
4	3

- Fit prediction rule: $H(p) = \frac{t_{NS}}{p} + t_s$

Example: fitting

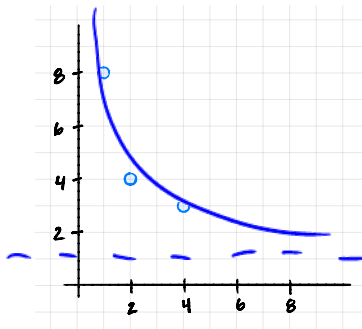
$$H(x) = w_1 \cdot \frac{1}{x_i} + w_0$$

$$\bar{z} = 7/4 \cdot 1/3 = 7/12$$

$$\bar{y} = 5$$

$$w_1 = \frac{\sum_{i=1}^n (z_i - \bar{z})(y_i - \bar{y})}{\sum_{i=1}^n (z_i - \bar{z})^2} = \frac{2}{\left(\frac{42}{144}\right)} = \frac{2 \cdot 144}{42} = \frac{48}{7}$$

$$w_0 = \bar{y} - w_1 \bar{z} = 5 - \frac{48}{7} \cdot \frac{7}{12} = 5 - 4 = 1$$



x_i	y_i	z_i	$(z_i - \bar{z})$	$(y_i - \bar{y})$	$(z_i - \bar{z})(y_i - \bar{y})$	$(z_i - \bar{z})^2$
1	8	1	$5/12$	3	$15/12$	$25/144$
2	4	$1/2$	$-1/12$	-1	$1/12$	$1/144$
4	3	$1/4$	$-4/12$	-2	$8/12$	$16/144$
\uparrow sum $24/12 = 2$					\uparrow sum $42/144$	

Example: Amdahl's Law

- ▶ We found: $t_{NS} = \frac{48}{7} \approx 6.88$, $t_S = 1$
- ▶ Therefore our prediction rule is:

$$\begin{aligned} H(p) &= \frac{t_{NS}}{p} + t_S \\ &= \frac{6.88}{p} + 1 \end{aligned}$$

Linear in the Parameters

- We can fit rules like:

$$w_1 x + w_0$$

$$w_1 \cdot \frac{1}{x} + w_0$$

$$w_1 x^2 + w_0$$

$$w_1 2^x + w_0$$

- We can't fit rules like:

$$e^{w_1 x} + w_0$$

$$\sin(w_1 x + w_0)$$

- Has to be **linear in the parameters**, or linear as a function of w_1, w_0 .

linear
combs
of
 w_i 's

Transformations

- Try rewriting functions to see if they can be expressed as linear functions in new variables.

- **Example**

$$H(x) = c_0 x^{c_1}$$

$$y = c_0 x^{c_1}$$

$$\log_{10} y = \log_{10} c_0 + x^{c_1}$$

$$\log y = \log c_0 + \log x^{c_1}$$

$$\underbrace{\log y}_v = \underbrace{\log c_0}_{w_0} + \underbrace{c_1}_{w_1} \underbrace{\log x}_z$$

Transformations

$$\boxed{\begin{aligned} y &= c_0 x^{c_1} \\ \log y &= \log c_0 + c_1 \log x \end{aligned}}$$

$\checkmark \quad \checkmark \quad \checkmark \quad \checkmark$
 $\checkmark \quad w_0 \quad w_1 \quad z$

$$w_1 = \frac{\sum_{i=1}^n (\log x_i - \frac{1}{n} \sum_{i=1}^n \log x_i) (\log y_i - \frac{1}{n} \sum_{i=1}^n \log y_i)}{\sum_{i=1}^n (\log x_i - \frac{1}{n} \sum_{i=1}^n \log x_i)^2}$$

$$w_0 = \frac{1}{n} \sum_{i=1}^n \log y_i - w_1 \cdot \frac{1}{n} \sum_{i=1}^n \log x_i$$

These w_0, w_1 satisfy $\log y = w_0 + w_1 \log x$
 $\underline{w_1} = \underline{c_1} \quad \log_{10} c_0 = w_0 \Rightarrow \underline{c_0} = 10^{w_0}$

General Strategy

To fit a prediction rule of the form $\underbrace{g(y)}_v = \underbrace{w_1}_{z} \cdot \underbrace{f(x)}_z + \underbrace{w_0}_z$:

1. Create a new data set $(z_1, v_1), \dots, (z_n, v_n)$, where $z_i = f(x_i)$ and $v_i = g(y_i)$.
2. Fit $v = w_1 z + w_0$ using familiar least squares solutions:

$$w_1 = \frac{\sum_{i=1}^n (z_i - \bar{z})(v_i - \bar{v})}{\sum_{i=1}^n (z_i - \bar{z})^2} \qquad w_0 = \bar{v} - w_1 \cdot \bar{z}$$

where \bar{z} is the mean of the z_i 's, \bar{v} is the mean of the v_i 's.

3. If necessary, use w_0 and w_1 to find the parameters of the original prediction rule.

Summary

- ▶ We can sometimes fit nonlinear functions to data by thinking of these non-linear functions as linear functions in new variables.
- ▶ **Next Time:** Using linear algebra to do regression helps us fit even more non-linear functions to data and allows us to make predictions based on multiple features.
- ▶ E.g., experience, highest education level, GPA, number of internships, etc.