A fast Stokes inversion technique based on quadratic regression

Fei Teng and Yuan-Yong Deng

Key Laboratory of Solar Activity, National Astronomical Observatories, Chinese Academy of Sciences, Beijing 100012, China; *fteng@bao.ac.cn, dyy@bao.ac.cn*

Received 2015 May 27; accepted 2015 September 7

Abstract Stokes inversion calculation is a key process in resolving polarization information on radiation from the Sun and obtaining the associated vector magnetic fields. Even in the cases of simple local thermodynamic equilibrium (LTE) and where the Milne-Eddington approximation is valid, the inversion problem may not be easy to solve. The initial values for the iterations are important in handling the case with multiple minima. In this paper, we develop a fast inversion technique without iterations. The time taken for computation is only 1/100 the time that the iterative algorithm takes. In addition, it can provide available initial values even in cases with lower spectral resolutions. This strategy is useful for a filter-type Stokes spectrograph, such as SDO/HMI and the developed two-dimensional real-time spectrograph (2DS).

Key words: Sun: magnetic fields — methods: statistical — methods: data analysis — techniques: polarimetric

1 INTRODUCTION

Stokes inversion calculation is a key process in resolving polarization information on radiation from the Sun and calculating the associated vector magnetic fields, which are important for studying the evolution, radiation and the fundamental causes of various solar activities.

Even for the cases in local thermodynamic equilibrium (LTE) and where the Milne-Eddington approximation is valid, the inversion problem may not be easy to solve. The most commonly used method is nonlinear least squares fitting. Two main instruments for measurement of vectorial magnetic fields, the HINODE/spectro-polarimeter (SP) and Solar Dynamics Observatory(SDO)/Helioseismic and Magnetic Imager (HMI), both depend on this method. Related references can be found in Auer et al. (1977), Landolfi et al. (1984), Skumanich & Lites (1987), Lites et al. (1988) and Su & Zhang (2004).

According to the work by Wittmann (1974), Lites et al. (1988) and Rees et al. (1989), the radiative transfer equations for Stokes parameters can be written as

$$\frac{dI}{d\tau_c} = KI - j, \qquad (1)$$

where I denotes the four Stokes parameters I, Q, U and V, τ_c is the optical depth, K is the total absorption matrix, and j is the total emission vector.

In the simple cases of LTE and where the Milne-Eddington approximation is applicable, we can solve these equations analytically and generate Stokes profiles with the following parameters,

– Magnitude of magnetic field *B*;

- Doppler broadening $\Delta \lambda_{\rm D}$;
- The line-of-sight component of macro velocity v_{los} ;
- Inclination ψ ;
- Azimuthal angle ϕ ;
- Ratio between the absorption coefficient at the line center and the continuous extinction coefficient at the reference wavelength η_0 ;
- S_0 and S_1 in the representation of the source function $S_C = S_0 + S_1 \tau_c$ under the Milne-Eddington approximation.

For example, Figure 1 shows a comparison between synthetic and observed profiles for HINODE/SP data (Kosugi et al. 2007).

The profiles are extracted from a sunspot at the boundary between the penumbra and umbra of active region AR 10930, and the spectral position is close to 6302.493 Å. A pixel with $V/I \approx 60\%$ is selected to make V look clearer. Some experts on instruments may think this is impossible, but this truly comes from HINODE/SP data. Actually, the scale of V/I does not influence the inversion procedure much, which is the main topic of this paper, so we just trust this data point.

Now we need to consider the inversion procedure with respect to the synthetic process. The objective is to obtain the above input parameters from observed Stokes profiles. The least squares fitting method can be applied during the inversion process, with the objective function,

$$\chi^{2} = \sum_{i=1}^{4} \sum_{j=1}^{m} \frac{1}{w_{i}^{2}} \left[y_{ij} - F_{ij} \left(x \right) \right]^{2}, \qquad (2)$$



Fig.1 Comparison between synthetic and observed Stokes profiles.

where y represents the observed values, x stands for the unknown physical parameters as mentioned above, F signifies the synthetic Stokes parameters, and w stands for the weights of I, Q, U and V. Here index i represents I, Q, U and V, index j stands for the different wavelengths and p signifies the solar atmosphere parameters.

As an example, the program VFISV (Borrero et al. 2011) which is adopted by SDO/HMI uses the modified Levenberg-Marquardt algorithm to solve the nonlinear least squares problem. This algorithm is one kind of iterative optimization method, most of which relies on local linearization. This kind of method usually ignores the global information during the iteration and does not behave so well for problems with multiple minima. The empirical evidence provided by Teng & Deng (2014) demonstrates that the most effective method to deal with the multi-minima problem so far is a random-jump strategy. When the iteration goes to a local minimum, by this strategy, the algorithm will allow a random jump to some point close to a potentially smaller minimum. However, a fixed amount of iterations are required to make the random-jump effective. They cost much more computing time than finding one local minimum.

So with the motive of developing a fast method to handle the huge amount of data coming from future instruments, another kind of method which relies on a set of training samples is under consideration. Rees et al. (2000) and Socas-Navarro et al. (2001) implemented principal component analysis to reduce the dimension of the Stokes profile. They considered the synthetic spectra associated with known model parameters as training data and used the nearest neighbor method to find the closest model which can fit the observation best. Carroll & Staude (2001), Socas-Navarro (2005) and Carroll et al. (2008) then applied artificial neural networks, which are another artificial intelligence technique, to the analysis of Stokes profiles. They trained the network using the inversion model for the magnetic field vector, velocity and temperature classification by a synthetic sample set. They applied a preprocessing stage to project the observed profile onto the hypersurface defined by the synthetic profile before substitution in the model.

In this paper, we will apply quadratic regression, which is another kind of sample-based approach, to this problem. The main structure in the magnetogram can be captured by this approach much faster than the iterative method. If more accuracy is needed, it can also provide available initial values for the iterations. The structure of this paper is depicted as follows. In Section 2, we give a short introduction to the iterative optimization algorithm that VFISV uses. In Section 3, we describe the quadratic regression methods and illustrate some comparisons between the results obtained by classical iterations and quadratic regressions. In Section 4, we apply the results gained by the quadratic regressions to the initial values for the iterative algorithm described in Section 2 and find an obvious improvement over the fixed initial values. Finally, we will conclude in Section 5 with a discussion of this approach.

2 THE ITERATIVE ALGORITHM

The iterative algorithm used by VFISV to solve the nonlinear least squares problem is the modified Levenberg-Marquardt algorithm listed below:

- (1) Set $x := x_0, \tau := \tau_0$ and n = 0.
- (2) If $n > n_{\text{max}}$, algorithm terminates.
- (3) Calculate the synthetic Stokes profile F(x) and its Jacobi matrix J(x) with respect to x.
- (4) Calculate the objective function χ² and its gradient g with respect to x.

- (5) Solve the equation $(J^T J + \tau D)s = -q$, where D is a diagonal matrix with the same diagonal element as $J^T J$.
- (6) Set $x^{\text{pre}} = x + s$, again calculate the objective function χ^{pre} at x^{pre} , and set n = 0.
 - If $\chi^{\rm pre} < \chi$, set $x = x^{\rm pre}$, decrease τ and go to step 2 for the next iteration.
 - If $n \ge n_{\max}$, set s as a random step in some range, set x = x + s, and go to step 2.
 - If $\chi^{\text{pre}} \geq \chi$, increase τ , set n = n + 1 and go to step 5.

During this algorithm, two additional limiting procedures are considered. For each step,

- If s is too long, shrink it to satisfy $|s| < s^{\lim}$ to prevent
- the solution from being too far from the initial point. If $x_i^{\text{pre}} < x_i^{\min}$ (or $x_i^{\text{pre}} > x_i^{\max}$), set $x_i = x_i^{\min}$ (or $x_i = x_i^{\min}$) and fix x_i in the following iterations until the next jumping to make the results meaningful in physics, where the limit values are shown in Table 1, in which I_c represents the maximal Stokes I value along the spectral line.

The classical Levenberg-Marquardt algorithm has a global convergence property which means the iterations must converge to a stable point, wherever we set the initial value. However, this parameter space has many local minima where the algorithm can converge. It does not perform so well for a problem with multi-minima. So, a randomjump strategy is considered at Step 6 in the modified algorithm, by which the iteration will more probably converge to a global minimum than the classical algorithm. However, it suffers from being time-consuming, because a fixed amount of iterations are required for each pixel.

Figure 2 is a comparison between the results of the classical and modified algorithm for a quiet region of the Sun. The classical algorithm may cause convergence to a local minimum and cause many bad pixels.

Figure 3 presents inversion results for an active region by the modified algorithm. The observed data come from HINODE/SP.

3 A FAST INVERSION METHOD BY USING QUADRATIC REGRESSION

In this section we will introduce a fast inversion method by using quadratic regression. This kind of approach is based on a database training process, which relies on some preselected samples with known inversion results and Stokes parameters. In this paper, selected inversion results by the iterative algorithm are taken into account as the training database. However, this approach does not directly depend on the physical model that is used to generate the Stokes profiles, and the properly selected samples are the only resources used to in the method.

The key point is to consider each unknown parameter x as a quadratic function of the whole observed Stokes profile y which can be expressed as

$$x = f_{H,a,b}(y) = y^T H y + a^T y + b.$$
 (3)

For example, if we consider 10 points along the wavelength of a spectral line, the vector y has 40 components, the coefficient vector a has 40 components, b is a scalar coefficient and the symmetric coefficient matrix H has 820 different components that are located in the upper triangular part.

A training process is needed to pre-calculate the coefficients in H, a and b. Least square errors are minimized between the selected training database and the parameters calculated by the quadratic expressions. Actually, we need to solve the optimization problem

$$\arg\min_{f\in P_2}\sum_{j=1}^m \left[f(y_j) - x_j\right]^2,\tag{4}$$

where P_2 is the function space composed of all the quadratic functions, and $\{(x_i, y_i), j = 1, \dots, m\}$ is the training set composed of independent and identically distributed samples and assumed to be randomly drawn from an unknown distribution of (x, y).

By the law of large numbers, the objective function in Equation (4) must converge to the expectation

$$E\left[\left(f\left(y\right)-x\right)^{2}\right],$$
(5)

called risk as $m \to \infty$. So, the statistical basis of this method is to find a suitable function f to obtain the minimal risk, which is based on Bayesian decision theory.

Notice that this training process is easy to solve, because it is a linear least squares problem. After the training process, the inversion results can be obtained within about 1/100 of the computing time required for the iterative algorithm.

To demonstrate this approach, we apply it to the data from two instruments - HINODE/SP (Kosugi et al. 2007) and SDO/HMI (Schou et al. 2012). Before the training process is applied, a sample set is required to be compiled. Because this approach does not directly rely upon the physical model when performing the iterative fitting method described in Section 2, the samples are only used for training the neural network. Properly selected samples are the key point for the successful implementation of this approach.

Figure 4 shows the distribution of magnetic fields as measured by all pixels within 13 active regions. The values contained in about 90% of the pixels are less than 500 G. In order to ensure samples include information about the full range of magnetic fields, we select a fixed amount of samples in each interval of 100 G. This composes a training set with a total of $18\,690$ samples, which account for 3%of the data from all the active regions.

For HINODE/SP, Figure 5 shows the inversion results by the quadratic regressions for the same region and data as shown in Figure 3.

Figures 6 - 8 illustrate the results for SDO/HMI. The main structure in the magnetogram can be captured by the regressions in both cases.

Table 1 Limits used by the Modified Levenberg-Marquardt Algorithm

| Parameter | Lower limit x^{\min} | Upper limit x^{\max} | Step limit s^{\lim} |
|---------------------------------|------------------------|------------------------|-----------------------|
| B (G) | -7000 | 7000 | 2000 |
| $\Delta \lambda_{\rm D}$ (cm) | 5×10^{-11} | 8×10^{-10} | 1×10^{-10} |
| $v_{\rm los}~({\rm cm~s^{-1}})$ | $-7 	imes 10^5$ | 7×10^5 | 1×10^5 |
| ψ | No limit | No limit | 20° |
| ϕ | No limit | No limit | 20° |
| η_0 | 1 | 300 | 25 |
| S_0 | $0.05I_{c}$ | $1.5I_c$ | $0.3I_c$ |
| S_1 | $0.05I_{c}$ | $1.5I_{c}$ | $0.2I_c$ |



Fig. 2 Comparison between classical (left) and modified (right) algorithm for a quiet region.



Fig.3 Inversion results using the modified Levenberg-Marquardt algorithm applied to AR 10930 with HINODE/SP data for B, ψ , ϕ and v_{los} in Fig. 5 from (1) to (4) respectively.



Fig. 4 Distribution of the scalar magnetic fields. The right panel is a partially enlarged view of the left panel.



Fig. 5 Inversion results by using quadratic regression applied to AR 10930 with HINODE/SP data for B, ψ , ϕ and v_{los} from (1) to (4) respectively.

For accuracy, we also calculated the average errors between the inversion and trained pictures for various ranges of magnetic fields, which are listed in Table 2. In order to estimate the error bounds caused by the disturbance of the sample values, the quadratic mean or root mean square (RMS)

$$\bar{e} = \sqrt{\frac{\sum_{i=1}^{n} \left[x_i - f_{H,a,b} \left(y_i \right) \right]^2}{n}}$$
(6)

is used here, where $x_i, i = 1, ..., n$ represents the original values obtained by the iterative algorithm using all the pixels in a region, and y_i represents the Stokes parameters in the same pixel as x_i .

In detail, let \tilde{x}_i denote a perturbed value with respect to x_i and assume

$$|\tilde{x}_i - x_i| < \varepsilon. \tag{7}$$

For brevity, let f denote the resulting function trained by samples extracted from $\{x_i\}$ and \tilde{f} denote that trained by



Fig. 6 Comparison between the real and trained results for region 1 with SDO/HMI data.



Fig.7 Comparison between the real and trained results for region 2 with SDO/HMI data.



Fig. 8 Comparison between the real and trained results for region 3 with SDO/HMI data.

(8)

samples extracted from $\{\tilde{x}_i\}$. If we use the notation

 $|u|^2 = \left(\sum_{i=1}^n u_i^2\right)^{\frac{1}{2}},$

f and \tilde{f} can be approximately written as

$$f = \arg\min_{f \in P_2} T(f) = \arg\min_{f \in P_2} |f(y) - x|, \qquad (9)$$

$$\tilde{f} = \arg\min_{f \in P_2} \tilde{T}(f) = \arg\min_{f \in P_2} |f(y) - \tilde{x}|.$$
(10)



Fig. 9 The derivatives of B with respect to I, Q, U and V at different wavelengths.



Fig. 10 Illustration for lower spectral resolution when Step = 3.

Table 2 RMS of the Errors between the Real and Trained Results for Different Ranges of Magnetic Fields (G)

| Region | Parameter | (0, 1000) | (1000, 2500) | > 2500 | Full range |
|--------|-----------|-----------|--------------|----------|------------|
| 1 - | В | 1.22E+02 | 1.21E+02 | 1.22E+02 | 1.21E+02 |
| | ψ | 1.11E+01 | 9.48E+00 | 6.85E+00 | 1.05E+01 |
| 2 - | В | 2.73E+02 | 3.70E+02 | 2.52E+02 | 3.20E+02 |
| | ψ | 2.40E+01 | 2.40E+01 | 1.17E+01 | 2.35E+01 |
| 3 - | В | 1.23E+02 | 1.50E+02 | 1.36E+02 | 1.34E+02 |
| | ψ | 9.91E+00 | 8.34E+00 | 7.76E+00 | 9.24E+00 |

Actually, in the training process defined by Equation (4), the summation extends over the selected training samples, but in Equations (9) and (10), it extends over all the pixels in a region. This difference leads to the *generalization error* described in the machine learning literature, which is not easy to estimate. Cucker & Smale (2002) gave a detailed description for this. In this paper, we just neglect this difference for simplicity.

Now for all $f \in P_2$,

$$\begin{aligned} \left| T(f) - \tilde{T}(f) \right| \\ &= \left| \left| f(y) - x \right| - \left| f(y) - \tilde{x} \right| \right| \\ &\leq \left| x - \tilde{x} \right| \\ &\leq \sqrt{\sum_{i=1}^{n} \left(x_i - \tilde{x}_i \right)^2} \\ &\leq \sqrt{n\varepsilon}, \end{aligned}$$
(11)

and Equations (9), (10) and (11) imply

$$\left|T(f) - \tilde{T}(\tilde{f})\right| \le \sqrt{n\varepsilon}.$$
 (12)

Then the quadratic mean with respect to the differences between the predicted results by f and \tilde{f} can be calculated by

$$\overline{f(y) - \tilde{f}(y)} = \frac{1}{\sqrt{n}} \left| f(y) - \tilde{f}(y) \right| \\
\leq \frac{1}{\sqrt{n}} \left(T(f) + \tilde{T}(\tilde{f}) + |x - \tilde{x}| \right) \\
\leq \frac{1}{\sqrt{n}} \left(2T(f) + \left| T(f) - \tilde{T}(\tilde{f}) \right| + |x - \tilde{x}| \right). \quad (13)$$

Finally, according to Equations (6), (7) and (12),

$$\overline{f(y) - \tilde{f}(y)} \le 2\bar{e} + 2\varepsilon.$$
(14)

For example, if we assume the error of B observed by HMI is bounded by 50 G, and consider the perturbed error of B in region 1, we can obtain

$$\varepsilon = 50 \,\mathrm{G},$$
 (15)

$$\bar{e} = 121 \,\mathrm{G},\tag{16}$$



Fig. 11 Inversion results of B by using the modified Levenberg-Marquardt algorithm on decreasing spectral resolutions from (1) to (6) with fixed initial values.



Fig. 12 Inversion results of B by the modified Levenberg-Marquardt algorithm on decreasing spectral resolutions from (1) to (6) with quadratic regression using initial values.



Fig. 13 Average errors between the maximum resolution and decreasing resolutions for two initial value strategies, with solid lines for the fixed initial values and dashed lines for the quadratic regressions, and for B, ψ , ϕ and v_{los} from (1) to (4) respectively.

from Table 2, and the bound on the perturbed error is $(2 \times 121 + 2 \times 50 = 342)$ G.

We can also derive other information from the quadratic expression. For example, the derivative of each obtained parameter with respect to I, Q, U and V at each wavelength can be calculated by

$$\nabla x = 2Hy + a. \tag{17}$$

Figure 9 illustrates the derivatives of the magnetic field *B*.

4 APPLICATION OF THE FAST INVERSION METHOD TO GENERATING INITIAL VALUES

As an application of the quadratic regression, we can use it to generate initial values for the above modified Levenberg-Marquardt algorithm. One purpose of this paper is to analyze data from the instrument "Two-Dimensional Real-Time Spectrograph (2DS)" that is being developed (Deng & Zhang 2009). It is a filter-type magnetograph which can observe the polarized filtergram at eight wavelength points within a spectral line simultaneously.

Because of its low spectral resolution, the above iterative algorithm needs to be validated with data that have decreasing resolutions. So, we define a step in wavelength that decreases the resolution of HINODE/SP data.

Figure 10 is an illustration when step = 3. We tried steps from 1 to 6. Figure 11 shows the results for different spectral resolutions using a fixed initial value near the value around the center.

If we use the quadratic regression as the initial value, we can obtain the results shown by Figure 12. They obviously behave better than the fixed initial value for low spectral resolutions. Finally, from the average errors shown in Figure 13, we can more clearly see the advantage of using the initial values derived by quadratic regression. These four graphs illustrate the average errors between the maximal resolution (step = 1) and decreasing resolutions (various steps along the horizontal axis), with solid lines for a fixed initial value and dashed lines for the quadratic regression. For example, with step = 5 and a fixed initial value, an average error over 3000 G makes the magnetic fields unavailable, but with values from the quadratic regression used as initial values, the average error is about 100 G and is only reasonable for 1/5 of the sample points along the wavelength.

5 CONCLUSIONS

In this paper, we propose a fast inversion technique based on quadratic regressions. The CPU time consumed for this method is only about 1/100 what the iterative algorithm takes. If more accuracy is required, this technique can also provide initial values for the iterative algorithm.

Because this approach does not directly rely upon the physical model in the process used for the iterative fitting method, properly selected samples are essential to the training results. In order to ensure the samples include information about the full range of magnetic fields, a fixed amount of samples in each interval of 100 G are selected in this paper, composing a training set with 3% of the points in both the training and validation datasets.

For the case with low spectral resolutions, bad pixels may even occur with the random-jump strategy if an unreasonable initial value is adopted. Although the quadratic regression cannot substitute the random-jump strategy, it can provide good initial values in this situation.

Acknowledgements We thank the referee for a constructive report and Prof. Zhong-Quan Qu for great advice on the data retrieval and analysis. We also acknowledge the data provided by the space telescopes SDO/HMI and HINODE/SP. All the results in this paper are computed on a workstation funded by the Key Laboratory of Solar Activity of Chinese Academy of Sciences and the National Science Foundation. This work was supported by the National Natural Science Foundation of China (Grant Nos. 11178005 and 11427901) and the Strategic Priority Research Program of the Chinese Academy of Sciences (XDB09040200).

References

- Auer, L. H., House, L. L., & Heasley, J. N. 1977, Sol. Phys., 55, 47
- Borrero, J. M., Tomczyk, S., Kubo, M., et al. 2011, Sol. Phys., 273, 267
- Carroll, T. A., Kopf, M., & Strassmeier, K. G. 2008, A&A, 488, 781
- Carroll, T. A., & Staude, J. 2001, A&A, 378, 316

Cucker, F., & Smale, S. 2002, Bull. Am. Math. Soc., 39, 1

- Deng, Y., & Zhang, H. 2009, Science in China: Physics, Mechanics and Astronomy, 52, 1655
- Kosugi, T., Matsuzaki, K., Sakao, T., et al. 2007, Sol. Phys., 243, 3
- Landolfi, M., Landi Degl'Innocenti, E., & Arena, P. 1984, Sol. Phys., 93, 269
- Lites, B. W., Skumanich, A., Rees, D. E., & Murphy, G. A. 1988, ApJ, 330, 493
- Rees, D. E., Durrant, C. J., & Murphy, G. A. 1989, ApJ, 339, 1093
- Rees, D. E., López Ariste, A., Thatcher, J., & Semel, M. 2000, A&A, 355, 759
- Schou, J., Borrero, J. M., Norton, A. A., et al. 2012, Sol. Phys., 275, 327
- Skumanich, A., & Lites, B. W. 1987, ApJ, 322, 473
- Socas-Navarro, H. 2005, ApJ, 621, 545
- Socas-Navarro, H., López Ariste, A., & Lites, B. W. 2001, ApJ, 553, 949
- Su, J.-T., & Zhang, H.-Q. 2004, ChJAA (Chin. J. Astron. Astrophys.), 4, 365
- Teng, F., & Deng, Y.-Y. 2014, RAA (Research in Astronomy and Astrophysics), 14, 1469
- Wittmann, A. 1974, Sol. Phys., 35, 11