

ECE 331 Project 2

Machine Learning (Linear Regression), Matrix Multiplication, and
Cache Simulation in MIPS assembly language
(Due Friday, November 20th 11:59 PM)

Objectives:

- Learn weights for a linear regression function that predicts the performance of a processor based on empirical data
- Use learned weights in a MIPS assembly program to perform matrix multiplication on a matrix of test data
- Perform cache simulations to understand cache performance across different cache architectures

A Crash Course on Machine Learning and Linear Regression

In this project, you will be implementing a rudimentary supervised machine learning algorithm, called linear regression. In previous algebra courses, you have already encountered a simple form of linear regression where you can fit a linear function to a set of values using the formula:

$$y(x) = mx + b \quad (1)$$

Where m is the slope of a line and b is the y-intercept. A more general form of Linear Regression is of the form:

$$y(x) = w_1x_1 + w_2x_2 + \cdots + w_nx_n + d \quad (2)$$

Where there are n linear weights that get multiplied by n input terms and are summed together along with a bias value d .

When linear regression is used in the context of machine learning, we can think of the terms $x_1 \dots x_n$ as a set of input parameters for a particular example and weights

$w_1 \dots w_n$ and bias term d as a learned set of weights that are used to predict the output y .

After determining the weights for the linear function defined above, to determine the predicted output y for a sequence of input examples x , we can formulate this computation as a matrix multiplication:

$$Y = W \cdot X + d \quad (3)$$

Where W is a $1 \times n$ matrix and X is an $n \times m$ matrix, Y is a $1 \times m$ matrix of output predictions, and d is a simple scalar value. The dot denotes the dot product operation, which is the operation performed in equation 2 above.

There a number of algorithms that can be used to estimate the best values for these weights such as ordinary least squares or gradient descent, but a study of these techniques is outside of the scope of this course. The goal of these techniques is to minimize the error of predictions after fixing the weight parameters. In this project, we will be using the tool Weka to learn these weights for a dataset that maps CPU parameters to ranked processor performance.

To learn the weights needed for our matrix multiplication, the first step is to download the open source data mining tool Weka:

<https://www.cs.waikato.ac.nz/ml/weka/downloading.html>

Weka is a data mining tool that can be used to learn machine learning model representations of data. It is a powerful tool, but we will be using it to perform a simple linear regression. In order to use Weka for machine learning tasks, it must be provided a set of training examples that use the .arff file extension. This is similar to a comma separated value format file (.csv), but has extra header information that describes the format of the input data and whether to consider each parameter a discrete set of classes, or continuous values that are integer or floating point values.

For linear regression, we need a .arff file that uses a set of continuous values that give us input values x and the actual output value y . The data for this project is found on Piazza under 'General Resources' and is the file `cpu.arff`.

If opened in a text editor, we can view the header information at the top of the file, which describes the dataset. We can see that there are 6 input parameters that quantify the processor cycle time and different memory parameters (we haven't

covered channel width in class, but this is the width of the databus that can be used to talk to a cache or main memory in a computer). The last parameter is the relative performance of the CPU with these parameters. Looking at the rest of the file, we can see there are 209 training examples for different CPU configurations that provide is a complete set of parameters.

To load this data into Weka for analysis, select 'Explorer'. In the following screen select 'open file' and locate the `cpu.arff` file. After loading the file, we can see that there are 7 parameters and in the lower right hand corner, we can see statistical information about the input and output parameters.

To learn the weights for our linear regression, click the 'classify' tab. Press the 'classifier' button and select functions -> linear regression. The right of this box, we can see parameters related to how the weights are computed. Make sure decimal places is set to 4 or greater. To learn the weights, Weka partitions the data into training and test sets – for our training, we will use the default approach of creating 10 'folds' of data, where random permutations of the data are split into 10 pieces. 9 of these pieces are used for training, while 1 is used for test. These values are rotated, and the weights are learned that minimize error in the output, in this case 'relative CPU performance'.

You can find the weight values in the classifier output window as the values after 'class =' that multiply each parameter. These are the weights we will use for our matrix multiplication algorithm. Since we want to avoid the use of floating point values, multiply each of these values by 1000 and round them (if interested in learning more about fixed point math, more details can be found here: https://en.wikipedia.org/wiki/Fixed-point_arithmetic). After a result is computed using matrix multiplication with the weights, we can divide the result by 1000 to determine the final result.

Deliverables:

-Screenshot of Weka output that shows the weights computed when training using a 10-fold cross validation. Include the performance results, which are included after the linear function.

Matrix Multiplication in MARS

For the purposes of this cache simulation assignment, we will be working with the MIPS assembly language using the MARS simulator which can be downloaded here:

<http://courses.missouristate.edu/kenvollmar/mars/>

This is a self-contained Java JAR file and can be run using the command `'java -jar Mars45.jar'`. Make sure you have a recent version of the Java runtime installed – *there is no need to use the Virtual Box virtual machine for this project and it can be run using your native operating system.*

The fixed point weights you computed in the previous section can be input directly into the MIPS assembly at the line `'learned_weight_vector'`. The 6 x 100 training examples from the .arff file are hard coded into the simulator's memory in the starter code in another memory location. These samples are identical to what you trained on in the .arff file – only the relative performance output is excluded. The key instructions that will affect cache performance are the `lw` and `sw` that load a 32-bit word from memory and store a 32-bit to memory, respectfully. The format `lw $s0, 32($s1)`, for example, loads a value from a memory address specified by the value held in register `$s1`, adds the offset 32 and places the result in register `$s0`. A store word instruction behaves the same way, but instead stores a register value in memory. The corresponding LEGv8 instructions that operate on 64-bit words are `ldur` and `stur`.

We provide you with the MIPS code that implements matrix multiplication – for your reference we give the equivalent C implementation of matrix multiplication below. It is your task to predict the relative performance using the weights provided from Weka and store the results in a fixed memory region. The place to put predictions is provided in the starter code. To compute these predictions, you need to perform matrix multiplication between the weight vector and the matrix of training examples.

C Code for Matrix Multiplication:

```
#include <stdio.h>

int main()
{
    int m, n, p, q, c, d, k, sum = 0;
    int first[10][10], second[10][10], multiply[10][10];

    printf("Enter number of rows and columns of first matrix\n");
    scanf("%d%d", &m, &n);
    printf("Enter elements of first matrix\n");

    for (c = 0; c < m; c++)
        for (d = 0; d < n; d++)
            scanf("%d", &first[c][d]);

    printf("Enter number of rows and columns of second matrix\n");
    scanf("%d%d", &p, &q);

    if (n != p)
        printf("The matrices can't be multiplied with each other.\n");
    else
    {
        printf("Enter elements of second matrix\n");

        for (c = 0; c < p; c++)
            for (d = 0; d < q; d++)
                scanf("%d", &second[c][d]);

        for (c = 0; c < m; c++) {
            for (d = 0; d < q; d++) {
                for (k = 0; k < p; k++) {
                    sum = sum + first[c][k]*second[k][d];
                }

                multiply[c][d] = sum;
                sum = 0;
            }
        }

        printf("Product of the matrices:\n");

        for (c = 0; c < m; c++) {
            for (d = 0; d < q; d++)
                printf("%d\t", multiply[c][d]);

            printf("\n");
        }
    }

    return 0;
}
```

Deliverables:

-Using the results from Weka, compute \hat{Y} (the prediction) for 3 training examples (i.e. multiply weights by 3 rows of data from Weka). Show all work.

-Provide a screenshot of the MARS memory map that shows the fixed point approximation for these 3 training examples. Comment on any error introduced by our fixed point approximation. Note: you will find the option 'Show Labels Window' under 'Settings' is extremely helpful in finding the memory offset where the predictions are located in data memory.

-Compare the predicted result to the ground truth relative performance value in the .arff file. How accurately did the linear regression function predict the value?

Optimizing Cache Performance Using a Cache Simulator

Matrix multiplication will be slow if we don't provide a cache organization that can ensure a) the weights used in matrix multiplication stay in the cache and b) the spatial locality of each set of input data is accounted for.

In MARS, it is simple to perform these cache simulations. Under 'tools' select 'Data Cache Simulator'. At the bottom left of this window, select 'connect to MIPS'. With the simulator connected, run your matrix multiplication code for a given cache configuration. In the cache performance section of this window, the tool will report the hit or miss rate.

For each cache simulation for a given organization, click the 'reset' button to reset the cache contents and the hit or miss statistics.

Deliverables:

-Try 12 different cache configurations that use different sizes, block lengths, and levels of associativity. Report the configuration results you tried in a table.

-What was the hit rate for these different configurations? Ensure that at least one of your configurations achieved $\geq 90\%$ hit rate.

-Provide intuition into why each configuration achieves its performance.

-Provide a screenshot of the cache configuration that achieved $\geq 90\%$ hit rate.