

## Choose a Multilayer Neural Network Training Function

It is very difficult to know which training algorithm will be the fastest for a given problem. It depends on many factors, including the complexity of the problem, the number of data points in the training set, the number of weights and biases in the network, the error goal, and whether the network is being used for pattern recognition (discriminant analysis) or function approximation (regression). This section compares the various training algorithms. Feedforward networks are trained on six different problems. Three of the problems fall in the pattern recognition category and the three others fall in the function approximation category. Two of the problems are simple "toy" problems, while the other four are "real world" problems. Networks with a variety of different architectures and complexities are used, and the networks are trained to a variety of different accuracy levels.

The following table lists the algorithms that are tested and the acronyms used to identify them.

| Acronym | Algorithm                | Description                                   |
|---------|--------------------------|---|
| LM      | <a href="#">trainlm</a>  | Levenberg-Marquardt                           |
| BFG     | <a href="#">trainbfg</a> | BFGS Quasi-Newton                             |
| RP      | <a href="#">trainrp</a>  | Resilient Backpropagation                     |
| SCG     | <a href="#">trainscg</a> | Scaled Conjugate Gradient                     |
| CGB     | <a href="#">traincgb</a> | Conjugate Gradient with Powell/Beale Restarts |
| CGF     | <a href="#">traincgf</a> | Fletcher-Powell Conjugate Gradient            |
| CGP     | <a href="#">traincgp</a> | Polak-Ribière Conjugate Gradient              |
| OSS     | <a href="#">trainoss</a> | One Step Secant                               |
| GDX     | <a href="#">traingdx</a> | Variable Learning Rate Backpropagation        |

The following table lists the six benchmark problems and some characteristics of the networks, training processes, and computers used.

| Problem Title | Problem Type           | Network Structure | Error Goal | Computer            |
|---------------|------------------------|-------------------|------------|---------------------|
| SIN           | Function approximation | 1-5-1             | 0.002      | Sun Sparc 2         |
| PARITY        | Pattern recognition    | 3-10-10-1         | 0.001      | Sun Sparc 2         |
| ENGINE        | Function approximation | 2-30-2            | 0.005      | Sun Enterprise 4000 |
| CANCER        | Pattern recognition    | 9-5-5-2           | 0.012      | Sun Sparc 2         |
| CHOLESTEROL   | Function approximation | 21-15-3           | 0.027      | Sun Sparc 20        |
| DIABETES      | Pattern recognition    | 8-15-15-2         | 0.05       | Sun Sparc 20        |

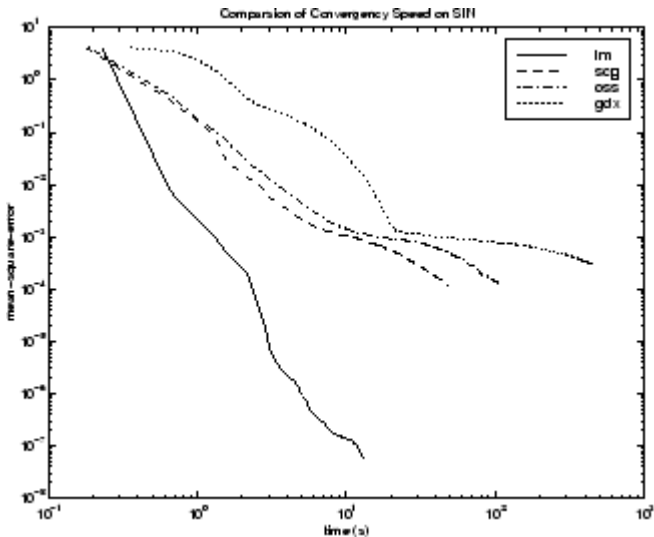
### SIN Data Set

The first benchmark data set is a simple function approximation problem. A 1-5-1 network, with [tansig](#) transfer functions in the hidden layer and a linear transfer function in the output layer, is used to approximate a single period of a sine wave. The following table summarizes the results of training the network using nine different training algorithms. Each entry in the table represents 30 different trials, where different random initial weights are used in each trial. In each case, the network is trained until the squared error is less than 0.002. The fastest algorithm for this problem is the Levenberg-Marquardt algorithm. On the average, it is over four times faster than the next fastest algorithm. This is the type of problem for which the LM algorithm is best suited—a function approximation problem where the network has fewer than one hundred weights and the approximation must be very accurate.

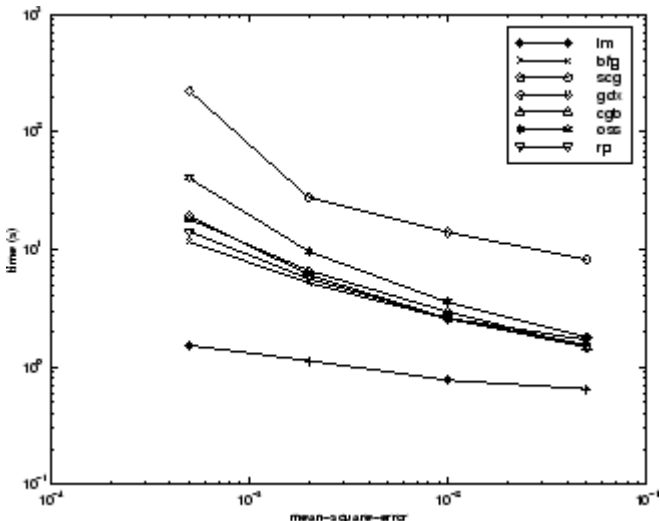
| Algorithm | Mean Time (s) | Ratio | Min. Time (s) | Max. Time (s) | Std. (s) |
|-----------|---------------|-------|---------------|---------------|----------|
| LM        | 1.14          | 1.00  | 0.65          | 1.83          | 0.38     |
| BFG       | 5.22          | 4.58  | 3.17          | 14.38         | 2.08     |
| RP        | 5.67          | 4.97  | 2.66          | 17.24         | 3.72     |
| SCG       | 6.09          | 5.34  | 3.18          | 23.64         | 3.81     |

| Algorithm | Mean Time (s) | Ratio | Min. Time (s) | Max. Time (s) | Std. (s) |
|-----------|---------------|-------|---------------|---------------|----------|
| CGB       | 6.61          | 5.80  | 2.99          | 23.65         | 3.67     |
| CGF       | 7.86          | 6.89  | 3.57          | 31.23         | 4.76     |
| CGP       | 8.24          | 7.23  | 4.07          | 32.32         | 5.03     |
| OSS       | 9.64          | 8.46  | 3.97          | 59.63         | 9.79     |
| GDX       | 27.69         | 24.29 | 17.21         | 258.15        | 43.65    |

The performance of the various algorithms can be affected by the accuracy required of the approximation. This is shown in the following figure, which plots the mean square error versus execution time (averaged over the 30 trials) for several representative algorithms. Here you can see that the error in the LM algorithm decreases much more rapidly with time than the other algorithms shown.



The relationship between the algorithms is further illustrated in the following figure, which plots the time required to converge versus the mean square error convergence goal. Here you can see that as the error goal is reduced, the improvement provided by the LM algorithm becomes more pronounced. Some algorithms perform better as the error goal is reduced (LM and BFG), and other algorithms degrade as the error goal is reduced (OSS and GDX).



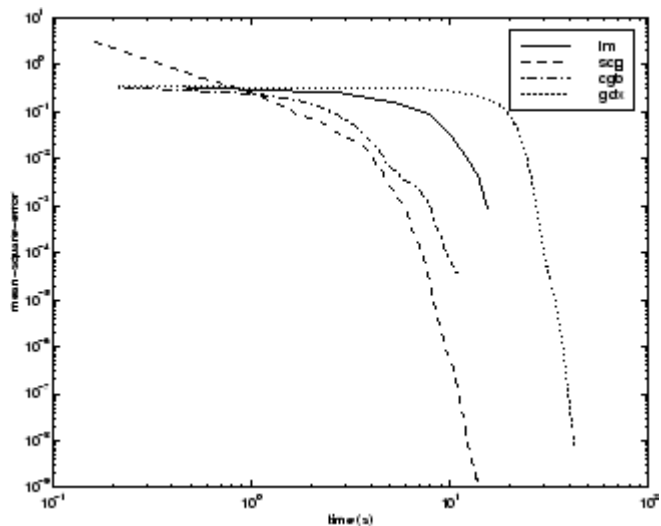
## PARITY Data Set

The second benchmark problem is a simple pattern recognition problem—detect the parity of a 3-bit number. If the number of ones in the input pattern is odd, then the network should output a 1; otherwise, it should output a -1. The network used for this problem is a 3-10-10-1 network with tansig neurons in each layer. The following table summarizes the results of training this network with the nine different algorithms. Each entry in the table represents 30 different trials, where different random initial weights are used in each trial. In each case, the network is trained until the squared error is less than 0.001. The fastest algorithm for this problem is the resilient backpropagation algorithm, although the conjugate gradient algorithms (in particular, the scaled conjugate gradient algorithm) are almost as fast. Notice that the LM algorithm does not perform well on this problem. In general, the LM algorithm does not perform as well on pattern

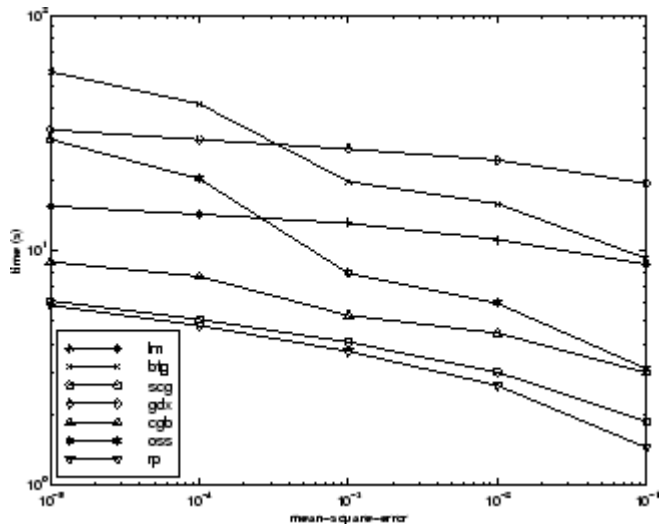
recognition problems as it does on function approximation problems. The LM algorithm is designed for least squares problems that are approximately linear. Because the output neurons in pattern recognition problems are generally saturated, you will not be operating in the linear region.

| Algorithm | Mean Time (s) | Ratio | Min. Time (s) | Max. Time (s) | Std. (s) |
|-----------|---------------|-------|---------------|---------------|----------|
| RP        | 3.73          | 1.00  | 2.35          | 6.89          | 1.26     |
| SCG       | 4.09          | 1.10  | 2.36          | 7.48          | 1.56     |
| CGP       | 5.13          | 1.38  | 3.50          | 8.73          | 1.05     |
| CGB       | 5.30          | 1.42  | 3.91          | 11.59         | 1.35     |
| CGF       | 6.62          | 1.77  | 3.96          | 28.05         | 4.32     |
| OSS       | 8.00          | 2.14  | 5.06          | 14.41         | 1.92     |
| LM        | 13.07         | 3.50  | 6.48          | 23.78         | 4.96     |
| BFG       | 19.68         | 5.28  | 14.19         | 26.64         | 2.85     |
| GDX       | 27.07         | 7.26  | 25.21         | 28.52         | 0.86     |

As with function approximation problems, the performance of the various algorithms can be affected by the accuracy required of the network. This is shown in the following figure, which plots the mean square error versus execution time for some typical algorithms. The LM algorithm converges rapidly after some point, but only after the other algorithms have already converged.



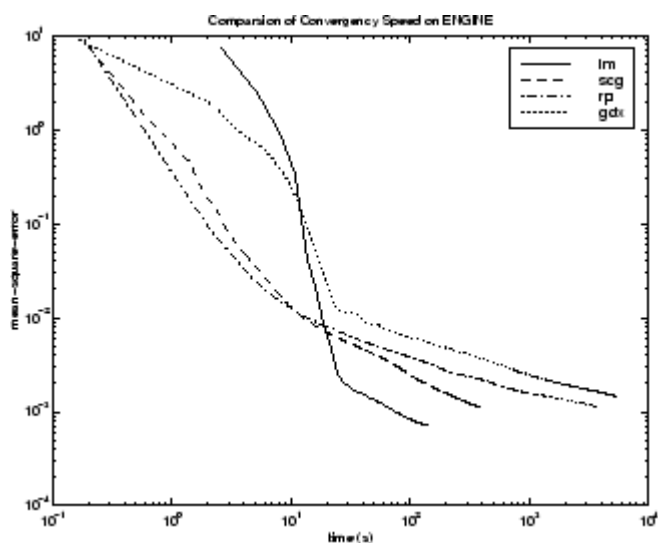
The relationship between the algorithms is further illustrated in the following figure, which plots the time required to converge versus the mean square error convergence goal. Again you can see that some algorithms degrade as the error goal is reduced (OSS and BFG).



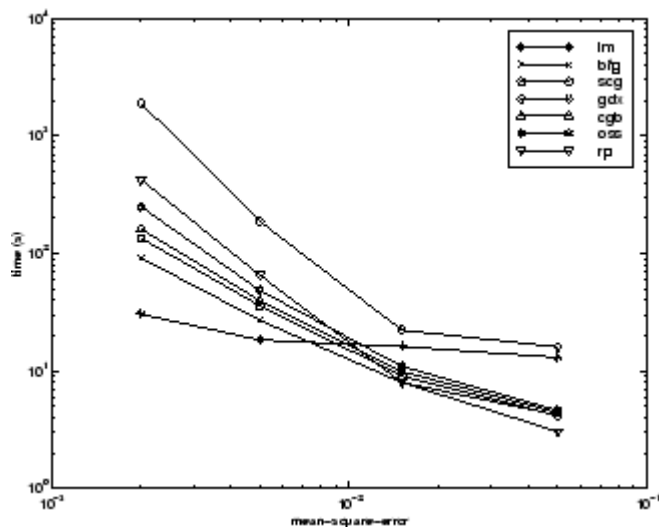
The third benchmark problem is a realistic function approximation (or nonlinear regression) problem. The data is obtained from the operation of an engine. The inputs to the network are engine speed and fueling levels and the network outputs are torque and emission levels. The network used for this problem is a 2-30-2 network with tansig neurons in the hidden layer and linear neurons in the output layer. The following table summarizes the results of training this network with the nine different algorithms. Each entry in the table represents 30 different trials (10 trials for RP and GDX because of time constraints), where different random initial weights are used in each trial. In each case, the network is trained until the squared error is less than 0.005. The fastest algorithm for this problem is the LM algorithm, although the BFGS quasi-Newton algorithm and the conjugate gradient algorithms (the scaled conjugate gradient algorithm in particular) are almost as fast. Although this is a function approximation problem, the LM algorithm is not as clearly superior as it was on the SIN data set. In this case, the number of weights and biases in the network is much larger than the one used on the SIN problem (152 versus 16), and the advantages of the LM algorithm decrease as the number of network parameters increases.

| Algorithm | Mean Time (s) | Ratio | Min. Time (s) | Max. Time (s) | Std. (s) |
|-----------|---------------|-------|---------------|---------------|----------|
| LM        | 18.45         | 1.00  | 12.01         | 30.03         | 4.27     |
| BFG       | 27.12         | 1.47  | 16.42         | 47.36         | 5.95     |
| SCG       | 36.02         | 1.95  | 19.39         | 52.45         | 7.78     |
| CGF       | 37.93         | 2.06  | 18.89         | 50.34         | 6.12     |
| CGB       | 39.93         | 2.16  | 23.33         | 55.42         | 7.50     |
| CGP       | 44.30         | 2.40  | 24.99         | 71.55         | 9.89     |
| OSS       | 48.71         | 2.64  | 23.51         | 80.90         | 12.33    |
| RP        | 65.91         | 3.57  | 31.83         | 134.31        | 34.24    |
| GDX       | 188.50        | 10.22 | 81.59         | 279.90        | 66.67    |

The following figure plots the mean square error versus execution time for some typical algorithms. The performance of the LM algorithm improves over time relative to the other algorithms.



The relationship between the algorithms is further illustrated in the following figure, which plots the time required to converge versus the mean square error convergence goal. Again you can see that some algorithms degrade as the error goal is reduced (GDX and RP), while the LM algorithm improves.



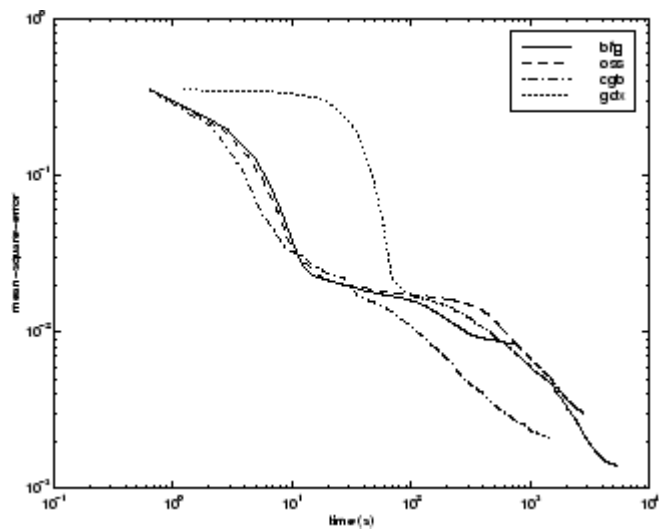
## CANCER Data Set

The fourth benchmark problem is a realistic pattern recognition (or nonlinear discriminant analysis) problem. The objective of the network is to classify a tumor as either benign or malignant based on cell descriptions gathered by microscopic examination. Input attributes include clump thickness, uniformity of cell size and cell shape, the amount of marginal adhesion, and the frequency of bare nuclei. The data was obtained from the University of Wisconsin Hospitals, Madison, from Dr. William H. Wolberg. The network used for this problem is a 9-5-5-2 network with tansig neurons in all layers. The following table summarizes the results of training this network with the nine different algorithms. Each entry in the table represents 30 different trials, where different random initial weights are used in each trial. In each case, the network is trained until the squared error is less than 0.012. A few runs failed to converge for some of the algorithms, so only the top 75% of the runs from each algorithm were used to obtain the statistics.

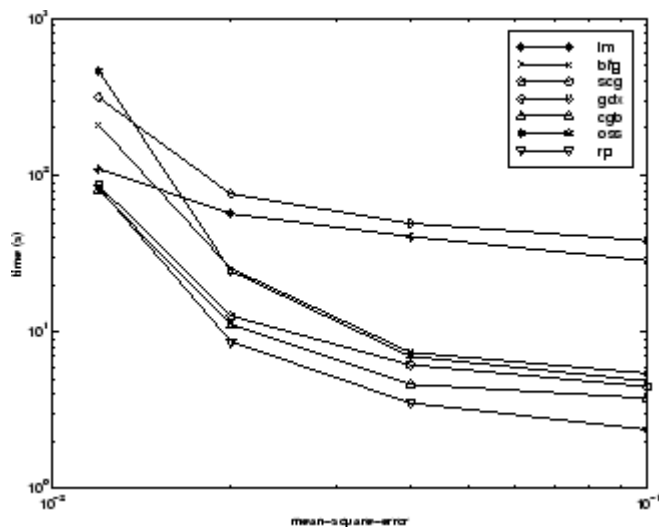
The conjugate gradient algorithms and resilient backpropagation all provide fast convergence, and the LM algorithm is also reasonably fast. As with the parity data set, the LM algorithm does not perform as well on pattern recognition problems as it does on function approximation problems.

| Algorithm | Mean Time (s) | Ratio | Min. Time (s) | Max. Time (s) | Std. (s) |
|-----------|---------------|-------|---------------|---------------|----------|
| CGB       | 80.27         | 1.00  | 55.07         | 102.31        | 13.17    |
| RP        | 83.41         | 1.04  | 59.51         | 109.39        | 13.44    |
| SCG       | 86.58         | 1.08  | 41.21         | 112.19        | 18.25    |
| CGP       | 87.70         | 1.09  | 56.35         | 116.37        | 18.03    |
| CGF       | 110.05        | 1.37  | 63.33         | 171.53        | 30.13    |
| LM        | 110.33        | 1.37  | 58.94         | 201.07        | 38.20    |
| BFG       | 209.60        | 2.61  | 118.92        | 318.18        | 58.44    |
| GDX       | 313.22        | 3.90  | 166.48        | 446.43        | 75.44    |
| OSS       | 463.87        | 5.78  | 250.62        | 599.99        | 97.35    |

The following figure plots the mean square error versus execution time for some typical algorithms. For this problem there is not as much variation in performance as in previous problems.



The relationship between the algorithms is further illustrated in the following figure, which plots the time required to converge versus the mean square error convergence goal. Again you can see that some algorithms degrade as the error goal is reduced (OSS and BFG) while the LM algorithm improves. It is typical of the LM algorithm on any problem that its performance improves relative to other algorithms as the error goal is reduced.



## CHOLESTEROL Data Set

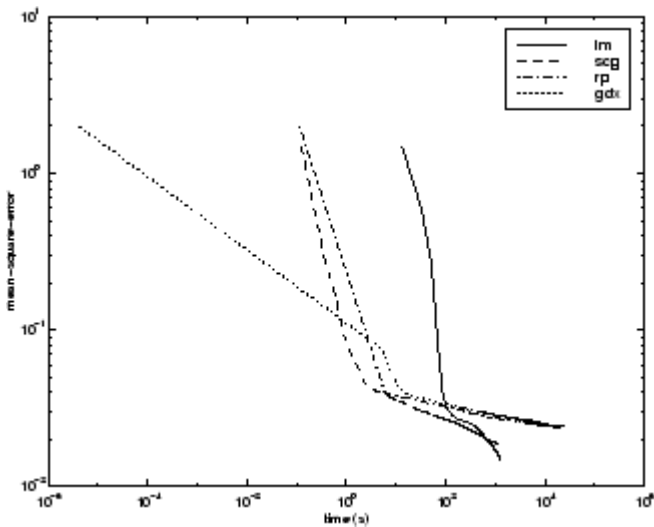
The fifth benchmark problem is a realistic function approximation (or nonlinear regression) problem. The objective of the network is to predict cholesterol levels (ldl, hdl, and vldl) based on measurements of 21 spectral components. The data was obtained from Dr. Neil Purdie, Department of Chemistry, Oklahoma State University [PuLu92]. The network used for this problem is a 21-15-3 network with tansig neurons in the hidden layers and linear neurons in the output layer. The following table summarizes the results of training this network with the nine different algorithms. Each entry in the table represents 20 different trials (10 trials for RP and GDX), where different random initial weights are used in each trial. In each case, the network is trained until the squared error is less than 0.027.

The scaled conjugate gradient algorithm has the best performance on this problem, although all the conjugate gradient algorithms perform well. The LM algorithm does not perform as well on this function approximation problem as it did on the other two. That is because the number of weights and biases in the network has increased again (378 versus 152 versus 16). As the number of parameters increases, the computation required in the LM algorithm increases geometrically.

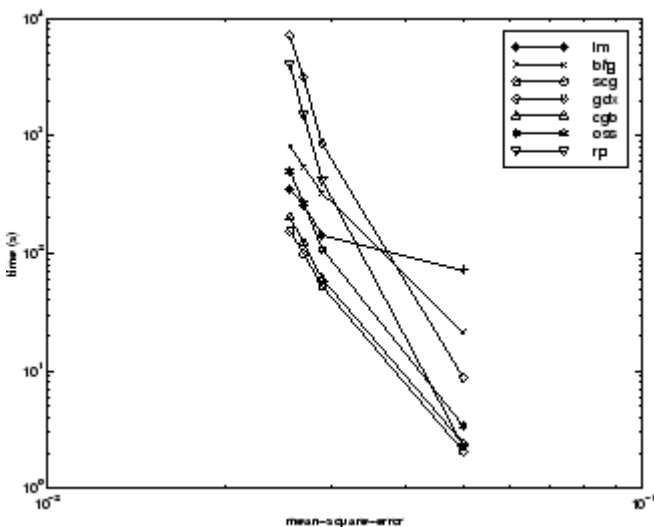
| Algorithm | Mean Time (s) | Ratio | Min. Time (s) | Max. Time (s) | Std. (s) |
|-----------|---------------|-------|---------------|---------------|----------|
| SCG       | 99.73         | 1.00  | 83.10         | 113.40        | 9.93     |
| CGP       | 121.54        | 1.22  | 101.76        | 162.49        | 16.34    |
| CGB       | 124.06        | 1.2   | 107.64        | 146.90        | 14.62    |
| CGF       | 136.04        | 1.36  | 106.46        | 167.28        | 17.67    |
| LM        | 261.50        | 2.62  | 103.52        | 398.45        | 102.06   |
| OSS       | 268.55        | 2.69  | 197.84        | 372.99        | 56.79    |

| Algorithm | Mean Time (s) | Ratio | Min. Time (s) | Max. Time (s) | Std. (s) |
|-----------|---------------|-------|---------------|---------------|----------|
| BFG       | 550.92        | 5.52  | 471.61        | 676.39        | 46.59    |
| RP        | 1519.00       | 15.23 | 581.17        | 2256.10       | 557.34   |
| GDX       | 3169.50       | 31.78 | 2514.90       | 4168.20       | 610.52   |

The following figure plots the mean square error versus execution time for some typical algorithms. For this problem, you can see that the LM algorithm is able to drive the mean square error to a lower level than the other algorithms. The SCG and RP algorithms provide the fastest initial convergence.



The relationship between the algorithms is further illustrated in the following figure, which plots the time required to converge versus the mean square error convergence goal. You can see that the LM and BFG algorithms improve relative to the other algorithms as the error goal is reduced.



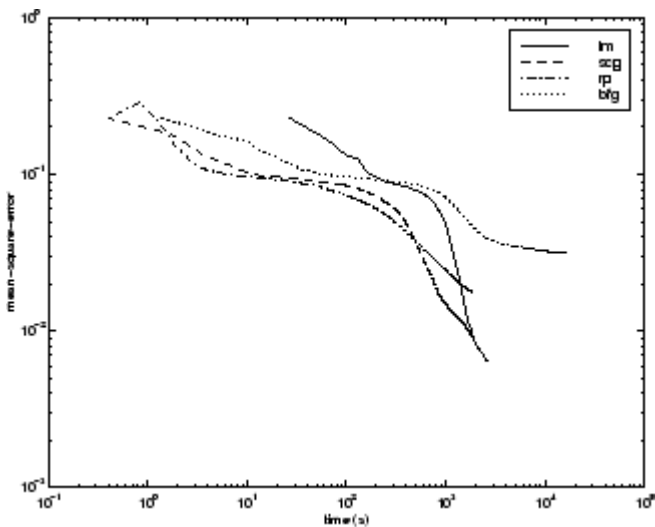
## DIABETES Data Set

The sixth benchmark problem is a pattern recognition problem. The objective of the network is to decide whether an individual has diabetes, based on personal data (age, number of times pregnant) and the results of medical examinations (e.g., blood pressure, body mass index, result of glucose tolerance test, etc.). The data was obtained from the University of California, Irvine, machine learning data base. The network used for this problem is an 8-15-15-2 network with tansig neurons in all layers. The following table summarizes the results of training this network with the nine different algorithms. Each entry in the table represents 10 different trials, where different random initial weights are used in each trial. In each case, the network is trained until the squared error is less than 0.05.

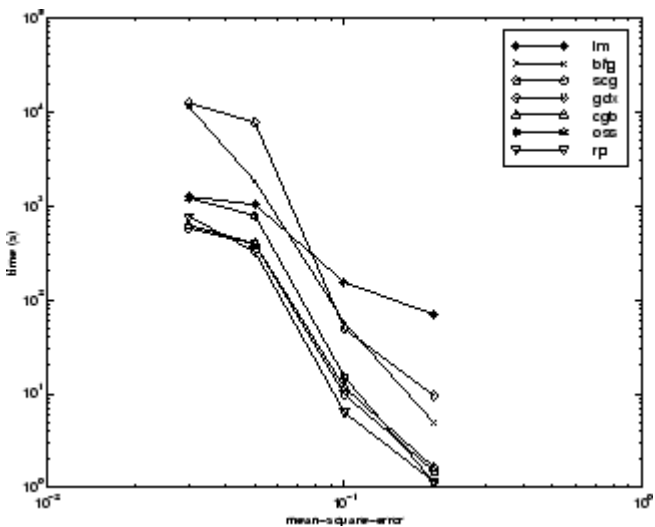
The conjugate gradient algorithms and resilient backpropagation all provide fast convergence. The results on this problem are consistent with the other pattern recognition problems considered. The RP algorithm works well on all the pattern recognition problems. This is reasonable, because that algorithm was designed to overcome the difficulties caused by training with sigmoid functions, which have very small slopes when operating far from the center point. For pattern recognition problems, you use sigmoid transfer functions in the output layer, and you want the network to operate at the tails of the sigmoid function.

| Algorithm | Mean Time (s) | Ratio | Min. Time (s) | Max. Time (s) | Std. (s) |
|-----------|---------------|-------|---------------|---------------|----------|
| RP        | 323.90        | 1.00  | 187.43        | 576.90        | 111.37   |
| SCG       | 390.53        | 1.21  | 267.99        | 487.17        | 75.07    |
| CGB       | 394.67        | 1.22  | 312.25        | 558.21        | 85.38    |
| CGP       | 415.90        | 1.28  | 320.62        | 614.62        | 94.77    |
| OSS       | 784.00        | 2.42  | 706.89        | 936.52        | 76.37    |
| CGF       | 784.50        | 2.42  | 629.42        | 1082.20       | 144.63   |
| LM        | 1028.10       | 3.17  | 802.01        | 1269.50       | 166.31   |
| BFG       | 1821.00       | 5.62  | 1415.80       | 3254.50       | 546.36   |
| GDX       | 7687.00       | 23.73 | 5169.20       | 10350.00      | 2015.00  |

The following figure plots the mean square error versus execution time for some typical algorithms. As with other problems, you see that the SCG and RP have fast initial convergence, while the LM algorithm is able to provide smaller final error.



The relationship between the algorithms is further illustrated in the following figure, which plots the time required to converge versus the mean square error convergence goal. In this case, you can see that the BFG algorithm degrades as the error goal is reduced, while the LM algorithm improves. The RP algorithm is best, except at the smallest error goal, where SCG is better.



## Summary

There are several algorithm characteristics that can be deduced from the experiments described. In general, on function approximation problems, for networks that contain up to a few hundred weights, the Levenberg-Marquardt algorithm will have the fastest convergence. This advantage is especially noticeable if very accurate training is required. In many cases, `trainlm` is able to obtain lower mean square errors than any of the other algorithms tested.



However, as the number of weights in the network increases, the advantage of `trainlm` decreases. In addition, `trainlm` performance is relatively poor on pattern recognition problems. The storage requirements of `trainlm` are larger than the other algorithms tested. By adjusting the `mem_reduc` parameter, discussed earlier, the storage requirements can be reduced, but at the cost of increased execution time.

The `trainrp` function is the fastest algorithm on pattern recognition problems. However, it does not perform well on function approximation problems. Its performance also degrades as the error goal is reduced. The memory requirements for this algorithm are relatively small in comparison to the other algorithms considered.

The conjugate gradient algorithms, in particular `trainscg`, seem to perform well over a wide variety of problems, particularly for networks with a large number of weights. The SCG algorithm is almost as fast as the LM algorithm on function approximation problems (faster for large networks) and is almost as fast as `trainrp` on pattern recognition problems. Its performance does not degrade as quickly as `trainrp` performance does when the error is reduced. The conjugate gradient algorithms have relatively modest memory requirements.

The performance of `trainbfg` is similar to that of `trainlm`. It does not require as much storage as `trainlm`, but the computation required does increase geometrically with the size of the network, because the equivalent of a matrix inverse must be computed at each iteration.

The variable learning rate algorithm `traingdx` is usually much slower than the other methods, and has about the same storage requirements as `trainrp`, but it can still be useful for some problems. There are certain situations in which it is better to converge more slowly. For example, when using early stopping you can have inconsistent results if you use an algorithm that converges too quickly. You might overshoot the point at which the error on the validation set is minimized.

---