# On the Construction of New Stellar Classification Templates Library for LAMOST Spectra Analysis Pipeline

Peng Wei[1,2], Ali* Luo[1,3], Yinbi Li[1], Jingchang Pan[3], Fengfei Wang[1], Jiannan Zhang[1], Liangping Tu[1,4], Bin Jiang[3], Yongheng Zhao[1], Jianjun Chen[1,2], Xiaoyan Chen[1], Bing Du[1], Wen Hou[1,2], Ge Jin[6], Xiao Kong[1,2], Jie Liu[3], Juanjuan Ren[1,2], Yihan Song[1], Yue Wu[1], Haifeng Yang[1,2,5] and Zhenping Yi[1,2,3]

[1] Key Laboratory of Optical Astronomy,National Astronomical Observatories, Chinese Academy of Sciences, Beijing, 100012, China *lal@nao.cas.cn weipeng@nao.cas.cn*

[2] University of Chinese Academy of Sciences, Beijing, 100049, China

[3] School of Mechanical, Electrical and Information Engineering, Shandong University,Weihai, 264209, China

[4] School of Science, Liaoning University of Science and Technology, Anshan, 144051, China

[5] School of Computer Science and Technology, Taiyuan University of Science and Technology,Taiyuan 030024, China

[6] University of Science and Technology of China, Hefei 230026, China

**Abstract** The LAMOST spectra analysis pipeline is one of LAMOST softwares to produce and analyze the final spectra and its aim is to classify and measure the spectra observed in the survey. Through the pipeline, the observed stellar spectra are classified into different sub-classes by matching with spectra templates. Consequently, the performance of the stellar classification is greatly influenced by the quality of templates. A new LAMOST stellar spectral classification templates library is constructed, which is supposed to improve the precision and credibility of the stellar classification. About one million spectra are selected from LAMOST Data Release one (DR1) to construct the new stellar templates, and they are gathered in 251 groups by two criteria: I) pseudo g-r colors obtained by convolving the spectra with the SDSS *ugriz* filter response curve II) the subclass labeled by the pipeline. In each group, the template spectra are constructed within three steps: I) Outliers are excluded using Local Outlier Probabilities (LoOP) algorithnm, and then the Principal Component Analysis(PCA) method is applied to the remaining spectra of each group. About 5% outliers are ruled out from one million spectra. II) All remaining spectrum are reconstructed using by the first principal components of each group. III) The weighted average spectra are made as the template spectra in the groups. And we initially obtain stellar tempalte spectra in 216 groups. All template spectra are visually inspected, and 52 spectra are abadoned due to low

spectral quality. Furthermore, each template spectrum is manually labeled with a MK class by comparing with three libraries of label-known templates with known MK class. Meanwhile, some unlabeled or wrongly labeled spectra are relabeled or abandoned. And we finally obtain 164 new template spectra with 65 different MK classes. The template library is composed by the spectra left and the first version contains 164 spectra and 65 different MK classes.

**Key words:**  methods: data analysis, methods: statistical, surveys

# 1 INTRODUCTION

The Large Sky Area Multi-Object Fiber Spectroscopic Telescope LAMOST is a special reflecting Schmidt telescope with an effective aperture of 3.6-4.9m, a focal length of 20 m and a field of view (FOV) of $5°$ (Cui et al. 2012). Its unique design enables it to take 4000 spectra in a single exposure. Consequently, the LAMOST has a great potential to efficiently survey a large volume of space for stars and galaxies.

The LAMOST data are processed by data processing softwares written specifically for the LAMOST Spectral Survey. The LAMOST Spectra Analysis Pipeline (also called 1D pipeline) (Luo & Zhao 2001; Luo et al. 2004, 2008; Wang et al. 2010; Luo et al. 2012) is one of these softwares to produce and analyze final spectra. The pipeline performs out $\chi^2$ fits of the spectra to templates in wavelength space, fitting spectra with linear combinations of eigen-spectra and low-order polynomials. Through the pipeline, the observed stellar spectra are classified into different sub-classes. Consequently, the performance of the stellar classification greatly depends on the quality of templates. The current library used for stellar classification in LAMOST contains 36 classes plus 20 subclasses specially for A-type star. The first 36 template are constructed from a set of SDSS spectra (Wang et al. 2010). Other 20 A-type spectra in MILES library (Falcón-Barroso et al. 2011) are picked out to add into the templates. Although there the LAMOST stellar spectra are similiar with SDSS stellar spectra and MILES spectra, there are some differences which can not be ignored. Firstly, the LAMOST and SDSS spectral resolutions are 1800 and 2000 correspondingly. Secondly, different instrumental designs bring about differnt effects on the spectra. In addition, the processes of spectrum extraction, wavelength calribration and flux calribration (Bai 2012) are also different. Considering these issues, it is very necessary to construct a new template library based on the spectra observed and processed by LAMOST.
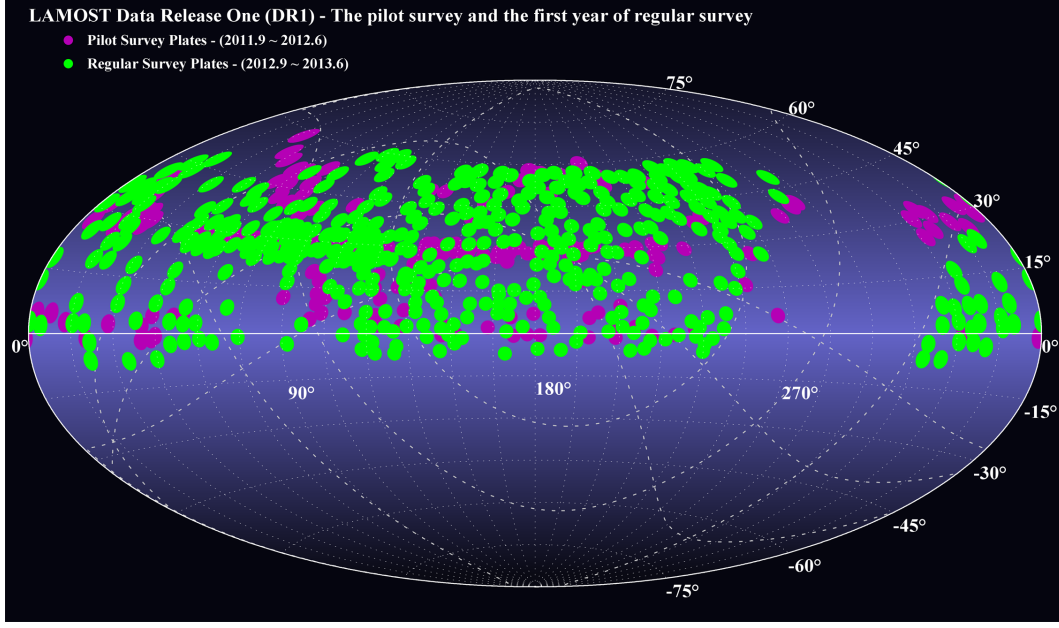
In this paper, we described in detail the construction of the new LAMOST stellar classification template library. The paper is organized as follows: 2 detailedly describe the construction process of the LAMOST stellar template library. The results and discussions are given in sedction 3. A brief summary is given in section 4.

# 2 THE CONSTRUCTION OF THE TEMPLATES LIBRARY

## 2.1 The Spectra From LAMOST Data Release One (DR1)

The first data release (DR1)[1] of LAMOST survey contains the spectra in the pilot survey and the first year of general survey. The pilot survey of LAMOST was launched on Oct 2011 , and ended in June 2012.The first

---

[1] http://data.lamost.org/dr1/

**Fig. 1**   The   LAMOST   DR1   skycoverage   (http://data.lamost.org/u/img/
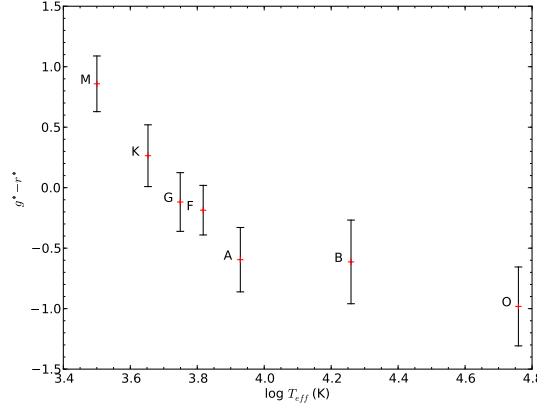dr1-full.png

year of LAMOST regular survey began from September 2012 and ended on June 2013. The DR1 totally contains 2,204,860 spectra, including 717,660 spectra of pilot survey and 1,487,200 spectra of regular survey. The sky coverage of LAMOST DR1 is shown in Fig 1. In addition, the atmospheric parameters of 1,085,404 stars are calculated, which becomes the largest stellar spectral parameters catalog in the world at present.

We exclude those spectra in the Galactic Anti Center and M31 to avoid the high effect of interstellar dust extinction and then about 742,669 spectra are left. The spectral resolution R is about 1800 around g band with a 2/3 silt width (Wang et al. 2013) and the wavelength coverage is from 3700 Å to 9100 Å. To extract spectra from raw observation data, the raw data have been reduced with LAMOST 2D pipeline (Bai 2012) including bias subtraction, cosmic-ray removal, spectral trace and extraction, flat-fielding, wavelength calibration sky subtraction, and combination. Then the 1D pipeline (Wang et al. 2010; Luo et al. 2012) gives spectral type and redshift (radial velocity for stellar spectra).

## 2.2  Group Dividing

### 2.2.1  The pseudo g-r color

LAMOST is a spectroscopic survey oriented telescope and doesn't have its own photometric data. The photometric data of the objects are from different catalogs of other surveys. Meanwhile, the flux calibration is relative not absolute (Bai 2012). Consequently, we can not get accurate and uniform colors for LAMOST spectra. To overcome this problem, we propose a pseudo g-r color(hereafter $g^*$-$r^*$) obtained by convolving each observed spectra with the SDSS $ugriz$ filter response curves . We describe the calculation in detail as follows:

**Fig. 2** The average value and standard deviation of $g^* - r^*$ for each class. For each class, the X-value are the median effective temperature in theory. Meanwhile, the center of each error-bar is the average $g^* - r^*$ color of spectra classified as the corresponding class and the half length is the standard deviation.

1. Suppose that the sampling points of SDSS $g$ & $r$ filter response curves are $P_g$, $P_r$ respectively and the response curve values are $C_g$, $C_r$ respectively .

2. Interpolate the flux of the observed spectra in the points of $P_g$ and $P_r$ to get $F_g$ and $F_r$ respectively.

3. Get the pseudo color $g^*$-$r^*$:

$$g^* - r^* = -2.5 * log\frac{F_g \bigotimes C_g}{\sum C_g} + 2.5 * log\frac{F_r \bigotimes C_r}{\sum C_r} \tag{1}$$

The g-r color is a very good indicator of stellar surface effective temperature (Teff) (Lee et al. 2008; Željko Ivezić et al. 2008). To check whether there are some relationships between the $g^* - r^*$ and the Teff, we select these objects with SDSS $ugriz$ filter magnitude (from different surveys) and sinal to noise ration (SNR) > 20. For these spectra, the average value and standard deviation of $g^* - r^*$ for different spectral types (O, B, A, F, G, K, M-type) are also calculated. As shown in Figure 2, the $g^* - r^*$ color varies obviously in each class.

For these objects, the diagram of g-r color and $g^* - r^*$ is shown in Fig.3. There is a obvious linear relationship between these two colors. We derive the best-fit expression as:
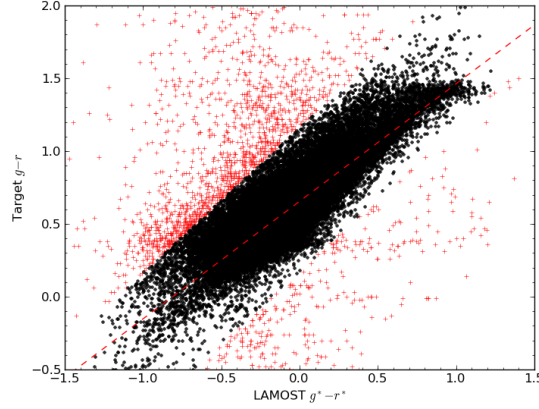
$$g - r = 0.807 * (g^* - r^*) + 0.655 \tag{2}$$

For SDSS spectra, Željko Ivezić et al. (2008) derived a relation between the Teff and the color g-r in the $-0.3 < g - r < 1.3$ color range:

$$Log_{10}(T_{eff}/K) = 0.0283 * (g - r)^3 + 0.0488 * (g - r)^2 - 0.316 * (g - r) + 3.882 \tag{3}$$
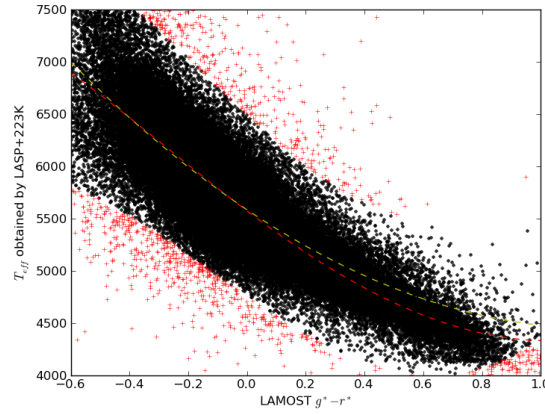
We can then derive a expression between effective temperature $T_{eff}$ and the color g*-r* from the formula 3 and the formula 2 as:

$$Log_{10}(T_{eff}/K) = 0.0283 * (g^* - r^*)^3 + 0.0318 * (g^* - r^*)^2 - 0.203 * (g^* - r^*) + 3.696 \tag{4}$$

For some spectra classified as A, F, G, K-type, the effective temperatures , surface gravities and metallicities determined by the LASP (LAMOST Stellar Parameter pipeline, see Wu et al. (2011)) are provided.

**Fig. 3** The diagram of g-r color and $g^* - r^*$. The X-axis is our proposed $g^* - r^*$ and the Y-axis is the g-r color obtained from target catalogue. The red line is our derived best-fit expression as formula 3. And the points in red are excluded outliers while deriving the expression.
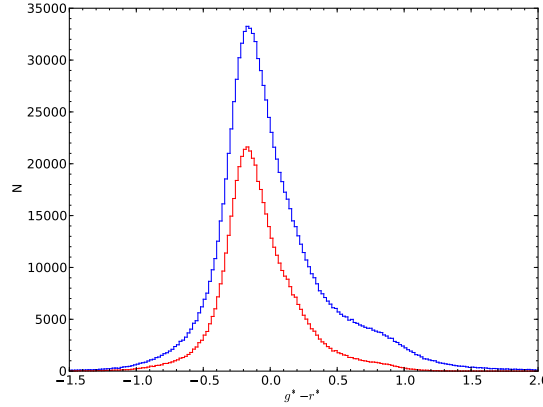


**Fig. 4** The relationship between $g^* - r^*$ color and $T_{eff}$. The $T_{eff}$ is added by 223K to decrease the system inconsistency between SSPP and LASP (Wu et al. 2011). The yellow line is the expression as formula 4. The red line is our derived best-fit expression as formula 5. And the points in red are excluded outliers while deriving the expression.

The relationship between $g^* - r^*$ color and $T_{eff}$ is shown in Figure 4. We also derive a best fit 3-order polynomial expression as:

$$Log_{10}(T_{eff}/K) = 0.0432 * (g^* - r^*)^3 + 0.0107 * (g^* - r^*)^2 - 0.165 * (g^* - r^*) + 3.746 \quad (5)$$

As shown in Fig.4, these two formulas 4 and 5 nearly coincide with each other in the Teff range [5500K,7000K]. Consequently, we can infer that our defined $g^* - r*$ color can also be a good indicator of Teff. And then the $g^* - r*$ color is applied into the group divding of the slected spectra.

**Fig. 5** The number distribution of spectra in each $g^*$-$r^*$ bin. The blue line is the distribution of all spectra while the red line is the distribution of the spectra with $SNR > 10$.

### 2.2.2 Group dividing criteria

To construct different kinds of templates, we gather these spectra in 251 different groups by the proposed $g^* - r^*$ color and the subclass labeled by the pipeline.

As discussed above, the proposed $g^* - r^*$ color is an good indicator of Teff. Consequently, we select these spectra with $g^*$-$r^*$ in the range[-1.5,2.0] and divide all spectra into 175 groups with 0.02 mag width interval. The number distribution is shown in Figure 5. These groups are marked with group-id from 1 to 175.

Meanwhile, other 75 groups are formed by the subclass labeled by current pipeline. After the automated processing of LAMOST Spectra Analysis Pipeline and visual inspection, there are 76 different stellar subclasses in our selected spectra. Some A-type spectra in MILES library (Falcón-Barroso et al. 2011) are picked out to add into the templates. Consequently, there are more A-type subclasses. These groups are marked with group-id from 176 to 251. Their distribution is as shown in Table 1. Yi et al. (2013) presented a spectroscopic catalog of 67,082 M dwarfs from LAMOST Pilot Survey. All spectra are divided into 10 subclasses from M0 to M9. We have divided these spectra into groups with group id as shown in Table 1. Zhao et al. (2013) presented a spectroscopically identified catalog of 70 DA white dwarfs (WDs). Meanwhile, Zhang et al. (2013) identified 230 other DA white dwarfs. We combine these two catalogs and add the spectra into group 250. Jiang et al. (2013) reported the identification of 10 cataclysmic variables. We add these 10 spectra into group 207.

## 2.3 The Construction of template spectra

### 2.3.1 LOcal Outlier Probabilities (LoOP)

Kriegel et al. (2009) proposed a local outlier factor (LOF, see Breunig et al. (2000)) based outlier detection method. LoOP is a local density based method that uses statistical concepts to output the final score. The LoOP score represents the probability that a particular point is a local density outlier. The LoOP is calculated as follows (Kriegel et al. 2009):

**Table 1** The number distribution of different subclasses

| Group ID | Subclass | Amount | Group ID | Subclass | Amount | Group ID | Subclass | Amount |
|---|---|---|---|---|---|---|---|---|
| 176 | A0 | 67 | 202 | B9 | 443 | 227 | - | - |
| 177 | A0I | 27 | 203 | Binary | 170 | 228 | - | - |
| 178 | A0III | 407 | 204 | Carbon | 168 | 229 | - | - |
| 179 | A0p | 27 | 205 | CarbonWD | 6 | 230 | - | - |
| 180 | A1IV | 653 | 206 | Carbon_lines | 10 | 231 | - | - |
| 181 | A1V | 527 | 207 | CV | 17 | 232 | - | - |
| 182 | A2I | 10 | 208 | EM | 42 | 233 | - | - |
| 183 | A2IV | 1692 | 209 | Emission | 21 | 234 | M0 | 19139 |
| 184 | A2V | 5761 | 210 | F0 | 27808 | 235 | M0V | 13 |
| 185 | A3I | 37 | 211 | F2 | 44192 | 236 | M1 | 19953 |
| 186 | A3IV | 2084 | 212 | F5 | 119328 | 237 | M2 | 17231 |
| 187 | A3V | 2080 | 213 | F9 | 292830 | 238 | M2V | 12 |
| 188 | A4III | 926 | 214 | G0 | 47697 | 239 | M3 | 9749 |
| 189 | A4V | 773 | 215 | G2 | 92229 | 240 | M4 | 3860 |
| 190 | A5 | 26 | 216 | - | - | 241 | M5 | 855 |
| 191 | A5I | 213 | 217 | G5 | 81202 | 242 | M6 | 412 |
| 192 | A5V | 1253 | 218 | G7 | 3650 | 243 | M7 | 259 |
| 193 | A6IV | 1240 | 219 | - | - | 244 | M8 | 47 |
| 194 | A6V | 322 | 220 | K0 | 1998 | 245 | M9 | 66 |
| 195 | A7III | 4033 | 221 | K1 | 85218 | 246 | Non | 2 |
| 196 | A7V | 647 | 222 | K3 | 77164 | 247 | O | 79 |
| 197 | A9 | 4 | 223 | K5 | 73045 | 248 | OB | 16 |
| 198 | A9V | 2170 | 224 | - | - | 249 | T2 | 11 |
| 199 | B | 15 | 225 | K7 | 45839 | 250 | WD | 535 |
| 200 | B0 | 1 | 226 | K9 | 4 | 251 | WDmagnetic | 14 |
| 201 | B6 | 492 | | | | | | |

Notes: '-' in subclass and amount means the corresponding group is neglected.

1. ($k$-$distance$ of an object $p$) For any positive integer $k$, the $k$-distance (The distance measurement function used in our work is the cosine distance $d = 1 - \frac{A*B}{|A|*|B|}$) of object $p$, denoted as $k$-$distance(p)$, is defined as the distance $d(p, o)$ between $p$ and an object $o \in D$ such that: (i) For at least $k$ objects $o' \in D \setminus p$, it holds that $d(p, o') \leq d(p, o)$. (ii )For at most $k$-1 objects $o' \in D \setminus p$, it holds that $d(p, o') < d(p, o)$.

2. ($k$-$distance$ neighborhood of an object $p$) Given the $k$-$distance$ of $p$, the $k$-$distance$ neighborhood of $p$ contains every object whose distance from $p$ is not greater than the $k$-$distance$, i.e. $N_{k-distance(p)}(p) = \{q \in D \setminus p \mid d(p, q) \leq k\text{-}distance(p)\}$. These objects $q$ are called the $k$-$nearest$ neighbors of $p$. Simplify the notation to use $N_k(p)$ as a shorthand for $N_{k-distance}(p)$.

3. The standard distance. $\sigma(p, N_k(p))$ is defined as the standard deviation of the distance around $p$:

$$\sigma(p, N_k(p)) = \sqrt{\frac{\sum_{s \in N_k(p)} d(p, s)^2}{|S|}} \quad (6)$$

4. The probabilistic set distance. $pdist(\lambda, p, N_k p)$ is defined as follows:

$$pdist(\lambda, p, N_k p) = \lambda * \sigma(p, N_k(p)) \quad (7)$$

5. The Probabilistic Local Outlier Factor PLOF. $PLOF_{\lambda,S(p)}(p)$ represents the ratio of the density estimation:

$$PLOF_{\lambda,S(p)}(p) = \frac{pdist(\lambda, p, N_k(p)) * |N_k(p)|}{\sum_{s \in N_k(p)} pdist(\lambda, s, N_k(s))} - 1 \tag{8}$$

6. The aggregate Probabilistic Local Outlier Factor nPLOF. This is the scaling factor that makes the score independent from any distribution:

$$nPLOF = \lambda * \sqrt{\frac{\sum_{p \in D} PLOF_{\lambda,S(p)}(p)^2}{|D|}} \tag{9}$$

7. The Local Outlier Probability:

$$LoOP_{N_k p}(p) = max(0, erf \frac{PLOF_{\lambda,S}(p)}{nPLOF * \sqrt{2}}) \tag{10}$$

Although the spectra are divided into different groups by g$^*$-r$^*$ or by the classification, there are also some outliers in each group. The LoOP method is used to exclude these outliers in each group.

### 2.3.2 Principal Component Analysis (PCA)

The Principal Component Analysis (PCA) (Jolliffe 2002) is a mathematical procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables into a set of values of linearly uncorrelated variables called principal components. The number of principal components is less than or equal to the number of original variables. This transformation is defined in such a way that the first principal component has the largest possible variance (that is, accounts for as much of the variability in the data as possible), and each succeeding component in turn has the highest variance possible under the constraint that it be orthogonal to (i.e., uncorrelated with) the preceding components. The detailed description of the PCA method is discussed in Whitney (1983); Jolliffe (2002).

As a viable tool, Principal Component Analysis (PCA) has been applied in the classification of spectra (Whitney 1983; Bailer-Jones et al. 1998; Yip et al. 2004; Almeida & Prieto 2013). In addition, Tu et al. (2009, 2010); Jiang et al. (2013); Wei et al. (2013) have used PCA to do the data dimension reduction in finding relatively rare objects. McGurk et al. (2010) applied PCA to 200,000 stellar spectra obtained by the Sloan Digital Sky Survey (SDSS). They discussed correlations of eigen-coefficients with metallicity and gravity estimated by the Sloan Extension for Galactic Understanding and Exploration Stellar Parameters Pipeline (SSPP) (Lee et al. 2008).

### 2.3.3 The construction steps

The construction steps is carried out as follows:

1. For the groups with more than 5,000 spectra, only first 5,000 spectra with the largest SNR are selected.
2. De-redshift the spectra and unify the wavelength to 3800Å-9000Å with fixed step 1Å(the amount of all sampling points is N=5201) and then get the unified flux $F$ for each spectrum.
3. Exclude these spectra existing $F \leq 0$ and normalize the remaining spectra $F$ with:

$$F_i = \frac{F_i}{\sqrt{\sum_{j=1}^{N} F_i^2}} \tag{11}$$

4. Calculate LoOP in each group

5. These spectra with $LoOP \geq 0.4$ are excluded.

6. Do PCA of the remaining spectra in each group to get a feature matrix $T$ and the respective eigen values $\lambda$.

7. Select the first $k$-th principal components (eigen spectra) while the variance contribution rate $\mu$ :

$$\mu = \frac{\sum\limits_{i=1}^{k} \lambda_i}{\sum\limits_{i=1}^{N} \lambda_i} > \theta \tag{12}$$

where $\theta$ is a fixed given threshold (0.99 is used in our work). $k$ is set to 2 when $k = 1$.

8. Reconstruct each remaining spectra using obtained first $k$ principal components

9. Calculate LoOP of remaining reconstructed spectra in each group again and exclude these spectra with $LoOP \geq 0.2$

10. Get the SNR weighted average spectrum as the template spectrum .

Following the above steps, the template spectra are successfully constructed in 216 groups (nearly 86%). Other 35 groups fail mainly because of lacking enough high quality spectra.

### 2.3.4 MK Class Labeling

Each spectrum should be labeled a subclass for latter usage in classification. We compare these spectra with three libraries and then label each spectrum a subclass.

Danks & Dennefeld (1994) presented spectra for MK standards in the wavelength range 5800Å-10200Å. The stars cover the normal spectral types from O to M and luminosity types I, III, and V. The projected slit width along the dispersion is about 4Å and the resolution R is about 1200. Two wavelength ranges [7500Å,7700Å] and [6800Å,7000Å] are masked to get rid of the strong telluric lines left in the spectra. We decrease the resolution of our templates to R 1200 by convolving a gaussian function. All template spectra and standard spectra are unified into the wavelength range [6100Å,9000Å] with a fixed step 4Å. For each template spectrum, the first four closest are chosen and the corresponding spectra are drawn together with the template spectra. Those figures are used for following visual inspection.

Bolton et al. (2012) described the detail of the pipeline for SDSS III and published the template used on the web page[2]. For stellar spectral classification, 123 templates created from the full database of Indo-U.S. spectra are provided. Each spectrum are labeled a MK class by matching with POLLUX database. The resolution R of these 123 spectra is about 2000 and the wavelength coverage is from 3500Å to 11200Å. These spectra are unified into the wavelength range [3800Å,9000Å] with a fixed step 1Å similar with the spectra in the library . Similarly, for each template spectrum, the first four closest are chosen and the corresponding spectra are drawn together with the template spectra. And those figures are also used for following visual inspection.
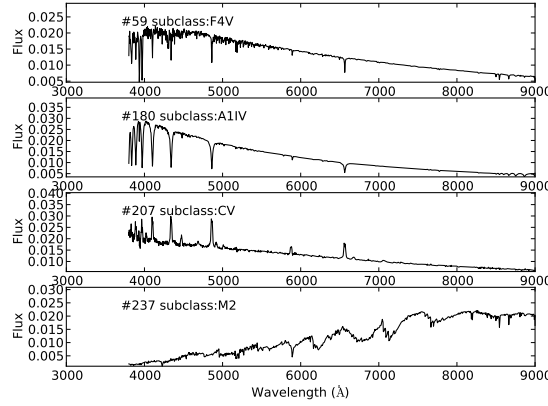
As introduced above, the current library used for stellar classification in LAMOST contains 36 classes plus 20 subclasses specially for A-type star. The resolution R of these 56 spectra is about 2000 and the wavelength coverage is from 3800Å to 9200Å. These spectra are unified into the wavelength range

---

[2] http://www.sdss3.org/svn/repo/idlspec2d/tags/v5_4_45/templates/

**Table 2** The main information of groups 59,180, 207 and 237

| Group ID | All spectra | used spectra | Subclass | Subclass1 | Subclass2 | Subclass3 |
|---|---|---|---|---|---|---|
| 59 | 18534 | 3048 | F4V | F3V/F5V | F8III-IV | F5 |
| 180 | 653 | 381 | A1IV | A4V/A1V | A2V | A1V |
| 207 | 27 | 13 | CV | - | - | - |
| 237 | 17231 | 325 | M2 | M1/M0 | M1.5V/M3V | M2/M1 |

Notes: Subclass is the finally labeled MK class. Subclass1 is the best fit Mk class with Bolton et al. (2012). Subclass2 is the best fit Mk class with Danks & Dennefeld (1994). Subclass3 is the best fit Mk class with Luo et al. (2013).



**Fig. 6** The template spectra of groups 59,180, 207 and 237.

[3800Å,9000Å] with a fixed step 1Å similar with the spectra in the library . Similarly, for each template spectrum, the first four closest are chosen and the corresponding spectra are drawn together with the template spectra. And those figures are also used for following visual inspection.

Each template spectrum is visually inspected by checking the three figures drawn above. And then each spectrum is labeled a MK class. Meanwhile, those template spectra with bad data or low S/N R are excluded. Finally, there are 164 spectra and 65 different MK classes are left in the template library.
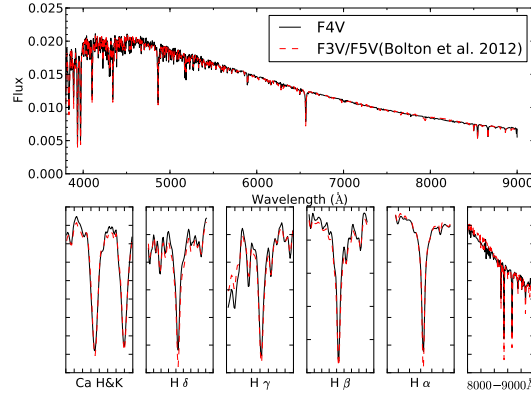
## 3 RESULTS AND DISCUSSIONS

A updated library is formed after adding some spectra of new subclasses and replacing some spectra. These spectra of the types not existing in our library are left. The library of the current version (V1.0) is publicly available on the web site[3]. The current library has been used in the new new version of LAMOST 1D pipeline for spectra after data release one.
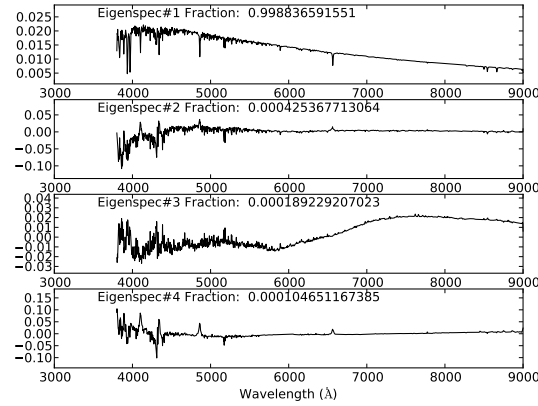
### 3.1 Examples

Here we choose four typical groups (Group 59,180 ,207 and 237) to discuss in detail. The main information of these groups is shown in Table.2. The MK classes are F6V, A1IV, CV and M2 respectively. The finally constructed template spectra of these groups are shown in Fig.6.

---

[3] http://sciwiki.lamost.org/lamost_sctl/v1

**Fig. 7** The comparison of the template spectrum in group 59 with F3V/F5V in Bolton et al. (2012). The black line is the spectrum constructed in our work. The red one is the closest spectrum in Bolton et al. (2012).
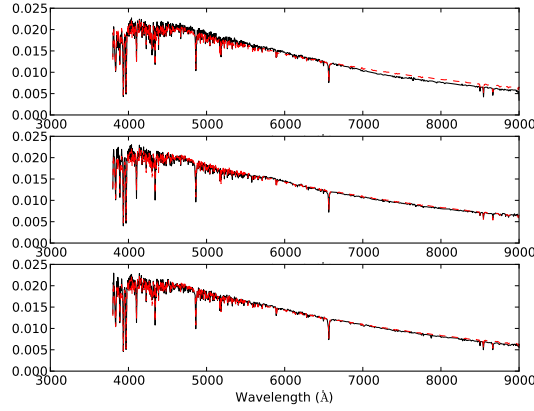


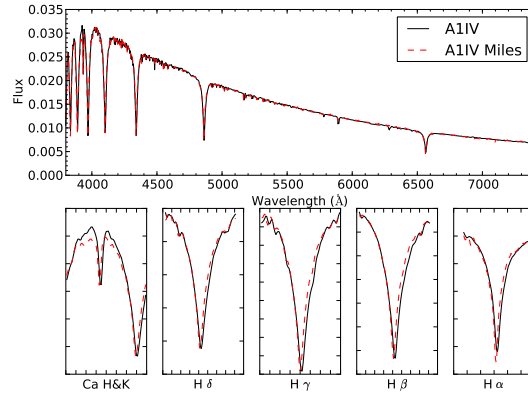**Fig. 8** The first four eigen spectra (principal components) of group 59.

**Group:59** This group contains the spectra in the color $g^*$-$r^*$ range [-0.34,-0.32]. There are totally 18,534 spectra and 3,048 spectra are selected from the first 5,000 spectra with the highest SNR. As shown in Fig 7, the spectrum is very close to F3V/F5V in Bolton et al. (2012). Consequently, the template spectrum is labeled as 'F4V'.

As shown in Fig 8, the variance of the first principal component exceeds more than 99% of the total variance of the original data. That is due to the high similarity of the spectra in the group. Consequently, the reconstructed spectra using first two principal components are nearly similar to the origin spectra (see Fig 9).

**Group:180** This group contains the spectra classified as 'A1IV' by pipeline. There are totally 653 spectra and 381 spectra are selected. The spectrum is labeled as 'A1IV' following the group selection criteria. As shown in Fig 10, the SNR of the template is a little larger than the template in Luo et al. (2013) while these two spectrum are nearly close to each other.

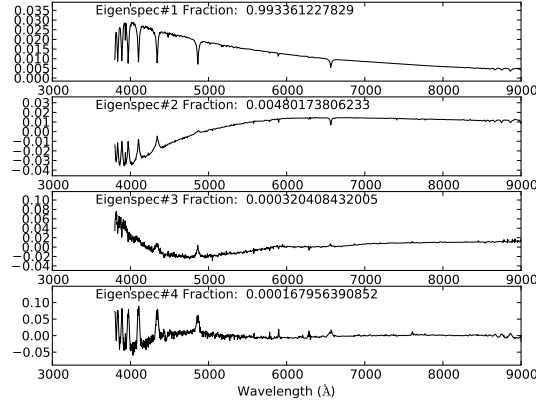**Fig. 9** Three examples of reconstructed spectra in group 59



**Fig. 10** The comparison of the template spectrum in group 59 with A1IV in Luo et al. (2013). The black line is the spectrum constructed in our work. The red one is the closest spectrum in Luo et al. (2013).

Similar with group 59, the variance of the first principal component also exceeds more than 99% of the total variance of the original data (see Fig 11). However, there are not as many spectra as in group 59. Consequently, some spectra are not well reconstructed (as shown in Fig 12). In spite of this, the template spectrum is well constructed after excluding these badly reconstructed spectra.
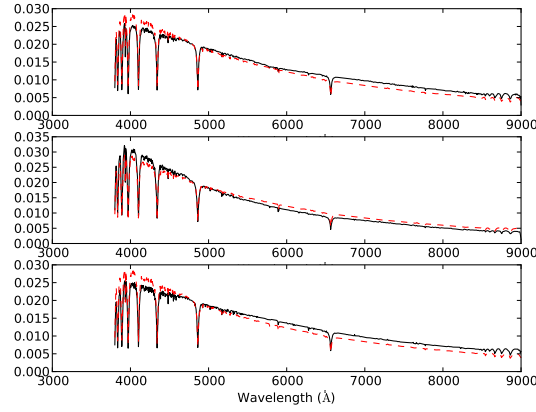
**Group:207** This group contains the spectra classified as 'CV'. There are totally 27 spectra and 13 spectra are selected. The spectrum is labeled as 'CV' following the group selection criteria.

Compared with normal stars, the spectra of CV stars are these with strong hydrogen Balmer and helium emission lines that typically signify ongoing accretion. As shown in Fig 13, the first two principal components show obvious and strong emission lines and the sum of the variances of these two principal components exceeds more than 99% of the total variance of the original data. Compared to normal stars misclassified as 'CV', the spectra of CV stars are almost faultlessly reconstructed (see Fig 14). And then these misclassified spectra are excluded in the next following steps.
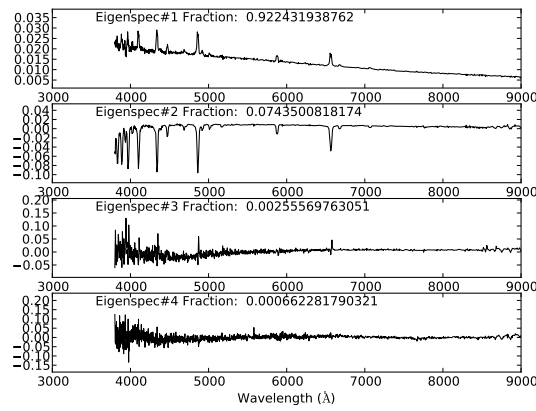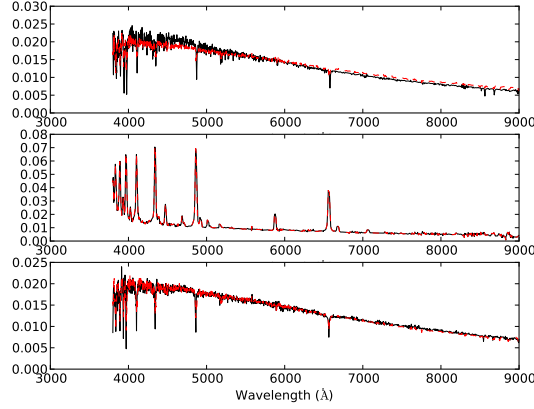
**Group:237**

**Fig. 11** The first four eigen spectra (principal components) of group 180.
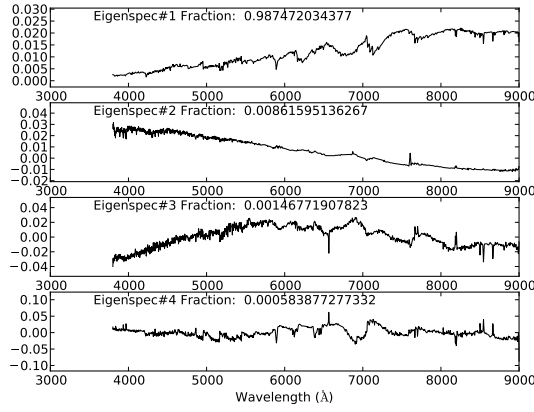


**Fig. 12** Three examples of reconstructed spectra in group 180



**Fig. 13** The first four eigen spectra (principal components) of group 207. Note that the strong lines in eigen spectra 2 are emission lines not absorption lines.

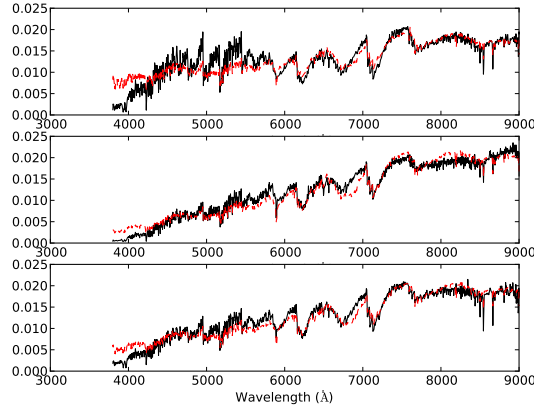**Fig. 14** Three examples of reconstructed spectra in group 207.



**Fig. 15** The first four eigen spectra (principal components) of group 237.
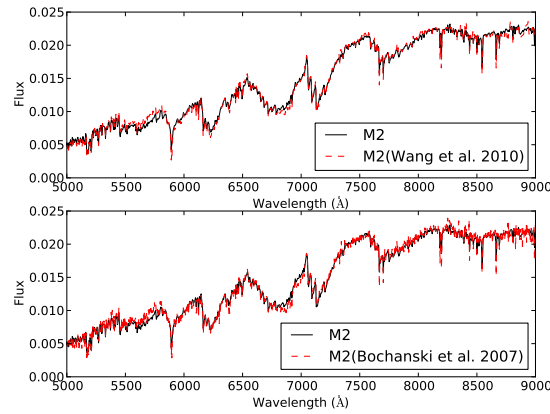
This group contains the spectra classified as 'M2'. There are totally 17,231 spectra and 325 spectra are selected from the first 5,000 spectra with the highest SNR. Due to the existence of wavelength points with $flux \leq 0$, a large amount of spectra are excluded in this group. The spectrum is labeled as 'M2' following the group selection criteria.

As shown in Fig 15, the sum of the variances of the first two principal components exceeds more than 99% of the total variance of the original data. The selected spectra are not well reconstructed in the blue arm (as shown in Fig 16). In spite of this, the template spectrum is also well constructed.

To check the quality, we choose the M2-type template spectrum in the current template library (Wang et al. 2010) and compare it with the template spectrum of group 237 (see Fig 17 upper panel). Bochanski et al. (2007) presented template spectra of low-mass (M0-L0) dwarfs derived from over 4000 Sloan Digital Sky Survey spectra. We choose the M2-type template spectrum and alsoe compare it with the template spectrum of group 237 (see Fig 17 bottom panel). As shown in Fig 17, we can infer that our constructed M2-type spectrum is a little better than these two spectra.

**Fig. 16** Three examples of reconstructed spectra in group 237.



**Fig. 17** The comparison of the template spectrum in group 59 with M2 in Bochanski et al. (2007). The black line is the spectrum constructed in our work. The red one is the closest spectrum in Bochanski et al. (2007).

## 3.2 Discussions

### 3.2.1 Comparison with *McGurk et al. (2010)*

McGurk et al. (2010) applied PCA to about 100,000 SEGUE spectra by dividing all spectra into 55 different bins. For each bin, the first four eigenspectra are published and the first one is a high SNR mean spectra. For the template spectra in all groups, to check the difference, we use similar method as MK class labeling to find three closest mean spectra in McGurk et al. (2010). As our color range is wider than McGurk et al. (2010), not all template spectra are similar with these in McGurk et al. (2010). As shown in Table 3, the groups with groups id from 38 to 99 (totally 60 groups, not including group 71) cover nearly all mean spectra constructed by McGurk et al. (2010). The finally subclasses are from A3 to K3, which coincides with sayings in McGurk et al. (2010).

**Table 3** The comparison with McGurk et al. (2010)

| Group ID | ESID1 | ESID2 | ESID3 | Group ID | ESID1 | ESID2 | ESID3 | Group ID | ESID1 | ESID2 | ESID3 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 38 | 1 | 2 | 3 | 59 | 19 | 20 | 18 | 80 | 39 | 38 | 40 |
| 39 | 1 | 2 | 3 | 60 | 19 | 20 | 21 | 81 | 40 | 41 | 39 |
| 40 | 2 | 3 | 1 | 61 | 20 | 23 | 22 | 82 | 41 | 40 | 42 |
| 41 | 3 | 2 | 4 | 62 | 24 | 23 | 25 | 83 | 42 | 41 | 43 |
| 42 | 4 | 5 | 3 | 63 | 25 | 24 | 26 | 84 | 43 | 42 | 44 |
| 43 | 5 | 6 | 7 | 64 | 26 | 25 | 27 | 85 | 44 | 43 | 45 |
| 44 | 6 | 7 | 5 | 65 | 27 | 26 | 28 | 86 | 45 | 44 | 46 |
| 45 | 7 | 8 | 6 | 66 | 28 | 27 | 29 | 87 | 45 | 46 | 44 |
| 46 | 8 | 9 | 7 | 67 | 28 | 29 | 27 | 88 | 46 | 47 | 45 |
| 47 | 9 | 10 | 11 | 68 | 29 | 30 | 28 | 89 | 47 | 48 | 46 |
| 48 | 11 | 10 | 12 | 69 | 30 | 29 | 31 | 90 | 48 | 47 | 49 |
| 49 | 12 | 11 | 10 | 70 | 30 | 31 | 29 | 91 | 49 | 48 | 50 |
| 50 | 12 | 13 | 11 | 71 | - | - | - | 92 | 50 | 49 | 51 |
| 51 | 13 | 14 | 12 | 72 | 32 | 33 | 31 | 93 | 51 | 50 | 52 |
| 52 | 14 | 15 | 13 | 73 | 33 | 32 | 34 | 94 | 52 | 51 | 53 |
| 53 | 15 | 14 | 16 | 74 | 34 | 35 | 33 | 95 | 53 | 52 | 54 |
| 54 | 16 | 15 | 17 | 75 | 35 | 34 | 36 | 96 | 53 | 54 | 52 |
| 55 | 16 | 17 | 15 | 76 | 35 | 36 | 37 | 97 | 54 | 55 | 53 |
| 56 | 17 | 18 | 16 | 77 | 36 | 37 | 35 | 98 | 55 | 54 | 53 |
| 57 | 18 | 17 | 19 | 78 | 37 | 38 | 36 | 99 | 55 | 54 | 53 |
| 58 | 18 | 19 | 20 | 79 | 38 | 37 | 39 | | | | |

Notes: The ESID1, ESID2 and ESID3 are bin ID of the first three closest eigen spectra in McGurk et al. (2010) respectively.

### 3.2.2 Comparison with current templates of LAMOST spectra analysis pipeline and Bolton et al. (2012)

As supplements to current templates, our constructed templates replace most spectra in the current library including the A-type stellar spectra. From the comparisons discussed in section 3.1, we can infer that newly constructed template spectra are a little better than current ones.

There are 123 stellar subclasses in Bolton et al. (2012). These template spectra are individual spectra in Indo-U.S. database. We notice that there are 80 subclasses which contain less than 500 spectra in SDSS DR9. There are totally 12,897 spectra (about 1.66% in all 773,275 spectra) and 1,062 spectra (about 0.22% in all 475694 spectra) with SNR$> 10$. Meanwhile, the average SNR of these spectra is about 4.93 which is much less the one of all spectra. In other words, there are mainly 43 stellar subclasses containing most spectra especially these spectra with high SNR.

Considering the differences between LAMOST spectra with spectra in other survey, these newly constructed spectra are more reliable and more similar to the spectra observed in LAMOST survey. That is because these template spectra are constructed from a healthy sum of spectra from LAMOST DR1.

### 3.2.3 Remaining problems

The result shows that our constructed template spectra can be used in the classification of observed stellar spectra in LAMOST survey. However, there are also some problems needing to solve.

1. We notice that most of our template spectra are main sequence stars. To construct the template spectra of some other rare types such as K-type giants, DC and DZ white dwarfs, some spectra need to be labeled manually or picked out by other methods.

2. We use three libraries to label each template. However, how to label the template spectra better is also a remaining problem.

3. In addition, there are some outliers excluded in each group while constructing the templates. It is also worth of studying these objects and finding rare types even new types of star.

## 4 SUMMARY

To improve the precision and credibility of the stellar classification, A new LAMOST stellar spectral classification templates library is constructed. We select about 750,0000 stellar spectra from LAMOST Data Release One (DR1) and then we gather them in 251 different groups by proposed pseudo g-r colors and the subclass labeled by the pipeline. Following the proposed contruction steps, including excluding outliers using LoOP, spectral PCA reconstruction etc., the weighted average spectra are constructed as the template spectra in the groups. Afterwards, each template spectrum is labeled with a MK class by comparing with three libraris and visual inspection. Some low-quality spectra are excluded afetr visual inspection . Meanwhile, some unlabeled or wrongly labeled spectra are relabeled or abandoned. The template library is composed by the spectra left and the first version contains 164 spectra and 65 different MK classes. The new templates library has been used in new version of LAMOST Spectra Analysis Pipeline and is published on the website [4].

## References

Almeida, J. S., & Prieto, C. A. 2013, ApJ, 763, 50 8

Bai, Z. 2012, Proceedings of the International Astronomical Union, 8, 189 2, 3

Bailer-Jones, C. A., Irwin, M., & Hippel, T. V. 1998, MNRAS, 298, 361 8

Bochanski, J. J., West, A. A., Hawley, S. L., & Covey, K. R. 2007, AJ, 133, 531 14, 15

Bolton, A. S., Schlegel, D. J., Aubourg, É., et al. 2012, AJ, 144, 144 9, 10, 11, 16

Breunig, M. M., Kriegel, H.-P., Ng, R. T., & Sander, J. 2000, in ACM Sigmod Record, vol. 29, 93–104 (ACM) 6

Cui, X.-Q., Zhao, Y.-H., Chu, Y.-Q., et al. 2012, RAA, 12, 1197 2

---

[4] http://sciwiki.lamost.org/lamost_sctl/v1

Danks, A. C., & Dennefeld, M. 1994, PASP, 382–396 9, 10

Falcón-Barroso, J., Sánchez-Blázquez, P., Vazdekis, A., et al. 2011, A&A, 532, A95 2, 6

Jiang, B., Luo, A., Zhao, Y., & Wei, P. 2013, MNRAS, 430, 986 6, 8

Jolliffe, I. T. 2002, Principal component analysis (Springer verlag) 8

Kriegel, H.-P., Kröger, P., Schubert, E., & Zimek, A. 2009, in Proceedings of the 18th ACM conference on
    Information and knowledge management, 1649–1652 (ACM) 6

Lee, Y. S., Beers, T. C., Sivarani, T., et al. 2008, AJ, 136, 2022 4, 8

Luo, A.-L., Wu, Y., Zhao, J., & Zhao, G. 2008, in Proc. of SPIE Vol, vol. 7019, 701935–1 2

Luo, A.-L., Zhang, H.-T., Zhao, Y.-H., et al. 2012, RAA, 12, 1243 2, 3

Luo, A.-L., Zhang, Y.-X., & Zhao, Y.-H. 2004, in Astronomical Telescopes and Instrumentation, 756–764
    (International Society for Optics and Photonics) 2

Luo, A.-L., & Zhao, Y.-H. 2001, Chinese Journal of Astronomy and Astrophysics, 1, 563 2

Luo, A.-L., et al. 2013, in preparation 10, 11, 12

McGurk, R. C., Kimball, A. E., & Željko Ivezić 2010, AJ, 139, 1261 8, 15, 16

Tu, L., Luo, A., Wu, F., & Zhao, Y. 2010, Science China Physics, Mechanics and Astronomy, 53, 1928 8

Tu, L.-P., Luo, A.-L., Wu, F.-C., Wu, C., & Zhao, Y.-H. 2009, RAA, 9, 635 8

Željko Ivezić, Sesar, B., Jurić, M., et al. 2008, ApJ, 684, 287 4

Wang, F., Luo, A., & Zhao, Y. 2010, in SPIE Astronomical Telescopes and Instrumentation: Observational
    Frontiers of Astronomy for the New Decade, 774031–774031 (International Society for Optics and
    Photonics) 2, 3, 14

Wang, F., Zhang, H., Luo, A.-L., et al. 2013, arXiv preprint arXiv:1306.1600 3

Wei, P., Luo, A., Li, Y., et al. 2013, MNRAS, 431, 1800 8

Whitney, C. 1983, Astronomy and Astrophysics Supplement Series, 51, 443 8

Wu, Y., Luo, A.-L., Li, H.-N., et al. 2011, RAA, 11, 924 4, 5

Yi, Z., Luo, A., Song, Y., et al. 2013, arXiv preprint arXiv:1306.4540 6

Yip, C., Connolly, A., Berk, D. V., et al. 2004, AJ, 128, 2603 8

Zhang, Y.-Y., Deng, L.-C., Liu, C., et al. 2013, AJ, 146, 34 6

Zhao, J., Luo, A., Oswalt, T., & Zhao, G. 2013, AJ, 145, 169 6