

## Intuition

Linear regression is an excellent method for predicting results, using a hypothesis that is computed from training data. It can filter noisy measurements as well as interpolate or extrapolate unknown datapoints. And best: It's easy to implement. This flyer will show you how.

- Use a suitable hypothesis (aka fitting function) with an initial guess for the parameters.
- Compare results to actual training data using a cost function. (This is called supervised learning.)
- Try to minimize the cost function iteratively by computing gradients and updating parameters accordingly.
- Alternatively compute the normal equation (useful for smaller systems.)

Given the limited space of the flyer, there are several (more or less fundamental) topics that could not be covered in detail. Aspects that were beyond the scope of the flyer, are

- Correlation, variance, covariance in normal equation formalism.
- Feature scaling, mean normalization.
- Implementing gradient descend; learning rate.
- Strategies for good initial guesses of parameters.
- Debugging strategies when encountering problems.

Some of these aspects will be subject of the Numerical Optimization flyer. In addition a lot of literature can be found on the net that provides a deeper insight, e.g.

- Wikipedia <sup>1</sup>
- Machine Learning @ The Stanford OpenClassroom <sup>2</sup>
- Calculating Regression Examples @ Wolfram Alpha <sup>3</sup>

The flyer can be distributed free of charge, as long as it is not altered in any way. If you have questions or comments contact the flyered science project.

fsadmin@flyered-science.org

## Numerical Solution (in Octave)

Using a numerical solver to compute  $\Theta$ , one has to define cost function (3) and gradients (5), and pass it along with an initial guess of parameters to a suitable linear solver.

In Octave the cost function for linear regression (without renormalization) can be written as follows.

```
function [J, grad] = linearCF(theta, X, y)
    J = 1/2/m * sum((theta'*X - y)*...
        (theta'*X - y)');
    grad = 1/m*(theta'*X - y)*X';
endfunction
```

Here  $\theta$  is a  $(n+1)$ -dimensional column vector,  $X$  is a  $(n+1) \times m$  matrix and  $y$  is a  $m$ -dimensional row vector.

The logistic regression cost function (without renormalization) can be written accordingly, using the logistic function.

```
function res = lFun(z)
    res = 1 ./ (1 + exp(-z));
endfunction

function [J, grad] = logisticCF(theta, X, y)
    J = -1/m * (log(lFun(theta'*X)*y') + ...
        log(1-lFun(theta'*X)) * (1-y)');
    grad = 1/m * (lFun(theta'*X)-y) * X';
endfunction
```

The optimal  $\Theta$  can then be computed using e.g. `fminunc`.

```
opt = optimset('GradObj', 'on', ...
    'MaxIter', 100);
[theta, cost] = ...
fminunc(@(t)(linearCF(t, X, y)), init, opt);
```

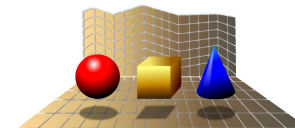
Here *init* is an initial guess for  $\Theta$ , and *opt* the solver options. (For more information on commands consult Octave help.)



# Linear Regression

An introduction to representation and algorithm,  
with implementation in Octave

The Flyered Science Series provides condensed factsheets on scientific topics. Persons new to the subject will find the intuitive introduction with links to further reading quite useful, while the experienced user benefits from the structured representation, formulary and implementation examples. Fits into every pocket, serves perfectly as bookmark in your textbooks.



<sup>1</sup>[http://en.wikipedia.org/wiki/Linear\\_regression](http://en.wikipedia.org/wiki/Linear_regression)

<sup>2</sup><http://openclassroom.stanford.edu>

<sup>3</sup>[http://www.wolframalpha.com/examples/](http://www.wolframalpha.com/examples/RegressionAnalysis.html)

## Theory

### Problem

Given a set of  $m$  datapoints  $(X^{(t)}, y^{(t)})$ , where  $y^{(t)}$  is the result of the  $i$ th measurement ( $0 \leq i \leq m$ ), depending on  $n$  variables  $X^{(t)} = (x_{(t)}^1, x_{(t)}^2, \dots, x_{(t)}^n)$ , find a linear function that best fits the data. The function can then be used to filter measurement noise or to determine unknown datapoints for predicting stock prices, weather, etc.

### Some Conventions

The following definitions are used throughout the flyer, that are common when talking about linear regression:

$m$  number of datapoints  $(X^{(t)}, y^{(t)})$   
 $n$  number of variables  $(X^{(t)} \in \mathcal{R}^n)$   
 (If  $n = 1$  the regression is called univariate, otherwise multivariate.)

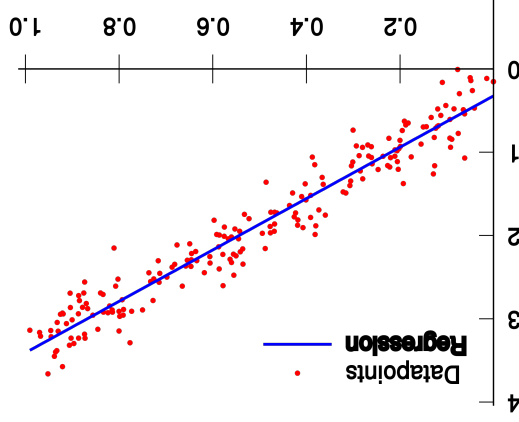


Figure 1: Example of a univariate linear regression function fitting the (red) datapoints  $(X^{(t)}, y^{(t)})$

Results  $y^{(t)}$  can be written as an  $m$ -dimensional vector  $y$ , while variables  $x_{(t)}^f$  are well represented as an  $n \times m$  dimensional matrix  $X$ . As we will see, this is often useful to simplify computation, e.g. when working with Octave.

## Representation

We will now introduce the following definitions:

$h_{\Theta}(X)$  hypothesis, i.e. fitting function  
 $\Theta$  fitting parameters ( $\Theta \in \mathcal{R}^{n+1}$ )

Ideally the hypothesis should match the results of the datapoints, i.e.

$$h_{\Theta}(X) = \Theta_0 + \Theta_1 X_1 + \dots + \Theta_n X_n = \Theta^T X \equiv y \quad (1)$$

Here  $\Theta^T$  is the transposed of  $\Theta$ . We also introduced an extra dimension for  $X$  to take into account the offset  $\Theta_0$ .  $X$  becomes an  $(n+1) \times m$  matrix.

$$X = \begin{pmatrix} 1 & x_{(1)}^1 & x_{(1)}^2 & \dots & x_{(1)}^n \\ 1 & x_{(2)}^1 & x_{(2)}^2 & \dots & x_{(2)}^n \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_{(m)}^1 & x_{(m)}^2 & \dots & x_{(m)}^n \end{pmatrix} \quad (2)$$

## Polynomial Regression

The term linear regression refers to the linear system of equations in (1). It does not mean that one cannot fit a polynomial to datapoints with linear regression.

Considering a univariate linear regression problem with variables  $X = (1, x_{(t)}^1)$  one can then construct a new variables matrix  $X^{new}$  taking into account higher orders of  $x_1$ , e.g.  $X^{new} = (1, x_{(t)}^1, x_{(t)}^2, x_{(t)}^3, \dots)$ , again applying the hypothesis of equation (1).

You can also use functions other than polynomials, like *sin*, *cos*, *exp* and so on, making the linear regression fitting very flexible.

## Binary and continuous results

So far we thought of the result vector  $y$  to hold continuous numbers (real, complex). The same representation can also be used for binary results (i.e. 0 or 1, Yes or No ...). We then speak of logistic regression rather than linear regression. There are some differences regarding how to solve logistic regression problems, that will be covered in the following.

## Cost function

The cost function  $J(\Theta)$  calculates the difference between hypothesis and actual result. Often this is implemented as a least squares approach, summing the squares of differences.

$$J(\Theta) = \frac{1}{2m} \sum_{t=1}^m (h_{\Theta}(X^{(t)}) - y^{(t)})^2 \quad (3)$$

The last line of equation (3) refers to linear regression problems. For logistic regression we choose a slightly different representation, using the logistic (or sigmoid) function.

$$h_{\Theta}(X^{(t)}) = \frac{1 + \exp(-\Theta^T X^{(t)})}{1} \quad (4)$$

The logistic function is 1 for large positive values of the exponent, while 0 for large negative values. These are the two possible results we expect for logistic regression problems. The parameters  $\Theta$  are chosen accordingly.

## Normal Equation

When solving for the optimal set of parameters  $\Theta$  one is looking for the minimum of the cost function. This implies the partial derivatives to vanish:

$$\frac{\partial J(\Theta)}{\partial \Theta_f} = 0 \quad (5)$$

This condition is met by the normal equation, where

$$\Theta = (X^T X)^{-1} X^T y \quad (6)$$

Calculating the inverse of  $X^T X$  is computationally costly, so this approach is useful mainly for small systems.

## Regularization

Sometimes  $h_{\Theta}$  fits the training data well, but is not smooth enough between datapoints. To avoid this a penalty term is introduced.

$$J(\Theta) \rightarrow J(\Theta) + \lambda \sum_{n=1}^f \Theta_n^2 \quad (7)$$

With this, small values of the parameters  $\Theta$  are favoured, causing a smoother approximation and less overfitting.