

## 02-512 Assignment 02

Karan Sikka

ksikka@cmu.edu

October 1, 2014

---

1

---

(a)

For one combination of substitutions made, define the following variables:

Let  $\vec{x}$  be an  $n$ -dimensional vector of 1s and 0s where  $x_i$  is 1 if the  $i^{th}$  substitution was made, and 0 if not.

Let  $\vec{k}$  be an  $n$ -dimensional vector of weights, where  $k_i$  is the amount that making the  $i^{th}$  substitution contributes to the expression level.

Let  $c$  be the baseline expression level.

Let  $y$  be the expression level.

Then:

$$\vec{k} \cdot \vec{x} + c = y$$

If we have  $n+1$  such equations, one for each combination, we can represent the linear system in the following matrix: ( $x_i$  will now be  $x_{j,i}$  where  $j$  is the ordinal number of the combination/equation, and apply a similar transformation to the indexes of  $k$  and  $y$ )

$$\begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,n} & 1 \\ x_{2,1} & x_{2,2} & \cdots & x_{2,n} & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{n+1,1} & x_{n+1,2} & \cdots & x_{n+1,n} & 1 \end{bmatrix} \begin{bmatrix} k_1 \\ k_2 \\ \vdots \\ k_n \\ c \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_{n+1} \end{bmatrix}$$

We know the eqns are independent since the combinations of substitutions are distinct, and may not be scalar multiples of one another because all values are zero or one.

The problem is in the canonical  $A\vec{x} = \vec{b}$  form (but we'll call it  $X\vec{k} = \vec{y}$ ). The  $A$  matrix is full-rank (number of rows equals number of columns), and you'll find exact solutions for  $\vec{k}$  using Gaussian Elimination.

(b) Now we have  $n^2$  equations rather than  $n+1$  and the number of rows is greater than the number of columns.

Therefore the system is overdetermined. Given our overdetermined system  $X\vec{k} = \vec{y}$ , we find  $\vec{k}$  which minimizes the sum of least-squares by solving the linear system  $(X^T X)\vec{k} = (X^T \vec{y})$

(c) We add an entry for each pair of substitutions ( $\binom{n}{2}$  pairs) to our bit vector  $x$  and correspondingly our weight vector  $k$  for each experimental test. The bit vector  $x$  had  $n+1$  entries before, now we add  $\binom{n}{2} = \frac{1}{2}(n^2 - n)$  entries to result in a total of  $\frac{1}{2}(n^2 + n + 2)$  entries.

The number of entries is larger than  $n^2$  when  $n \geq 2$  (determined using Wolfram Mathematica), so the system is underdetermined. There are many valid solutions, but one way to get one solution

would be to find the pseudoinverse of the matrix formed by the linear systems (done in a manner similar to part A) using Singular Value Decomposition. Then multiply the pseudoinverse by the  $\vec{y}$  to obtain a solution for  $\vec{k}$  that has the most number of zeros in it.

## 2

(a)

$$\begin{aligned}
 F(x) &= \frac{d}{dx} E(x) \\
 &= \frac{d}{dx} E_1(x) + \frac{d}{dx} E_2(x) \\
 &= \frac{d}{dx} \frac{1}{4} x^4 + \frac{d}{dx} x^2 \\
 &= x^3 + 2x
 \end{aligned}$$

(b)

$$\begin{aligned}
 F(x) &= x^3 + 2x = 5 \\
 \implies x^3 + 2x - 5 &= 0
 \end{aligned}$$

Now we need to find the values of  $x$  for which the equation  $x^3 + 2x - 5 = 0$ . This is a zero-finding problem.

(c) Bisection method for  $f(x) = x^3 + 2x - 5$

Initialization:

$$x_{min} = 1 \text{ and } x_{max} = 2.$$

$$f_{min} = x_{min}^3 + 2x_{min} - 5 = 1 + 2 - 5 = -2$$

$$f_{max} = x_{max}^3 + 2x_{max} - 5 = 8 + 4 - 5 = 7$$

First round:

$$x_{mid} = \frac{1+2}{2} = 1.5$$

$$f_{mid} = f(x_{mid}) = 1.375$$

It's greater than zero, so:

$$x_{max} = 1.5$$

Second round:

$$x_{mid} = \frac{1+1.5}{2} = 1.25$$

$$f_{mid} = f(x_{mid}) = -0.546875$$

It's less than zero, so:

$$x_{min} = 1.25$$

Third round:

$$x_{mid} = \frac{1.25+1.5}{2} = 1.375$$

$$f_{mid} = f(x_{mid}) = 0.34961...$$

It's greater than zero, so:

$$x_{max} = 1.375$$

Return  $x = \frac{x_{min} + x_{max}}{2} = \frac{1.25 + 1.375}{2} = 1.3125$

The backward error is  $x_{max} - x_{min} = 1.375 - 1.25 = 1.125$

The forward error is  $f(x) = 0.34961...$

(d) Secant method for  $f(x) = x^3 + 2x - 5$

Initialization:

$x_{min} = 1$  and  $x_{max} = 2$ .

$f_{min} = x_{min}^3 + 2x_{min} - 5 = 1 + 2 - 5 = -2$

$f_{max} = x_{max}^3 + 2x_{max} - 5 = 8 + 4 - 5 = 7$

First round:

$x_{mid} = x_{min} - \frac{x_{max} - x_{min}}{f_{max} - f_{min}} f_{min} = 1 - \frac{2 - 1}{7 - (-2)} - 2 = 1$

$f_{mid} = f(x_{mid}) = 1.375$

It's greater than zero, so:

$x_{max} = 1.5$

Second round:

$x_{mid} = \frac{1 + 1.5}{2} = 1.25$

$f_{mid} = f(x_{mid}) = -0.546875$

It's less than zero, so:

$x_{min} = 1.25$

Third round:

$x_{mid} = \frac{1.25 + 1.5}{2} = 1.375$

$f_{mid} = f(x_{mid}) = 0.34961...$

It's greater than zero, so:

$x_{max} = 1.375$

Return  $x = \frac{x_{min} + x_{max}}{2} = \frac{1.25 + 1.375}{2} = 1.3125$

The backward error is  $x_{max} - x_{min} = 1.375 - 1.25 = 1.125$

The forward error is  $f(x) = 0.34961...$