



# Machine Learning and Data Science: Linear Regression Part 5

Wednesday July 4, 2017 by Dr Donald Kingham

Share: [f](#) [t](#) [g](#) [p](#) [e](#) [r](#)

## Linear Regression Part 5: Vectorization and Matrix Equations

We've covered a lot of fundamentals in the last 4 posts about Linear Regression and in this post we will cover another important idea, "vectorization". In this context we will rewrite the equations for Linear Regression in matrix/vector form and derive the optimized solution to find the model parameters by solving a simple matrix equation. We will also generalize the linear regression problem to multiple feature variables, **Multi-Varate Linear Regression**. This is a natural extension from expressing the equations in matrix form.

In the paragraph above I used the word "simple" in regard to the matrix equations. I should qualify that **Everything is simple when you know how to do it!** This post is a bit more advanced than the previous posts since it is using some linear algebra. **Linear algebra** is one of the important corner-stones of understanding machine learning, and actually, is a corner-stone of nearly all numerical computing.

The basic mathematics pre-requisites for understanding Machine Learning are Calculus, Linear Algebra, and Probability and Statistics. A computer science student that is interested in Machine Learning would be well advised to get a minor in Mathematics (or just get a degree in Mathematics instead). You don't need to know these subjects at the level of a mathematician (after more like an engineer). This is especially true for linear algebra. Linear algebra is a great first course for a theoretical/proving-course (not mathematics). What is needed for machine learning and a computer scientist is more like "applied numerical linear algebra".

You don't need to understand "all" of the math used in this post! The main take-aways are the final equations and the "idea" of rewriting formulas in matrix/vector form. You can think of vectors and matrices as "programming data structures" and then realize how they are used by libraries like numpy.

If you don't have much mathematics background don't be intimidated by the statements above! It's not that hard and you can learn a lot by just following along as best you can. I will explain everything in detail but will try to give insight into what's going on.

For background the previous 4 posts contain a lot of good stuff!

- [Linear Regression Part 1: Introduction](#)
- [Linear Regression Part 2: Getting and Evaluating Data](#)
- [Linear Regression Part 3: Model and Cost Function](#)
- [Linear Regression Part 4: Parameter Optimization by Gradient Descent](#)

These posts along with the current one were converted to Jupyter notebooks. The notebooks are available at <https://github.com/dkingham/blog-jupyter-notebooks>.

## Scalars, Vectors, Matrices and some Operations with them

Here is very brief (and inadequate!) introduction to some of the structures and operations in linear algebra.

### Scalar

A scalar is just a single number. For example,

$$a = 2$$

### Vector

A vector is a column array of scalars. It has 1 index ( $a_1$  is the first element of  $a$ , ... or maybe that could be  $a_0$ )

$$a = \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}$$

where  $a_1$  and  $a_2$  are scalars.

### Matrix

A matrix is a 2 dimensional array. The following matrix has 2 rows and 2 columns. You can think of the columns as vectors.

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

A matrix has 2 index values ( $i, j$ ). The " $i,j$ -element" of  $A$  is  $a_{ij}$ "

### Dimensions

We would say the matrix  $A$  is a "2 by 2" or  $2 \times 2$  matrix. In general it could be  $n \times m$ , i.e.  $n$  by  $m$ , where  $n$  and  $m$  are scalars. The vector  $a$  is "a 2 vector" i.e. has length 2. You could also think of it as a 2 by 1 matrix.

### Transpose

The transpose of a vector or matrix has its "indices switched". It is usual written as a superscript "T". For example,

$$a^T = [a_1 \quad a_2]$$

and

$$A^T = \begin{bmatrix} a_{11} & a_{21} \\ a_{12} & a_{22} \end{bmatrix}$$

The elements around the diagonal of  $A$  are switched, to general  $A^T \neq A$ .

### Multiplication

Scalars, vectors and matrices can be multiplied.

$$ca^T = [ca_1 \quad ca_2]$$

Matrices and vectors are multiplied so that each element of the result is a "row times column with the result summed". The following vector product is also called a vector "dot" product.

$$a^T b = [a_1 \quad a_2] \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} = a_1 b_1 + a_2 b_2$$

The "1 by 2" row vector  $a$  times the "2 by 1" column vector  $b$  is the "1 by 1" scalar  $a_1 b_1 + a_2 b_2$ .

$$AB = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} = \begin{bmatrix} a_{11}b_{11} + a_{12}b_{21} & a_{11}b_{12} + a_{12}b_{22} \\ a_{21}b_{11} + a_{22}b_{21} & a_{21}b_{12} + a_{22}b_{22} \end{bmatrix}$$

For matrices think "row times column summed" for each element of the equal matrix.

### Operation Dimensions

When you perform operations with vectors and matrices you need to keep their "dimensions" in mind. For conventional matrix multiplication the **inner dimensions must be conformable**. For example you can multiply a  $1 \times 2$  vector with a  $2 \times 3$  matrix and the result would be a  $1 \times 3$  dimensional vector. You can multiply a matrix or vector by a scalar element-wise and you can add a scalar to a matrix or vector element-wise. You can also add 2 vectors or matrices element-wise if they have the same dimensions.

You have to be careful when using matrices and vectors. It's always good to check your work by looking at the dimensions of all operations to be sure you haven't forgotten a transpose or something. A couple of other characteristics to be aware of is that matrix multiplication is not commutative, that is, in general  $AB \neq BA$  and the transpose or inverse of a product may not be what you think,  $(AB)^T = B^T A^T$  and  $(AB)^{-1} = B^{-1} A^{-1}$ .

That is not enough explanation! If you understand the above you are set, if not, then you might want to look at some of the [great Wikipedia](http://www.wikimedia.org) pages.

## Vectorizing the Linear Regression Model and Cost Function

### Model function in matrix/vector form

We have been thinking of our model function in two ways. One was as our "predictor equation",

$$\hat{h}(x) = a_0 + a_1 x$$

where, given the optimal values for  $a_0, a_1$ , we can use  $\hat{h}(x)$  as a predictor for the value of a house of size  $x$  sqft. The other way we were looking at  $\hat{h}(x)$  was in terms of our data set,

$$\hat{h}_i(x^{(i)}) = a_0 + a_1 x^{(i)}$$

where we were considering the  $i^{th}$  element in our training data set. Recall we use  $x$  for the number of elements in our dataset so the complete column of values  $x$  (and thus  $\hat{h}(x)$ ) are vectors with  $n$  elements.

We can consider  $\hat{h}(x)$  as a vector valued function of the complete column of  $n$  numbers that are house sizes in our dataset i.e.  $x$  is a vector of house sizes. We had also introduced  $a$  as being a 2-vector containing the parameters  $a_0, a_1$ . So, without using any special change in notation we can think of  $\hat{h}$  as a vector function of  $x$ ,

$$\hat{h}_i(x) = a_0 + a_1 x = \begin{bmatrix} 1 & x^{(i)} \\ 1 & x^{(i)} \\ \vdots & \vdots \\ 1 & x^{(n)} \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \end{bmatrix} = \begin{bmatrix} a_0 + a_1 x^{(1)} \\ a_0 + a_1 x^{(2)} \\ \vdots \\ a_0 + a_1 x^{(n)} \end{bmatrix} = \hat{X}a$$

**Matrix form of model across the full column of feature data,**

$$\hat{h}_i(X) = \hat{X}a$$

I have introduced the "**augmented matrix**"  $\hat{X}$  that has a column of 1's to take care of the constant term  $a_0$ . The model function  $\hat{h}$  on the full set of training data is now just a simple matrix vector multiplication!

### Cost function in matrix/vector form

We originally wrote our cost function as the sum of squared errors between the model function and the dataset values for the house prices  $y$ ,

$$J(a_0, a_1) = \frac{1}{2n} \sum_{i=1}^n (a_0 + a_1 x^{(i)} - y^{(i)})^2$$

A sum of squares is known as a "quadratic form" and we can write it in matrix form using the vector expression for  $\hat{h}_i(X)$  and the full column vector of house prices  $y$ . This will replace the summation,  $\sum$ , with matrix/vector multiplication.

**Matrix form of the cost function,**

$$J(a) = \frac{1}{2n} (Xa - y)^T (Xa - y) \\ = \frac{1}{2n} (y^T X^T Xa - 2y^T Xa + y^T y)$$

Note, the 3 terms are just quadratic forms, they are scalars i.e. just numbers. That means that each of these terms is equal to its transpose, just in case you derive it as a scalar where  $-2y^T Xa$  came from.

### Gradient of the cost function in matrix/vector form

In the last post [\[part 4\]](#) in order to find the optimal parameters  $a$  we needed the gradient of  $J(a)$  when we used the gradient descent algorithm. We discovered the gradient (or derivative) to find the "matrix" solution to the equations. **(I will do the derivation here but, this would be considered advanced even though it is a relatively easy problem ... I just happen to know matrix differential calculus :))**

**The differential and gradient of the cost function  $J(a)$**

The differential is,

$$dJ(a) = \frac{1}{2n} (dX^T Xa + (X^T X)a da - 2da^T X^T y) \\ = \frac{1}{2n} (2(X^T X)a da - 2da^T X^T y) \\ = \frac{1}{n} ((X^T X)a - y^T X) da$$

This implies that the derivative is (the stuff in front of  $da$  ... a row vector)

$$\frac{\partial J}{\partial a} = \frac{1}{n} (X^T X)a - y^T X$$

and the gradient is (the transpose of the derivative ... a column vector)

**Matrix form of the cost function gradient,**

$$\nabla_a J(a) = \frac{1}{n} (X^T X)a - X^T y$$

We could use that matrix form of the gradient vector in an optimization calculation, like gradient descent, to find the best values of  $a$ . We can also now find the solution without doing an optimization!

## Solution of the Linear Regression Least Squares Equations

The cost function  $J(a)$  is a quadratic function (its bowl shaped). That means it is "convex" and that means that it has a single minimum and that minimum is the global minimum i.e. the "lowest point of the function". At the minimum the gradient is equal to zero. We have the gradient so we can set it equal to zero and solve for  $a$

$$\nabla_a J(a) = \frac{1}{n} (X^T X)a - X^T y = 0$$

Implies that,

$$X^T Xa = X^T y$$

applying the inverse of  $X^T X$  to both sides gives,

$$(X^T X)^{-1} X^T Xa = (X^T X)^{-1} X^T y$$

therefore we have,

$$a = (X^T X)^{-1} X^T y$$

**That's it!** That's the matrix solution to find the least squares linear regression parameters  $a$ .

The term

$$(X^T X)^{-1} X^T$$

is known as the **Moore-Penrose (pseudo) inverse** of  $X$  and is sometimes written as  $X^+$ . Most numerical programming libraries have an efficient function to find that inverse. It's often something like `pinv(X)`.

## Multi-Varate Linear Regression.

Now that we have the regression equations in matrix form it is trivial to extend linear regression to the case where we have more than one feature variable in our model function. For 1 feature our model was a straight line. For 2 features the model would be a flat plane. For more than 2 features the model would be a "hyper-plane" (higher dimensional plane/surface). It is much more difficult to visualize multi-variate linear regression but all of the characteristics of the simple 1 variable case remain. The multi-variate extension can be very powerful for complicated models!

Generalizing the matrix/vector equations for the case where we have  $n$  feature variables would look like the following.

$$\hat{h}_i(X) = a_0 + a_1 x_1 + a_2 x_2 + \dots + a_n x_n$$

$$= \begin{bmatrix} 1 & x_1^{(i)} & x_2^{(i)} & \dots & x_n^{(i)} \\ 1 & x_1^{(i)} & x_2^{(i)} & \dots & x_n^{(i)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & x_2^{(n)} & \dots & x_n^{(n)} \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$$

### Multi-Varate version of the equations in matrix form

$$\hat{h}_i(X) = Xa \\ J(a) = \frac{1}{2n} (Xa - y)^T (Xa - y) \\ \nabla_a J(a) = \frac{1}{n} (X^T X)a - X^T y \\ a = (X^T X)^{-1} X^T y$$

**These formulas are the same as for the single feature variable versions!** For the multi-variate version the matrix  $X$  has a column for each feature and the vector  $a$  has a parameter for each feature (plus the constant term  $a_0$ ).

## Matrix solution applied to the King County zipcode 98019 house price dataset

Let's load up some Python tools and our dataset and give this a try.

```

In [5]: import pandas as pd      # data handling
import numpy as np             # mathematical computing
import matplotlib.pyplot as plt # plotting core
import seaborn as sns          # higher level plotting tools
import sys, os, random
sns.set()

```

```

In [6]: df_98019 = pd.read_csv("df_98019.csv")
X = df_98019[["sqft_living"]]
y = df_98019["price"]
n = len(X) # Number of data points

```

```

In [7]: x_mean = x.mean()
x_std = x.std()
X = [X - x_mean]/x_std

y_mean = y.mean()
y_std = y.std()
y = [y - y_mean]/y_std

```

The data for zipcode 98019 is now loaded and we have both the raw and mean-normalized data. This will allow us to check the results against the work that we did in the last post with the gradient descent code. We'll do the scaled data first.

```

In [8]: X = np.column_stack((np.ones((n,1)), X)) # construct the augmented matrix X
a = np.dot(np.linalg.pinv(X), y) # a = pinv(X)*y
print "a0 = %.4f, a1 = %.4f" % (a[0],a[1]) # print the values of a
a0 = -40.000000, a1 = 0.340000

```

Yes, these are the same values obtained with the gradient descent on the mean-normalized data. Now check the raw data.

```

In [9]: X = np.column_stack((np.ones((n,1)), X)) # construct the augmented matrix X
a = np.dot(np.linalg.pinv(X), y) # a = pinv(X)*y
print "a0 = %.4f, a1 = %.4f" % (a[0],a[1]) # print the values of a
a0 = -36058.062708, a1 = 422.882318

```

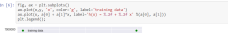
These are the values we obtained with the back transformation is the last post. Doing mean normalization is not as important when using the matrix solution.

Here's a plot of the data with these parameters for  $a$ .

```

In [10]: fig, ax = plt.subplots()
ax.plot(x,y, "o", markersize=10, label="training data")
ax.plot(X, a[0] + a[1]*x, label="f(x) = 5.24 + 5.24 * 10^4 * x[1]")
plt.legend()

```



To finish lets check the cost of an 8000 sqft house:

```

In [11]: def h(x,a): # The model function h
    h = a[0] + a[1]*x
    return h

```

```

In [12]: print "h(8000) = %f" % h(8000, a[0], a[1])
h(8000) = 455.17

```

Indeed it still can't afford it!

I hope the value of converting problems to matrix form is clear from this post. I find it is much easier to work with matrix equations and it makes implementing efficient code straight forward. It's generally the case that having formulas in matrix form and using optimized mathematical libraries will give excellent performance. In fact, the Python-numpy is compiled against Intel's excellent MKL (Math Kernel Library) this gives Python the potential for some of the performance capability of code built with lower level compiled compiled languages like C/C++ or Fortran.

In the next post I'll look up Linear Regression with a look at using non-linear feature variables (for polynomial regression). This will be a chance to illustrate the very important problem of over/under fitting data.

### Happy computing! :o)

Tags: Machine Learning, Data Science, Python, Jupyter notebooks, Programming

© Cambridge Puget Systems © Oregon Privacy Policy

Recommended: [f](#) [t](#) [g](#) [p](#) [e](#) [r](#)

Sort by: Oldest

Be the first to comment.