| 175 | 1 | 0.279614 | -0.052290 | 0 | 1 | 0 | 0 |
|-----|---|----------|-----------|---|---|---|---|
| 63  | 1 | 0.799020 | 1.208986  | 0 | 0 | 1 | 0 |
| 67  | 0 | 0.279614 | -0.236227 | 1 | 0 | 0 | 0 |
| 216 | 0 | -2.144282 | -1.287291 | 1 | 0 | 0 | 0 |
| 145 | 0 | -1.798011 | 0.105369  | 0 | 0 | 1 | 0 |
| 286 | 1 | 1.837832 | -0.446439 | 1 | 0 | 0 | 0 |
| 339 | 1 | 0.625884 | 0.210476  | 0 | 0 | 1 | 0 |

**Ten rows of the data after transformations.**

Now that the data is ready, we see that there are six input features: gre , gpa , and the four rank dummy variables.

## Mean Square Error

We're going to make a small change to how we calculate the error here. Instead of the SSE, we're going to use the **mean** of the square errors (MSE). Now that we're using a lot of data, summing up all the weight steps can lead to really large updates that make the gradient descent diverge. To compensate for this, you'd need to use a quite small learning rate. Instead, we can just divide by the number of records in our data, $m$ to take the average. This way, no matter how much data we use, our learning rates will typically be in the range of 0.01 to 0.001. Then, we can use the MSE (shown below) to calculate the gradient and the result is the same as before, just averaged instead of summed.

$$E = \frac{1}{2m} \sum_{\mu} (y^{\mu} - \hat{y}^{\mu})^2$$

Here's the general algorithm for updating the weights with gradient descent:

- Set the weight step to zero: $\Delta w_i = 0$