



Automated classification and stellar parameterization

Sunetra Giridhar

**Exploring Large Surveys for Galactic Astronomy
August 22, 2006**

**INDIAN INSTITUTE OF ASTROPHYSICS,
BANGALORE 560034, INDIA**

Introduction

The MK spectral types are classical description of stellar spectra. Although the two dimensional MK spectral type (SpT) and luminosity class (LC) are related to temperature and gravity of a star, the SpT are not assigned based upon these parameters but use visual appearance of stellar spectra. The MK classification involves comparing spectra to be classified with those of classification standards of defined class. The advantage of MK system has been that it is model independent and works well even with spectra of modest resolution. N. Houk and her collaborators have done a monumental work of determining SpT and LC for all stars in HD catalog ($\sim 12,000$ stars up to $V_{\text{mag}} \sim 11$) with RMS error of 0.6 in SpT and 0.25 in LC. This data has been used as reference for automated classifications (see for example von Hippel et al. 1994, Bailer-Jones et al. 1998).

Methods of automated spectral classification

The most commonly used automated spectral classification methods are based on (a) Minimum Distance Method (MDM) (b) Gaussian Probability Method (PDM), (c) Principal Component analysis (PCA) and (d) Artificial Neural Network (ANN). Quantitative methods involving measurement of equivalent widths of certain lines, line strength ratios etc and calibration of these quantities in terms of stellar parameters have also been used. For example, Stock and Stock (1999) used equivalent widths of 19 absorption lines, (B-V) colors and M_V derived from Hipparcos catalog for a sample of 487 stars for calibration of M_V . Their algorithms can predict M_V from these line strengths for spectral types O-M with an average error of 0.26 mag.

MDM

The classification is done by minimizing distance metric between the object to be classified and each member of a set of templates. The object is assigned the class of the template, which gives the smallest distance. If the star spectrum is represented by a vector $X = (x_1, x_2, \dots, x_i, \dots, x_N)$ and template c is represented by another vector $S = (s_{c_1}^c, s_{c_2}^c, \dots, s_{c_i}^c, \dots, s_{c_N}^c)$

the distance D_c is evaluated

$$D_c = \frac{1}{N} \left[\sum_{i=1}^{i=N} w_i^{(c)} |x_i - s_i^{(c)}|^p \right]^{1/p}$$

where $w_i^{(c)}$ is the weight assigned to spectral element i of the class c . The spectrum X is assigned class c for which D_c is minimum. The weights are assigned to spectral elements based on their relative importance in determining the spectral class. In this approach the number of templates used to define subclasses limit the accuracy of classification. Interpolation can be made to make intra-class assignment.

Katz et al. (1998) used this method with χ^2 weighing on high resolution Elodie spectra using a large number of reference stars of known T_{eff} , $\log g$ and $[M/H]$ to derive atmospheric parameters of target stars. These authors achieved accuracy of 86 K in T_{eff} , 0.28 in $\log g$ and 0.35 in $[M/H]$. Vanssevicius and Bridzius (1994) used MDS with χ^2 weighing to estimate SpT and M_V from Vilnius photometric indices. An accuracy of 0.7 was achieved for SpT and 0.8mag for M_V over spectral type range O5 to M5.

GPM

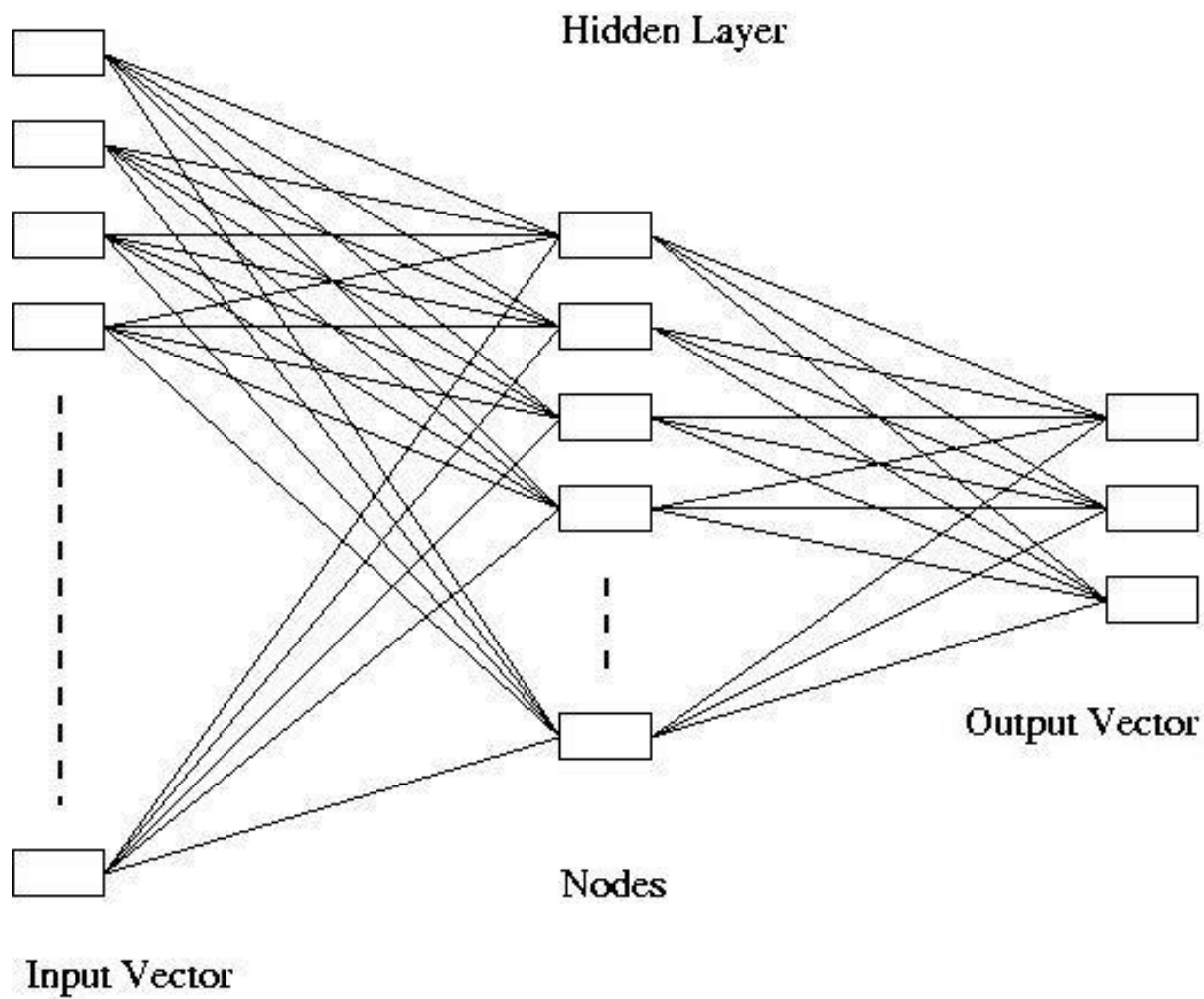
It is a statistical approach of multivariate classification explained very clearly in Bailer-Jones (2001). It was employed by Christlieb et al. (2002) to estimate T_{eff} , $\log g$ and $[M/H]$ for stars observed in objective prism Hamburg/ESO survey with an accuracy of 400 K in T_{eff} , 0.68 dex in $\log g$.

Principal Component Analysis

It is a method of representing a set of N dimensional data by means of their projection onto a set of optimally defined axes. Since these axes (Principal components) form an orthogonal set, a linear transformation of the data is achieved. Not all components are important. Components that represent large variance are important while those represent least variance can be ignored and data set can be replaced by significant components alone resulting in reduction of the data size. These compressed data sets are used as input for neural networks. Bailer-Jones et al.(1998) had demonstrated that precise calibration could be done using these compressed spectra and that the optimal compression also results in noise removal. Singh et al. (2006) have used a variation of PCA technique to restore missing data in a sample of 300 stars in Indo-US code` feed spectral library.

Neural Network

As explained very well in numerous papers of Bailer-Jones, Ted von Hippel and others, it is a computational method which can provide non-linear parameterized mapping between an input vector (a spectrum for example) and one or more outputs like SpT, LC or T_{eff} , $\log g$ and [M/H]. The method is generally supervised; it means that for the network to give required input-output mapping, it must be trained with the help of representative data patterns. These are stellar spectra for which classification or stellar parameters are well determined. The training procedure is a numerical least square error minimization method. The training proceeds by optimizing the network parameters (weights) to give minimum classification error. Once the network is trained the weights are fixed, the network can be used to produce output SpT, LC or T_{eff} , $\log g$ and [M/H] for an unclassified spectrum.



Neural Network Configuration

Neural Network

As shown in figure 1, the neural network has one input layer containing stellar spectrum. Each of the input nodes connects to every node in the next layer of nodes called the hidden layer. The neural network architecture may contain one or more hidden layers. Each of these connections has a weight w associated with it. A given node in hidden layer forms a weighted sum of its inputs. It then passes this sum through a non-linear sigmoid transfer function to give final output from this node. The outputs from nodes in the hidden layer serve as input to the node in the output layer, which again forms weighted sum of its inputs. The training takes place as follows. The weights are initially set with random values over a small range. When the spectrum is fed into the network, the output would also be random. By comparing this output with the target output we can adjust the weights to give an output that is closer to the target value.

Neural Network

The network is trained iteratively by successive passes of the training data through the network and on each pass the weights are perturbed towards their optimal value. The network training is performed by minimizing the least square error. Since the output from neural network is some non-linear function of all of the network input, it implies that the network output is based upon the appearance of whole spectrum. Depending upon the training data the network will learn which wavelength features are more significant than others in determining the correct spectral parameters and correspondingly would assign appropriate values to the network weights. Here the weights are updated backwards from the output layer through the hidden layer hence the algorithm is called back-propagation method.

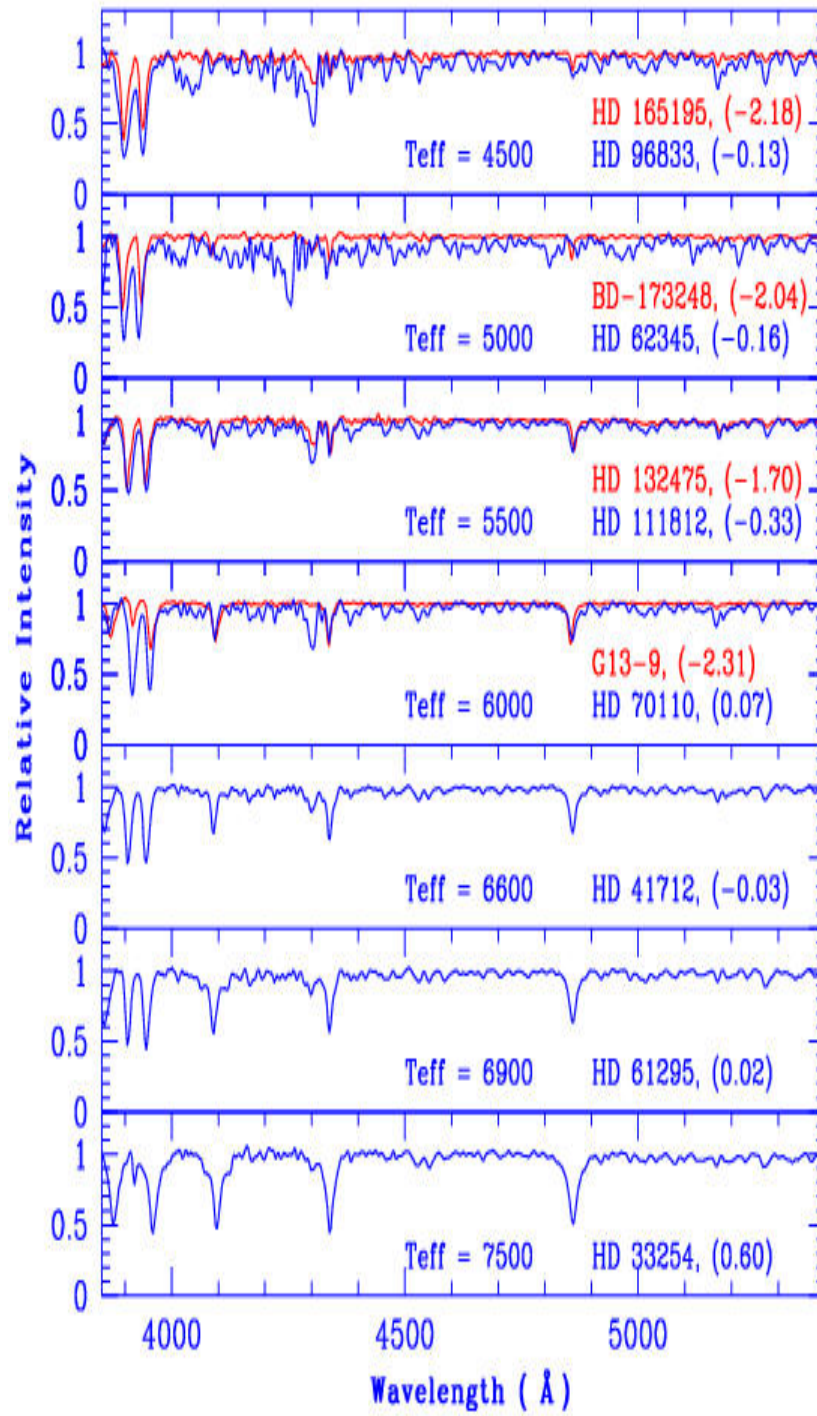
Neural Network : applications

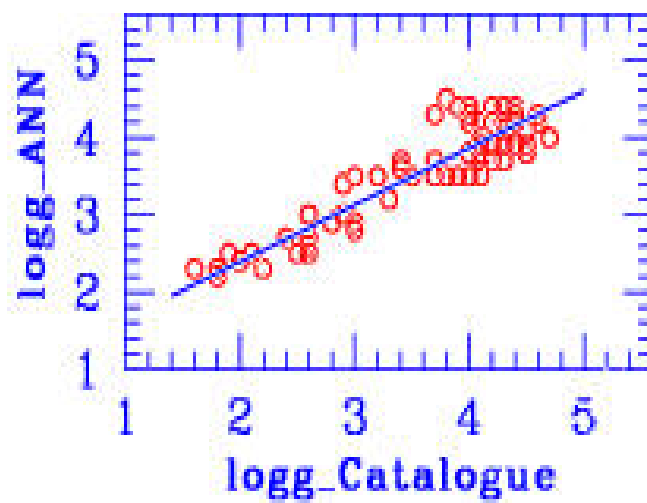
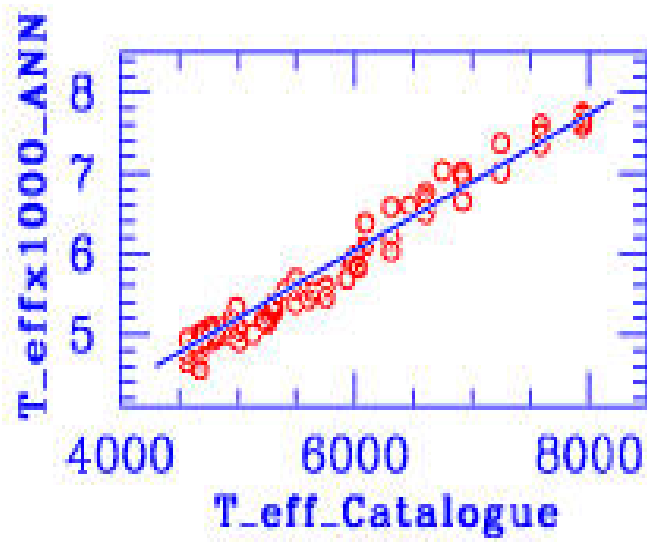
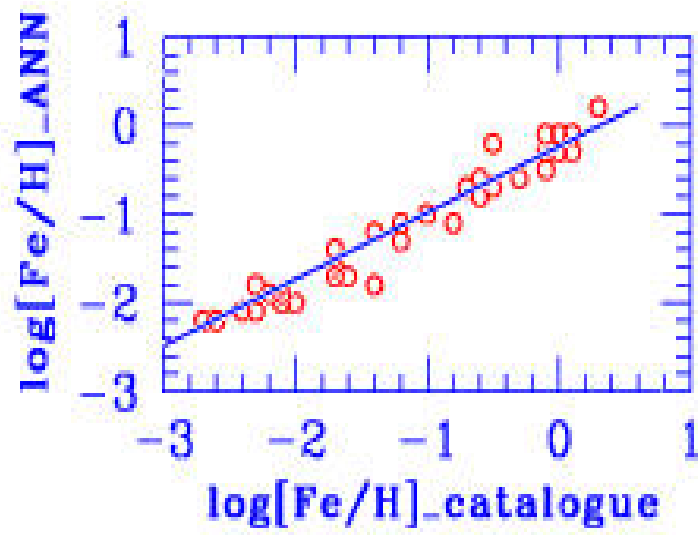
The ANN has been used in very large number of stellar applications. Vieira and Ponz (1995) have used ANN on low-resolution IUE spectra and have determined SpT with an accuracy of 1.1 subclass. Although they attempted classification also with MDM, the errors of classification were larger than that of ANN. In visual region Bailer-Jones, Irwin and von Hippel (1998) used ANN to classify spectra from Michigan Spectral Survey with an accuracy of 1.09 SpT. Visual-near IR spectra were classified by Weaver & Torres-Dodgen(1997) using a two step approach. At first a coarse classification is being done to get main spectral class say F, then it is further classified by more specialist network for that class. This approach results in an accuracy of 0.4 to 0.8 for SpT and 0.2 to 0.4 in LC.

Allende Prieto et al. (2000) used ANN in their search of metal-poor stars. Snijder et al. (2001) used ANN for the three dimensional classification of metal-poor stars. Willemsen et al. (2005) have used ANN to estimate metallicity for main sequence turn-off, subgiants and red giant stars in the globular clusters M55 and ω Centauri using the medium resolution spectra of cluster members.

Parameterization of OMR spectra

We have made a modest effort to use ANN for parameterization of a sample of stars in temperature range 4500 to 8000 K. We have used a medium resolution Cassegrain spectrograph with 2.3m Vainu Bappu Telescope at VBO, Kavalur, India. The spectrograph gives a resolution R (~ 1000) when used with a grating of 600 grooves/mm and a camera of focal length 150mm. The spectral coverage is 3800-6000 Å. The pre-processing of spectra was carried out following a procedure very similar to that of Snijder et al. (2001). We have observed stars from the list of Allende Prieto and Lambert (1999) and Snijder et al. (2001) to develop a library of stars with known temperatures, gravities and metallicities. These spectra were used for training and testing the network. We have also observed stars from the lists of metal-poor candidates to estimate the metallicity for them. Although more than 200 stars were observed, here we report results for 90 stars for which atmospheric parameters T_{eff} and $\log g$ are well determined; $[Fe/H]$ was known for 47 of them. We have used a software developed by B.D. Ripley based on back propagation technique. Figure 2 shows a few representative spectra. The preliminary results based on 680:11:3 architecture are presented in figure 3. The RMS error for $T_{eff} = 200$ K, $[Fe/H] = 0.3$ dex, and that of $\log g = 0.4$ dex. We propose to experiment with an architecture containing two hidden layers instead of one to further reduce the errors of these estimated parameters.





Future Goals

It is very important to envisage an approach that would give quick and reliable spectral classifications (or stellar parameters) for stars falling in all regions of the HR diagram. A single ANN architecture may not give the same desired accuracy over full range of spectral types and luminosity classes. A pilot program using the photometric inputs e.g. Strömgren indices, special photometric indices measuring the strengths of molecular bands for late type stars could serve as preprocessor and help in identifying a set of specialist networks which would lead to classification with desired accuracy. A specialist system also needs to be evolved for A-type star for a quick identification of chemically peculiar, magnetic or emission line stars. An expert system should also give strength of α elements or that of carbon using CH, CN bands.

Special network needs to be developed for objects displaying complex spectra such as Symbiotic stars, Novae and Supernovae. Here the network must be trained on flux calibrated spectra and must use emission line strength as well as shape and structure of the continuum (composite for symbiotic stars and novae) for classification purposes.

Acknowledgements

It is a pleasure to thank Ted Van Hippel for his help in using the software of B.D. Ripley and general encouragement in different stages of this project.