



**HAL**  
open science

# Principal component analysis of palaeomagnetic directions: converting a Maximum Angular Deviation (MAD) into an $\alpha_{95}$ angle

G Khokhlov, G Hulot

► **To cite this version:**

G Khokhlov, G Hulot. Principal component analysis of palaeomagnetic directions: converting a Maximum Angular Deviation (MAD) into an  $\alpha_{95}$  angle . Geophysical Journal International, 2015, 204 (1), pp.274-291. 10.1093/gji/ggv451 . insu-01409093

**HAL Id: insu-01409093**

**<https://hal-insu.archives-ouvertes.fr/insu-01409093>**

Submitted on 5 Dec 2016

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Principal component analysis of palaeomagnetic directions: converting a Maximum Angular Deviation (MAD) into an $\alpha_{95}$ angle

A. Khokhlov<sup>1,2,3</sup> and G. Hulot<sup>3</sup>

<sup>1</sup>IEPT Russian Academy of Sciences, 84/32, Profsoyuznaya, 117997 Moscow, Russia

<sup>2</sup>Institute of Physics of the Earth, Russian Academy of Sciences, 10, B. Gruzinskaya, 123242 Moscow, Russia

<sup>3</sup>Institut de Physique du Globe de Paris, Sorbonne Paris Cité, Université Paris Diderot, CNRS, F-75005 Paris, France. E-mail: [gh@ipgp.fr](mailto:gh@ipgp.fr)

Accepted 2015 October 15. Received 2015 October 12; in original form 2015 May 16

## SUMMARY

Directions recovered from palaeomagnetic samples are usually archived with some quantitative information about their precision, most often in the form of a so-called  $\alpha_{95}$  angle. Such angles are classically co-estimated with the recovered palaeomagnetic direction from a collection of samples providing individual estimates of this direction. In some instances, however, palaeomagnetic directions have to be inferred from a single sample in which case no  $\alpha_{95}$  angle can be recovered in this way. Fortunately, the progressive demagnetization techniques and principal component analysis universally used to recover directional information from single samples provide alternative measures of the error affecting the recovered direction, known as Maximum Angular Deviation (MAD) angles. These have so far only been considered as rough quality indicators. Here, however, we show that directions recovered in this way can be assumed to satisfy a Fisher distribution, and that the corresponding MAD angles can be rescaled into  $\alpha_{95}$  estimates by multiplying it by an appropriate factor, which only depends on the number of demagnetization steps used in the principal component analysis and on whether one relies on a standard or a so-called ‘anchored’ principal component analysis. These coefficients have been tabulated and practical recommendations for taking advantage of them outlined in the final section of the text. They provide simple means for users to produce much needed error bars on declination and inclination time series recovered from sedimentary long sequences.

**Key words:** Probability distributions; Magnetostratigraphy; Palaeomagnetic secular variation.

## 1 INTRODUCTION

Palaeomagnetic directions are usually measured and archived with some information about their precision, most often in the form of a so-called  $\alpha_{95}$  angle. This angle defines the size of the cone within which the true palaeomagnetic direction is expected to lie with 95 per cent confidence. It is traditionally co-estimated with the palaeomagnetic direction recovered from a collection of samples providing individual estimates of this direction. Such samples may for instance be collected from a lava flow expected to have uniformly recorded the same ‘instantaneous’ field (on time scales of up to a few months), the direction of which one wishes to recover. They may also be collected from an outcrop of sediments assumed to have been deposited over a much longer time period (say, a few tens of thousands of years), in which case the aim is to recover the direction of the local average paleofield over that period of time. In either case, the corresponding collection of directions is then dealt with as if drawn from a Fisher distribution (Fisher 1953), and best estimates of the direction to be recovered are computed together with an associated  $\alpha_{95}$  angle, using well-known formulae (Merrill *et al.* 1996; Tauxe *et al.* 2014). These  $\alpha_{95}$  can be used for many purposes,

such as deciding whether directions provided by different sampling sites are statistically different (e.g. McFadden & McElhinny 1990), or estimating the compatibility of large directional data sets with statistical models of the geomagnetic field (e.g. Khokhlov *et al.* 2006; Khokhlov & Hulot 2013).

In some instances, however, palaeomagnetic directions can only be recovered from one sample. This is particularly the case when one aims at reconstructing the behaviour of the past geomagnetic field from sediment cores. Since the field changes on time scales [down to centuries and decades, see e.g. Hulot & Le Mouél (1994), Hongre *et al.* (1998) or Lhuillier *et al.* (2011)] comparable to or shorter than the time elapsed between two consecutive sediment samples in such cores, each sample provides a unique palaeomagnetic direction. No  $\alpha_{95}$  angle can then be computed in the way described above. Fortunately, the demagnetization technique used to reconstruct directional information from individual samples still provides a measure of the uncertainty affecting the recovered direction. It is the possibility of converting this measure into an equivalent  $\alpha_{95}$  angle that we explore in this paper.

The technique used by palaeomagneticians to reconstruct directional information by progressive demagnetization of a sample

relies on the assumption that the total magnetic moment of such a sample is the sum of contributions from a population of magnetic carriers within the sample, and that only a fraction of this population is demagnetized at each demagnetization step (see e.g. Dunlop & Özdemir 2007). Between two successive steps, only the magnetic moment of this fraction of the population is removed. As a result, if a sequence of demagnetization steps happens to progressively demagnetize carriers originally magnetized within a common magnetic field, successive measurements of the magnetic moment between these demagnetization steps, plotted as a succession of points in an orthogonal vector plot [also known as a Zijderveld plot (Zijderveld 1967)], will statistically align themselves along the direction of the magnetizing field. If this direction is the one sought, it can then be recovered from the observed alignment, the dispersion of the points providing an indication of the directional uncertainty.

By far the most common way of estimating this direction and its uncertainty consists in relying on the so-called principal component analysis (Kirschvink 1980). This analysis consists in searching for a least-squares fit to the identified sequence of aligned points in the Zijderveld plot. Provided such a sequence of aligned points can indeed be found, it provides both a paleodirection estimate and the information needed to compute the Maximum Angular Deviation (MAD parameter), a measure of the uncertainty affecting this estimate. All modern software reconstructing paleodirections from demagnetization measurements provide such directional estimates with their MAD values (see e.g. Cogné 2003; Mazaud 2005; Tauxe *et al.* 2014, or the so-called PMGSC software, developed by R. Enkin at the Geological Survey of Canada, and directly available from its author).

Most often, however, MAD values are only used for sample selection purposes at the palaeomagnetic site level. When multiple samples are available, as is the case for the lava flow example given above, this makes sense. Samples displaying MAD values less than a certain threshold (typically  $5^\circ$ , see e.g. Johnson *et al.* 2008) can be selected to simultaneously infer the paleodirection and an  $\alpha_{95}$  angle, in which case individual MAD values do not particularly need to be published. But when only one sample is available, the MAD angle remains the only information quantifying the uncertainty with which the direction is recovered. Yet, also in such instances, and particularly in magnetostratigraphic studies, MAD angles are often just used as a mere qualitative indicator of the quality of the data (see e.g. Nagy & Valet 1993; Laj *et al.* 2006). As a result, errors truly affecting such data are either ignored (as is the case when only first order comparisons are being made between distant sediment cores, see e.g. Lund *et al.* 2005; Laj *et al.* 2006), arbitrarily assigned, and/or buried in some a posteriori model error (e.g. Leonhardt & Fabian 2007; Korte *et al.* 2009). Only in a few instances have MAD angles been tentatively used in a quantitative manner as a proxy for  $\alpha_{95}$  angles (e.g. Korte *et al.* 2005).

From a statistical point of view, however, MAD angles are not trivially related to  $\alpha_{95}$  angles. This drawback of the principal component analysis was recognized early on by Kent *et al.* (1983) who proposed an alternative method of recovering paleodirections from Zijderveld plots. But this alternative method failed to meet the success of its predecessor and today the vast majority of published palaeomagnetic studies still relies on the principal component analysis of Kirschvink (1980), keeping with the habit of providing uncertainties in terms of MAD angles. Fortunately, and as we shall explain in the rest of this paper, MAD angles can be converted into equivalent  $\alpha_{95}$  angles, in an approximate but very practical way.

To show this, we first introduce a simple statistical model of the way demagnetization techniques lead to Zijderveld plots, when

magnetic carriers can indeed be assumed to have been originally magnetized within a single common magnetic field (Section 2). We next use this model to generate Zijderveld plots of artificial samples assumed to have recorded known (imposed) palaeomagnetic direction. These plots are analysed with the standard principal component analysis and a variant known as the anchored principal component analysis (see e.g. Butler 1992; Mazaud 2005). A first useful result is then obtained, namely, that under the statistical assumptions made both analysis lead to recovered palaeomagnetic directions with close to Fisherian statistical distributions (Section 3). We next proceed and compare the statistical distribution of MAD estimates recovered from these two analysis with the statistical distribution of  $\alpha_{95}$  estimates recovered from a Fisher analysis of the same Zijderveld plots (Section 4). This reveals that MAD and  $\alpha_{95}$  estimates have different statistical properties and that these differences must be taken into account. In particular, we show that MAD estimates scale directional errors differently compared to  $\alpha_{95}$  estimates, and that this scaling depends on whether one relies on a standard or an anchored principal component analysis. This nevertheless also suggests that reconstructing an estimate of the  $\alpha_{95}$  associated with a direction recovered from a standard or an anchored principal component analysis is possible, at least in an approximate way, by simply multiplying the corresponding MAD estimates by some appropriate scaling factors. Still relying on the same statistical model of Zijderveld plots, we then derive asymptotic (when the number of demagnetization steps becomes large) ratios of the root mean square (rms) values inferred from the MAD and  $\alpha_{95}$  distributions, and derive candidate asymptotic scaling factors for both types of principal component analysis (Section 5). In practice, however, such asymptotic ratios are not appropriate enough for a variety of reasons we discuss (Section 6). This leads us to provide an explicit table of better suited scaling factors to be used to convert MAD estimates into  $\alpha_{95}$  estimates for any realistic number of demagnetization steps analysed in a given Zijderveld plot (Table 8). As an ultimate step, and because such scaling factors are appropriate only to the extent the Zijderveld plot analysed reasonably satisfy the statistical model we assumed, we also carry on a number of tests using various types of real data Zijderveld plots, extracted from representative published palaeomagnetic studies (Section 7). This finally leads us to conclude that the theoretically recommended factors of Table 8 can indeed safely be used to convert MAD angles into  $\alpha_{95}$  estimates when dealing with sediment data (to within 30 per cent possible relative error) and volcanic data (though in that case, the conversion is likely to provide a slightly overestimated  $\alpha_{95}$ , by up to 50 per cent). Last but not least, for readers not willing to get in the mathematical and technical details of this paper, we conclude with a summary and an illustration (Section 8) followed by some practical recommendations for producing not only  $\alpha_{95}$  estimates, but also 95 per cent confidence intervals on declinations and inclinations (Section 9). These recommendations can easily be implemented to routinely produce  $\alpha_{95}$  estimates and 95 per cent confidence intervals on declination and inclination time series recovered from sedimentary long sequences, whether using classical palaeomagnetic techniques on U-channels (Tauxe *et al.* 1983), or discrete samples.

## 2 RANDOM WALK MODEL OF ZIJDERVELD PLOTS

As already noted, the total magnetic moment of a sample is the sum of contributions from a population of magnetic carriers within the sample. Several techniques may be used to progressively

demagnetize these carriers and the sample [by progressively heating it, or submitting it to stronger and stronger demagnetizing alternative fields, or increasingly more powerful microwaves, see Hill *et al.* (2002); Dunlop & Özdemir (2007); Tauxe *et al.* (2014) for details]. Sediment cores can also be analysed using so-called ‘U-channel’ techniques that allow the core to be analysed section by section as it progressively slides through the apparatus, in which case the core does not have to be sliced into individual samples (see e.g. Tauxe *et al.* 1983; Nagy & Valet 1993). Independently of the demagnetization technique used, however, the way the magnetic moment of the sample (or section of a sediment core) behaves can always be approximated in the same manner. At each demagnetization step, only a fraction of the population of magnetic carriers is demagnetized. As a result, if all carriers were originally magnetized within a common magnetic field, successive measurements of the magnetic moment between these demagnetization steps, plotted as a succession of points in an orthogonal vector plot (Zijderveld plot) will statistically align themselves along the direction of this magnetizing field. The corresponding fraction of the plot, linking the final to the initial sample moments, may thus be viewed as a ‘walk’ produced by a succession of steps corresponding to different magnetic carrier’s moments, each providing an independent estimate of the direction to be recovered. It is the statistical behaviour of this ‘walk’ that we need to describe in a plausible way.

One of the difficulty in doing so is that, in practice, demagnetization steps are the result of some subjective decision of the operator. Palaeomagneticians often reduce these steps when approaching expected transitions between different populations of magnetic carriers (e.g. when approaching Curie temperatures in the case of thermal demagnetization). This is useful for assessing which type of magnetic carriers is most involved in the magnetization of the sample and for revealing differences in the directions recorded by each population. This results in a few smaller steps in the Zijderveld plots. However, when it comes to selecting data for a principal component analysis, palaeomagneticians rarely include all of these clustered parts of the Zijderveld plot, precisely because they usually correspond to changes in the recorded paleodirection. As a result, one may hope that the final practical behaviour of the data selected along a segment will indeed turn out to be close to that produced by adding independent random magnetic moments drawn from a common statistical distribution reflecting how faithfully these moments have recorded the paleodirection to be recovered.

Statistically modelling such a ‘walk’, however, is not trivial, particularly since errors in the measurements that helped produce the Zijderveld plots, and not only the physical dispersion of magnetic moments in the samples, must also be considered. A statistical model accounting for such issues has already been proposed by Kent *et al.* (1983). Their model, however, did not recognize the fundamental random walk nature of the Zijderveld plot produced by the dispersion of magnetic moments in the sample. In contrast, our model intrinsically assumes Zijderveld plots to result from such a walk. This characteristic, as we shall later see, turns out to be important in relating MAD angles to equivalent  $\alpha_{95}$  angles.

In the rest of this paper, and for the purpose of investigating the statistical properties of the principal component analysis and of MAD angles, we will thus assume that each step of the ‘walk’ along such a segment has been drawn from a common 3-D Gaussian distribution to produce what can be seen as an analogue of a Brownian motion with a known drift. We will describe each step  $\mathbf{S}_i$  in the 3-D Zijderveld plots as a systematic step of length  $\delta$  along the direction of the paleodirection to be recovered, which we arbitrarily assume to be the  $(1, 0, 0)$  direction, plus a random unbiased isotropic error

of  $(\alpha_1, \alpha_2, \alpha_3)_i$ , representing the uncertainty with which the corresponding magnetic carriers will have recorded this direction. In other words,

$$\mathbf{S}_i = \mathbf{D} + \mathcal{A}_i = (\delta, 0, 0) + (\alpha_1, \alpha_2, \alpha_3)_i, \quad (1)$$

where each  $\alpha_j$  is Gaussian with expectancy  $E(\alpha_j) = 0$  and for all  $j, j' = 1, 2, 3$ , components  $\alpha_j$  and  $\alpha_{j'}$  are uncorrelated with standard deviation  $\sigma(\alpha_j) = \sigma_\alpha$ .

Following Kent *et al.* (1983), we also introduce a statistical description of measurement errors. These measurement errors do not affect individual steps  $\mathbf{S}_i$  in the 3-D Zijderveld plots, but only their sum  $\mathcal{R}_k$  up to the latest step of interest (reflecting the fact that measurements are done on the total remaining magnetization after each demagnetization step). Formally, this means that the succession of points  $\mathcal{R}_k$  on the Zijderveld plot is being modelled as ( $\mathcal{R}_1$  and  $\mathcal{R}_n$  being, respectively, the remaining magnetization after the final demagnetization step, and the initial natural remanent magnetization):

$$\mathcal{R}_k = \sum_{i=1}^k \mathbf{S}_i + (\beta_1, \beta_2, \beta_3)_k, \quad (2)$$

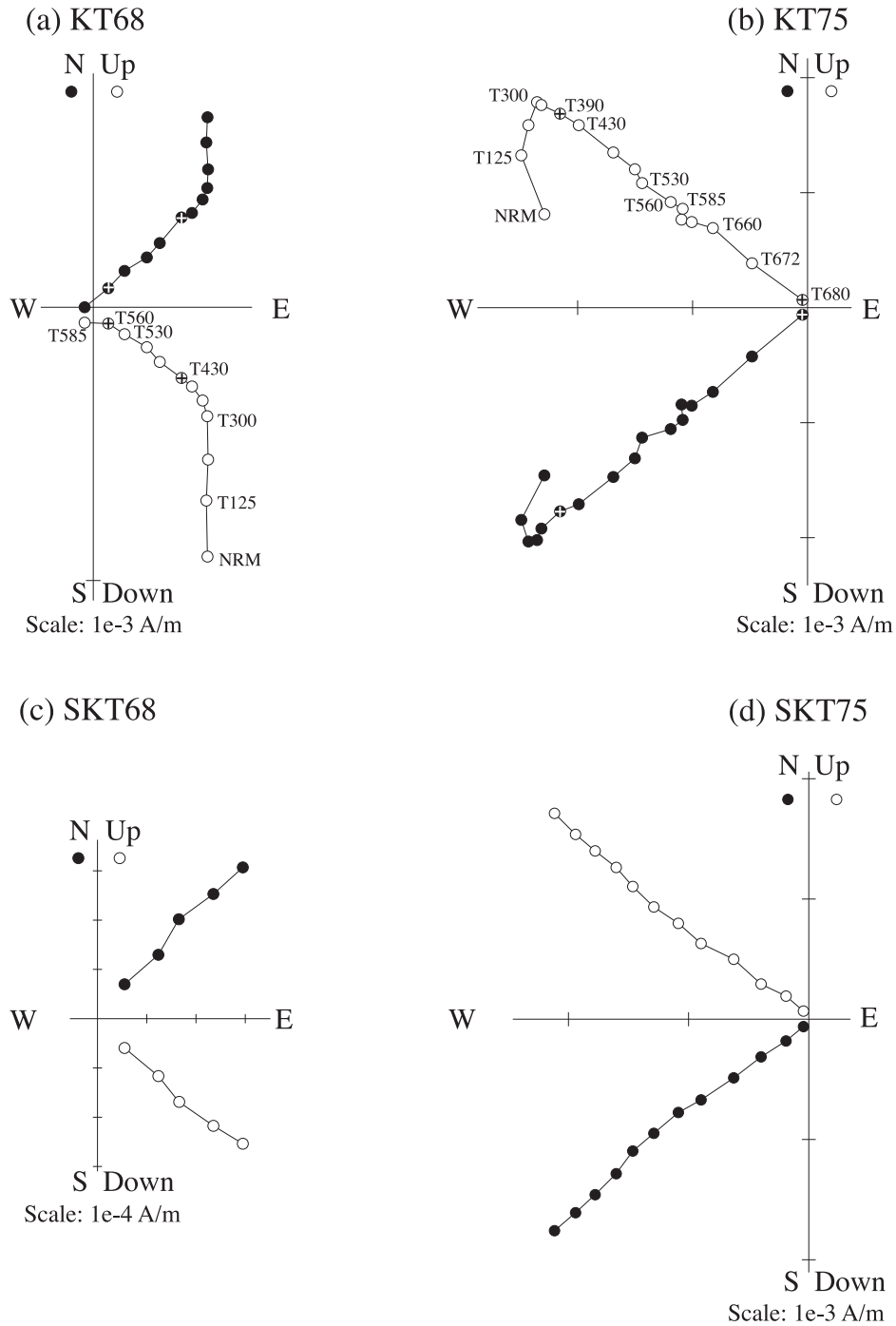
where each  $\beta_j$  is Gaussian with expectancy  $E(\beta_j) = 0$  and for all  $j, j' = 1, 2, 3$ , components  $\beta_j$  and  $\beta_{j'}$  are uncorrelated with standard deviation  $\sigma(\beta_j) = \sigma_\beta$ .

Producing synthetic Zijderveld plots then consists in choosing the amount  $n$  of steps to be produced,  $\delta$ ,  $\sigma_\alpha$  and  $\sigma_\beta$ . These parameters can be adjusted to produce plausible Zijderveld plots. Fig. 1 shows such synthetic plots together with real Zijderveld plots to illustrate this point when considering sediment data. For such data, and as we shall later see, Zijderveld plots can satisfyingly be described in terms of a pure random walk with no significant measurement errors (i.e. assuming  $\sigma_\beta = 0$ , as illustrated here). Because we first and foremost aim at converting MAD values derived from sediment data, this pure random walk model is in fact also the one we will mainly rely on in this paper (leaving the discussion of the impact of the Kent *et al.* (1983) measurement errors to the later Section 7). In what follows (and up to Section 7), we will thus assume  $\sigma_\beta = 0$ , and the relative size of the step  $\delta$  with respect to  $\sigma_\alpha$ ,  $d = \delta/\sigma_\alpha$ , will then be the only parameter controlling the directional error in the model. For convenience, we will set  $\sigma_\alpha = 1$  (and therefore  $\delta = d$ ) in all our simulations and mathematical derivations, and use  $d$  and  $n$  as the only two control parameters.

### 3 PRINCIPAL COMPONENT ANALYSIS OF SYNTHETIC ZIJDERVELD PLOTS

We now numerically investigate the statistical properties of the directions recovered with the standard principal component analysis of Kirschvink (1980), using sets of synthetic Zijderveld plots produced in the way just described (with  $\sigma_\beta = 0$ , since here we ignore measurement errors). The detailed justification for the principal component analysis can be found in Kirschvink (1980), and in classical books such as Tauxe *et al.* (2014). The practical of this analysis consists in the following steps.

Considering a given Zijderveld plot and having selected the  $n$  data points to be analysed, one first defines a Cartesian coordinate system within which each data point  $k$  is defined by its  $(x_{1k}, x_{2k}, x_{3k})$  coordinates. One next computes the ‘centre of mass’ of these points, of coordinates  $(\bar{x}_1, \bar{x}_2, \bar{x}_3)$ , where  $\bar{x}_j = \frac{1}{n} \sum_{k=1}^n x_{jk}$ , and the new coordinates  $(x'_{1k}, x'_{2k}, x'_{3k}) = (x_{1k} - \bar{x}_1, x_{2k} - \bar{x}_2, x_{3k} - \bar{x}_3)$  with respect



**Figure 1.** Comparing real Zijderveld plots with simulated Zijderveld plots produced using the statistical model assumed for this study. KT68 (a) and KT75 (b) correspond to two different sediment samples from the study of Pavlov & Gallet (2010); SKT68 (c) corresponds to a synthetic Zijderveld plot built with  $\sigma_\alpha = 1$ ,  $\sigma_\beta = 0$  and  $d = \delta/\sigma_\alpha = 1/0.12$  for  $n = 5$  steps to match the behaviour observed in KT68 between  $T = 430^\circ$  and  $T = 560^\circ$ ; likewise SKT75 (d) corresponds to a synthetic Zijderveld plot built with  $\sigma_\alpha = 1$ ,  $\sigma_\beta = 0$  and  $d = \delta/\sigma_\alpha = 1/0.12$  for  $n = 12$  steps to match the behaviour observed in KT75 between  $T = 390^\circ$  and  $T = 680^\circ$ .

to this centre of mass. This then makes it possible to compute the ‘orientation tensor’ (Scheidegger 1965)

$$T = \begin{pmatrix} \sum x'_{1k}x'_{1k} & \sum x'_{1k}x'_{2k} & \sum x'_{1k}x'_{3k} \\ \sum x'_{2k}x'_{1k} & \sum x'_{2k}x'_{2k} & \sum x'_{2k}x'_{3k} \\ \sum x'_{3k}x'_{1k} & \sum x'_{3k}x'_{2k} & \sum x'_{3k}x'_{3k} \end{pmatrix}, \quad (3)$$

where the sums are carried out over  $k = 1$  to  $n$ .

This matrix has three positive eigenvalues, which can be ranked arbitrarily as  $\lambda_1 \geq \lambda_2 \geq \lambda_3$ . The direction to be recovered is then the one provided by the eigenvector associated with  $\lambda_1$ , provided, as we assume throughout this paper, that  $\lambda_1$  is clearly greater than both  $\lambda_2$  and  $\lambda_3$ , thereby testifying for the fact that the data points analysed are indeed essentially aligned (see Kirschvink 1980, who also discusses the cases not considered here of more complex Zijderveld plots reflecting overlapping magnetization produced by



different magnetizing fields, in which case successive points in the Zijdeveld plot are no longer aligned, but either co-planar or even 3-D). For later purpose, we also introduce here the formula defining the associated MAD angle in such a case, namely,

$$\text{MAD} = \arctan \sqrt{\frac{\lambda_2 + \lambda_3}{\lambda_1}}. \quad (4)$$

A slightly different way of using the principal component analysis is also often used, which we will refer to as the ‘anchored’ principal component analysis (as in, e.g. Butler 1992; Mazaud 2005). This analysis differs from the previous one in that it explicitly requires the best fit to the selected data points to strictly go through the origin of the Zijdeveld plot. This can be achieved by simply introducing a symmetric data point (with respect to the origin of the plot) for each selected data point, before proceeding as in the standard analysis. This is equivalent to changing the definition (3) of the tensor  $T$  to

$$T = \begin{pmatrix} \sum x_{1k}x_{1k} & \sum x_{1k}x_{2k} & \sum x_{1k}x_{3k} \\ \sum x_{2k}x_{1k} & \sum x_{2k}x_{2k} & \sum x_{2k}x_{3k} \\ \sum x_{3k}x_{1k} & \sum x_{3k}x_{2k} & \sum x_{3k}x_{3k} \end{pmatrix}, \quad (5)$$

where the sums are again carried out over  $k = 1$  to  $n$ , the rest of the calculation being unchanged. In that case also, MAD angles can be calculated using eq. (4). For a given Zijdeveld plot, the MAD recovered from this anchored principal component analysis will usually be different from that recovered from the principal component analysis. From thereon, and for the purpose of clarity, MAD angles recovered from the anchored principal component analysis will thus be referred to as aMAD angles. Likewise, anchored principal component analysis will be referred to as an aPC analysis to avoid confusion with the principal component analysis, which we will refer to as a PC analysis.

In practice, the anchored analysis is used by palaeomagneticians only when the final steps of the Zijdeveld plot unambiguously tend towards the origin of the plot, indicative of a common so-called ‘primary’ component recorded by all the corresponding magnetic carriers. This is an ideal situation, and the use of the anchored analysis, which amounts to introduce some a priori information (a ‘primary’ component must indeed go through the origin), is expected to improve the recovery of the paleodirection. Otherwise, when the final (e.g. high temperature) demagnetization steps of a Zijdeveld plot are ambiguous or suggest a direction of different origin than the initial (low temperature) demagnetization steps, only the standard analysis can be used. Since both the anchored and standard analyses are used by palaeomagneticians, statistical properties of paleodirections recovered in both ways from our synthetic Zijdeveld plots will be investigated.

To produce the numerical examples of Zijdeveld plots that we need, realistic values have to be chosen for  $d$  and the number  $n$  of steps involved. We used  $d = 5$  and  $d = 10$ , and  $n = 3$  and  $n = 16$  as end cases (plus a limit  $n = 100$  case for checking convergence properties). We generated and investigated many sets of Zijdeveld plots combining values of  $d$  and  $n$  within these ranges.

Our investigation then proceeds in the following way. For each set of  $(d, n)$  values, we generate a minimum of 100 000 Zijdeveld plots (i.e. sets  $(\mathbf{S}_1, \dots, \mathbf{S}_n)$  of  $n$  steps as computed from eq. 1), each of which is analysed in four different ways: using the PC and aPC analysis that we wish to investigate, but also a Fisher analysis and what we will refer to as an Angular Gaussian analysis (Bingham 1983; Khokhlov *et al.* 2006). These last two analysis have known statistical behaviours and will prove useful for reference. In the case

of the Fisher analysis, the estimated direction is the direction of the resultant vector

$$\mathbf{R} = \sum_{i=1}^n \mathbf{s}_i, \quad (6)$$

where  $\mathbf{s}_i$  is the unit vector corresponding to  $\mathbf{S}_i$ . In the case of the Angular Gaussian analysis, the estimated direction is that of the sum of all steps (with  $\sigma_\beta = 0$ , since here we ignore measurement errors, recall eq. 2)

$$\mathcal{R} = \sum_{i=1}^n \mathbf{S}_i. \quad (7)$$

Each of these analysis thus produces an estimate of the direction to be recovered, and in each case, the angular distance  $\theta$  of this direction to the expected direction can be computed. This then makes it possible to construct empirical probability density functions (pdfs) from the over 100 000 Zijdeveld plots for the PC ( $f_{\text{PC}}(\theta)$ ), aPC ( $f_{\text{aPC}}(\theta)$ ), Fisher ( $f_{\text{F}}(\theta)$ ) and Angular Gaussian ( $f_{\text{AG}}(\theta)$ ) analysis. (Note that by construction, all these four analysis produce unbiased symmetric distributions about the expected direction, which is why one may reduce the pdfs to a function of  $\theta$ .)

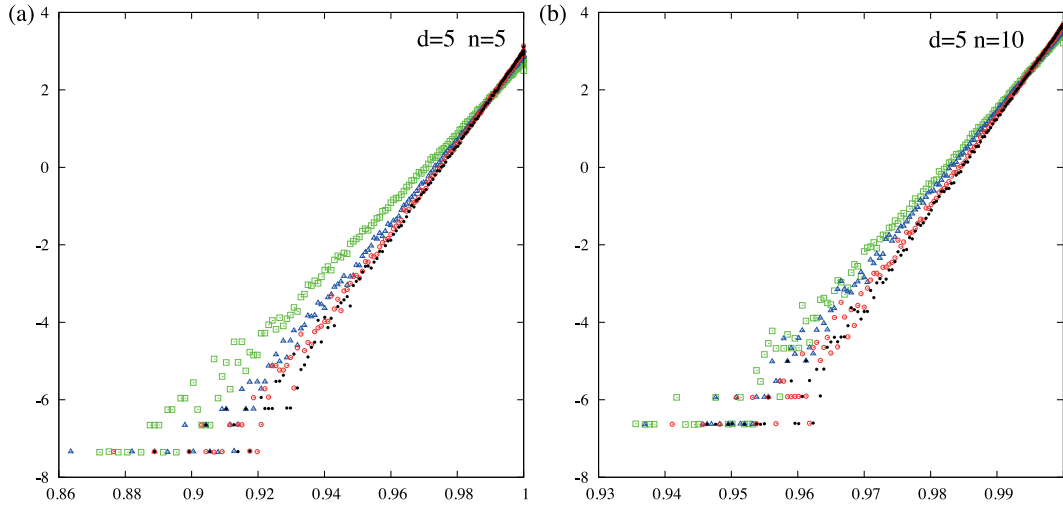
The behaviour of these empirical pdfs can next be checked against the one expected for a Fisher distribution, the pdf of which is of the form:

$$f_K(\theta) = \frac{K e^{K \cos \theta}}{2\pi (e^K - e^{-K})} \quad (8)$$

where  $K$  is the precision parameter (see e.g. Merrill *et al.* 1996; Tauxe *et al.* 2014). The prefactor being a simple normalization parameter, eq. (8) shows that checking the Fisherian behaviour of the  $f_{\text{PC}}(\theta)$ ,  $f_{\text{aPC}}(\theta)$ ,  $f_{\text{F}}(\theta)$  and  $f_{\text{AG}}(\theta)$  empirical pdfs can be done by checking the linear behaviour of the logarithm of these functions as a function of  $\cos \theta$ . Fig. 2 shows that such a linear behaviour is indeed observed for all four pdfs. The same behaviour has been found in all cases we investigated. This shows that all four analysis lead to estimates very close to being Fisherian. In the case of  $f_{\text{F}}(\theta)$  and  $f_{\text{AG}}(\theta)$  (the Fisher and Angular Gaussian analysis), this does not come as a surprise. Fisher and Angular Gaussian analysis of a set of vectors drawn from a 3-D Gaussian distribution are expected to produce a distribution very close to being Fisherian (see e.g. Bingham 1983; Khokhlov *et al.* 2006). What our results now show is that PC and aPC analysis of the same sets of vectors also produce  $f_{\text{PC}}(\theta)$  and  $f_{\text{aPC}}(\theta)$  distributions very close to being Fisherian.

Fig. 2, however, reveals that the PC analysis leads to a precision parameter  $K$  (slope of the pdfs, when plotted in this way) that can differ somewhat from those obtained with the aPC, Fisher and Angular Gaussian analysis (differences among these being weaker). Computing the empirical 0.95-quantile angle (defining the cone about the expected direction within which the recovered direction has a 95 per cent probability to lie) for each of these empirical pdfs confirms this slight discrepancy (Tables 1 and 2). These empirical values can also usefully be compared to the 0.95-quantile angles expected for the Angular Gaussian analysis. In that case indeed, the distribution of  $\mathcal{R}$  vectors as defined by eq. (7), knowing that each  $\mathbf{S}_i$  is produced as a result of eq. (1), is expected to be very well approximated by a Fisherian distribution with a precision parameter of  $K = nd^2$  (Khokhlov *et al.* 2006). Since eq. (8) leads to a 0.95 quantile angle of

$$Q_{95} = \arccos \left[ \frac{1}{K} \ln \left( e^K + 0.95(e^{-K} - e^K) \right) \right], \quad (9)$$



**Figure 2.** Natural logarithm of the  $f_{PC}(\theta)$  (green squares),  $f_{aPC}(\theta)$  (blue triangles),  $f_F(\theta)$  (red circles) and  $f_{\mathcal{AG}}(\theta)$  (black dots) empirical pdfs (to within some normalization factor, amounting to some arbitrary shift along the vertical axis in these plots) plotted as a function of  $\cos \theta$  for  $d = 5$  when considering  $n = 5$  (a) or  $n = 10$  (b). Note the clear linear behaviour of these plots (to within the dispersion that occurs for low probability values, due to the finite amount of realizations).

**Table 1.** Empirical 0.95-quantile angles  $Q_{95}$  for the principal component (PC), anchored principal component (aPC), Fisher (F) and Angular Gaussian ( $\mathcal{AG}$ ) analyses, considering  $n$  samples. Also provided,  $Q_{95}$  as computed from eq. (10) expressed in degrees. All results refer to simulations with  $d = 5$ . Note that empirical results are accurate only to within a few unit changes in the last digit.

$n$	$Q_{95}(\text{PC})$	$Q_{95}(\text{aPC})$	$Q_{95}(\text{F})$	$Q_{95}(\mathcal{AG})$	$Q_{95}$ of eq. (10)
3	20.12°	16.89°	16.56°	16.32°	16.19°
4	16.49°	14.78°	14.36°	14.12°	14.02°
5	14.43°	13.24°	12.83°	12.61°	12.54°
6	12.98°	12.15°	11.66°	11.46°	11.45°
7	11.88°	11.28°	10.79°	10.60°	10.60°
8	11.08°	10.59°	10.11°	9.94°	9.92°
9	10.40°	10.03°	9.55°	9.38°	9.35°
10	9.84°	9.55°	9.05°	8.90°	8.87°
11	9.39°	9.11°	8.63°	8.49°	8.46°
12	8.97°	8.72°	8.25°	8.11°	8.10°
13	8.57°	8.37°	7.93°	7.78°	7.78°
14	8.26°	8.10°	7.68°	7.53°	7.50°
15	7.98°	7.82°	7.39°	7.25°	7.24°
16	7.73°	7.59°	7.18°	7.05°	7.01°
100	3.08°	3.07°	2.86°	2.80°	2.80°

and since  $K$  is very large, this predicted value can be very well approximated by:

$$Q_{95} \sim \sqrt{\frac{2 \ln 20}{K}} \sim \frac{1}{d\sqrt{n}} \sqrt{2 \ln 20}, \quad (10)$$

where  $Q_{95}$  is expressed in radians. As can be seen in Tables 1 and 2, the Fisher and Angular Gaussian analysis indeed lead to empirical 0.95-quantile angles  $Q_{95}(\text{F})$  and  $Q_{95}(\mathcal{AG})$  very close to  $Q_{95}$ , whereas the aPC analysis leads to  $Q_{95}(\text{aPC})$  values slightly larger (within typically half of a degree), the PC analysis for small  $n$  leading to  $Q_{95}(\text{PC})$  values with the largest discrepancies.

Fortunately, that this should be the case can readily be understood. Contrary to the Fisher, Angular Gaussian and aPC analysis, the PC analysis does not rely on all of the information provided by the  $n$  points of a Zijderfeld plot. As explained above, a centre of mass is first removed, which essentially means that one degree of information is lost. Indeed, Tables 1 and 2 clearly show that a PC

**Table 2.** Same as Table 1, except that all results now refer to simulations with  $d = 10$ .

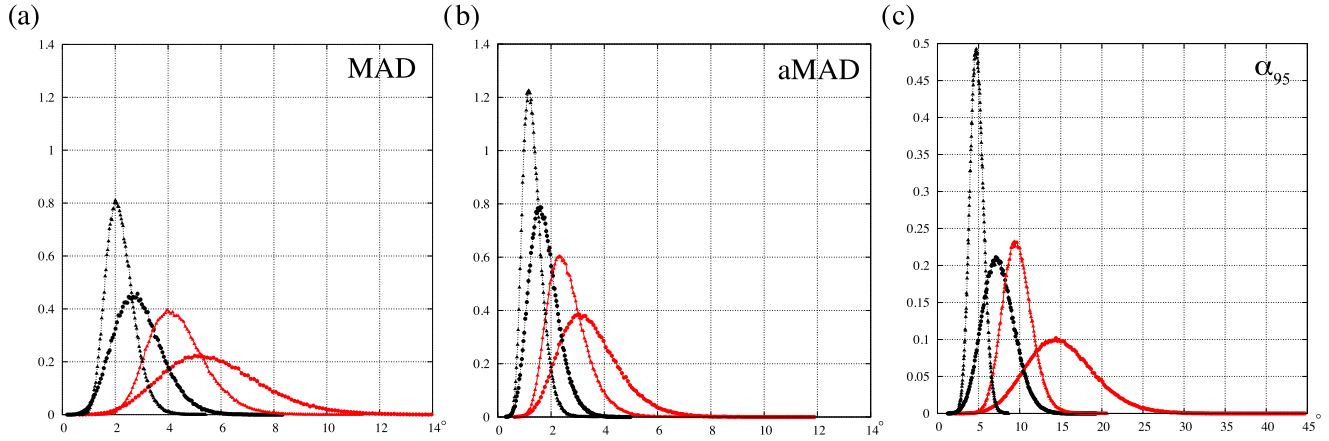
$n$	$Q_{95}(\text{PC})$	$Q_{95}(\text{aPC})$	$Q_{95}(\text{F})$	$Q_{95}(\mathcal{AG})$	$Q_{95}$ of eq. (10)
3	9.93°	8.39°	8.14°	8.12°	8.09°
4	8.21°	7.33°	7.04°	7.01°	7.01°
5	7.16°	6.63°	6.30°	6.28°	6.27°
6	6.44°	6.06°	5.74°	5.72°	5.74°
7	5.92°	5.64°	5.32°	5.30°	5.30°
8	5.52°	5.30°	4.98°	4.96°	4.96°
9	5.19°	5.01°	4.70°	4.68°	4.67°
10	4.90°	4.75°	4.46°	4.43°	4.43°
11	4.67°	4.53°	4.24°	4.22°	4.23°
12	4.48°	4.35°	4.06°	4.04°	4.05°
13	4.29°	4.18°	3.90°	3.89°	3.89°
14	4.13°	4.05°	3.77°	3.75°	3.75°
15	3.98°	3.91°	3.63°	3.61°	3.62°
16	3.86°	3.80°	3.53°	3.51°	3.51°
100	1.54°	1.54°	1.41°	1.41°	1.40°

analysis of a Zijderfeld plot with  $n$  points will lead to an empirical pdf with a 0.95-quantile angle  $Q_{95}(\text{PC})$  very close (to within half of a degree) to that of an empirical pdf produced by the analysis of a Zijderfeld plot with  $n - 1$  (rather than  $n$ ) points using any of the other methods.

We may thus conclude that, at least within the statistical assumptions made so far, both the PC analysis and aPC analysis of Zijderfeld plots will produce Fisherian estimates, with Fisherian pdfs (and thus  $\alpha_{95}$  angles characterizing them) very similar to those associated to Fisher and Angular Gaussian analysis of the same plots. Not surprisingly, our results also confirm that for a given Zijderfeld plot, the standard PC analysis will provide a slightly less accurate estimate of the paleodirection than the aPC analysis (which mainly affects Zijderfeld plots with very few points).

#### 4 COMPARING MAD, aMAD AND $\alpha_{95}$ STATISTICAL DISTRIBUTIONS

Since we now know that PC and aPC analysis of Zijderfeld plots both produce Fisherian estimates of the direction to be recovered



**Figure 3.** Empirical pdfs for the MAD (a), aMAD (b) and  $\alpha_{95}$  (c) angles when recovering directions from synthetic Zijderveld plots, using respectively a principal component, anchored principal component or Fisher analysis. In each plot, results for the same four sets of (over 100 000) synthetic Zijderveld plots are shown:  $d = 5, n = 5$  (red circles);  $d = 5, n = 10$  (red triangles);  $d = 10, n = 5$  (black circles);  $d = 10, n = 10$  (black triangles).

**Table 3.** Numerical estimations of the mean, median and rms values inferred from the MAD, aMAD and  $\alpha_{95}$  empirical pdfs, for various values of  $n$ , and  $d = 5$ . Note that these empirical results are accurate only to within a few unit changes in the last digit. Also provided are the MAD, aMAD and  $\alpha_{95}$  rms values predicted by the asymptotic formulae of eqs (29), (37) and (41) expressed in degrees.

$n$	MAD				aMAD				$\alpha_{95}$			
	mean	med.	rms	eq. (29)	mean	med.	rms	eq. (37)	mean	med.	rms	eq. (41)
3	5.79°	5.37°	6.58°	8.37°	4.02°	3.84°	4.34°	5.12°	23.78°	22.83°	25.52°	16.19°
4	6.00°	5.77°	6.44°	7.25°	3.75°	3.60°	3.97°	4.44°	17.99°	17.53°	18.87°	14.02°
5	5.73°	5.56°	6.03°	6.48°	3.48°	3.35°	3.66°	3.97°	15.13°	14.84°	15.68°	12.54°
6	5.40°	5.25°	5.64°	5.92°	3.25°	3.13°	3.40°	3.62°	13.33°	13.13°	13.72°	11.45°
7	5.10°	4.96°	5.29°	5.48°	3.06°	2.94°	3.19°	3.35°	12.07°	11.91°	12.36°	10.60°
8	4.82°	4.70°	4.98°	5.12°	2.89°	2.78°	3.00°	3.14°	11.10°	10.98°	11.34°	9.92°
9	4.59°	4.47°	4.74°	4.83°	2.75°	2.64°	2.86°	2.96°	10.35°	10.25°	10.54°	9.35°
10	4.37°	4.25°	4.50°	4.58°	2.61°	2.52°	2.71°	2.81°	9.72°	9.64°	9.88°	8.87°
11	4.19°	4.08°	4.31°	4.37°	2.51°	2.42°	2.60°	2.68°	9.20°	9.13°	9.34°	8.46°
12	4.02°	3.91°	4.13°	4.18°	2.41°	2.32°	2.49°	2.56°	8.76°	8.69°	8.87°	8.10°
13	3.87°	3.77°	3.98°	4.02°	2.32°	2.24°	2.40°	2.46°	8.37°	8.31°	8.47°	7.78°
14	3.74°	3.64°	3.83°	3.87°	2.24°	2.16°	2.32°	2.37°	8.03°	7.98°	8.12°	7.50°
15	3.62°	3.52°	3.71°	3.74°	2.17°	2.09°	2.25°	2.29°	7.74°	7.69°	7.82°	7.24°
16	3.50°	3.41°	3.60°	3.62°	2.10°	2.03°	2.18°	2.22°	7.46°	7.42°	7.54°	7.01°
100	1.42°	1.38°	1.45°	1.45°	0.86°	0.83°	0.88°	0.89°	2.88°	2.88°	2.89°	2.80°

(given, again, the statistical assumptions made so far), the next question to address is that of finding ways of recovering estimates of the corresponding  $\alpha_{95}$  from the only other piece of information provided by these analysis, the MAD and aMAD angles. To investigate this, we first compare the statistical distributions of the MAD, aMAD and  $\alpha_{95}$  estimates recovered from respectively the PC, aPC and Fisherian analysis of the same statistical samples, using the same numerical approach as in the previous section.

For each set of  $(d, n)$  values, we again consider the more than 100 000 Zijderveld plots already generated. For each of these plots, we consider the MAD angles produced as outcomes of the PC analysis (the MAD, using eqs 3 and 4) and aPC analysis (the aMAD, using eqs 5 and 4). In addition, we compute  $\alpha_{95}$  estimates associated with the Fisher analysis of each plot, using the well-known formula (which applies under very general conditions, see e.g. McFadden 1980; Merrill *et al.* 1996; Tauxe *et al.* 2014),

$$\alpha_{95} = \arccos \left[ 1 - \frac{n-R}{R} \left( \left( \frac{1}{0.05} \right)^{\frac{1}{n-1}} - 1 \right) \right] \quad (11)$$

where  $R$  is the length (norm) of  $\mathbf{R}$  defined by eq. (6). These estimates are then used to produce empirical pdfs for the MAD, aMAD, and

$\alpha_{95}$  angles. Plotting these pdfs for realistic  $(d, n)$  values reveals a number of interesting features (Fig. 3).

It first shows that all distributions display some significant width, particularly for samples with few steps (small  $n$ ) and large directional error (small  $d$ , recall Section 2). Not surprisingly, this width is a function of both  $d$  and  $n$ . All pdfs become sharper as  $n$  or  $d$  increases. The way the pdfs evolve with  $d$  and  $n$ , however, is not exactly the same for the three pdfs. In addition, both the MAD and aMAD distributions have a slightly longer ‘tail’ (i.e. are more dissymmetric with respect to their maximum values) than the  $\alpha_{95}$  distributions. Finally, and most important, is the fact that, for a given set of  $d$  and  $n$  values, the MAD, aMAD and  $\alpha_{95}$  pdfs do not peak at the same maximum values and do not lead to the same mean, median and rms (defined as the square root of the second moment of the pdf, without removing the mean) values (Tables 3 and 4). MAD pdfs, and even more so, aMAD pdfs, lead to much smaller values than  $\alpha_{95}$  pdfs. In particular, whereas  $\alpha_{95}$  estimates appear to be commensurate with the 0.95-quantile angles we previously inferred from the directions recovered simultaneously (recall Tables 1 and 2), this is not the case for the MAD and aMAD estimates. This clearly shows that MAD and aMAD estimates do not directly scale directional errors the way  $\alpha_{95}$  estimates do.



**Table 4.** Same as Table 3, but for  $d = 10$ .

$n$	MAD				aMAD				$\alpha_{95}$			
	mean	med.	rms	eq. (29)	mean	med.	rms	eq. (37)	mean	med.	rms	eq. (41)
3	2.93°	2.74°	3.31°	4.18°	2.02°	1.94°	2.17°	2.56°	11.67°	11.33°	12.44°	8.10°
4	3.03°	2.93°	3.24°	3.62°	1.88°	1.82°	1.98°	2.22°	8.86°	8.69°	9.25°	7.01°
5	2.89°	2.82°	3.03°	3.24°	1.75°	1.69°	1.83°	1.98°	7.46°	7.36°	7.70°	6.27°
6	2.72°	2.65°	2.83°	2.96°	1.63°	1.58°	1.70°	1.81°	6.57°	6.50°	6.74°	5.73°
7	2.56°	2.50°	2.65°	2.74°	1.53°	1.48°	1.59°	1.68°	5.94°	5.89°	6.07°	5.30°
8	2.42°	2.36°	2.50°	2.56°	1.45°	1.40°	1.50°	1.57°	5.47°	5.43°	5.57°	4.96°
9	2.30°	2.24°	2.37°	2.42°	1.37°	1.33°	1.42°	1.48°	5.09°	5.06°	5.17°	4.67°
10	2.19°	2.14°	2.26°	2.29°	1.31°	1.26°	1.36°	1.40°	4.79°	4.76°	4.86°	4.43°
11	2.10°	2.05°	2.16°	2.19°	1.25°	1.21°	1.30°	1.34°	4.53°	4.51°	4.59°	4.23°
12	2.02°	1.96°	2.07°	2.09°	1.21°	1.16°	1.25°	1.28°	4.31°	4.30°	4.36°	4.05°
13	1.94°	1.89°	1.99°	2.01°	1.16°	1.12°	1.20°	1.23°	4.12°	4.10°	4.17°	3.89°
14	1.87°	1.83°	1.92°	1.94°	1.12°	1.08°	1.16°	1.19°	3.96°	3.94°	4.00°	3.75°
15	1.81°	1.77°	1.86°	1.87°	1.09°	1.05°	1.12°	1.15°	3.81°	3.80°	3.85°	3.62°
16	1.76°	1.71°	1.80°	1.81°	1.05°	1.02°	1.09°	1.11°	3.68°	3.66°	3.71°	3.51°
100	0.71°	0.69°	0.72°	0.72°	0.43°	0.41°	0.44°	0.44°	1.42°	1.42°	1.42°	1.40°

**Table 5.** Numerical estimates of the ratios of the mean, median and rms values of the  $\alpha_{95}$  pdfs to the mean, median and rms values of the MAD and aMAD pdfs, for  $d = 5$ . Recall that these empirical results are expected to be accurate only to within a few unit changes in the last digit.

$n$	$\alpha_{95}$ to MAD ratios			$\alpha_{95}$ to aMAD ratios		
	$\frac{\text{mean}(\alpha_{95})}{\text{mean}(\text{MAD})}$	$\frac{\text{median}(\alpha_{95})}{\text{median}(\text{MAD})}$	$\frac{\text{rms}(\alpha_{95})}{\text{rms}(\text{MAD})}$	$\frac{\text{mean}(\alpha_{95})}{\text{mean}(\text{aMAD})}$	$\frac{\text{median}(\alpha_{95})}{\text{median}(\text{aMAD})}$	$\frac{\text{rms}(\alpha_{95})}{\text{rms}(\text{aMAD})}$
3	4.11	4.25	3.88	5.92	5.95	5.88
4	3.00	3.04	2.93	4.80	4.87	4.75
5	2.64	2.67	2.60	4.35	4.43	4.28
6	2.47	2.50	2.43	4.10	4.19	4.04
7	2.37	2.40	2.34	3.94	4.05	3.87
8	2.30	2.34	2.28	3.84	3.95	3.78
9	2.25	2.29	2.22	3.76	3.88	3.69
10	2.22	2.27	2.20	3.72	3.83	3.65
11	2.20	2.24	2.17	3.67	3.77	3.59
12	2.18	2.22	2.15	3.64	3.75	3.56
13	2.16	2.20	2.13	3.61	3.71	3.53
14	2.15	2.19	2.12	3.58	3.69	3.50
15	2.14	2.18	2.11	3.57	3.68	3.48
16	2.13	2.18	2.09	3.55	3.66	3.46
100	2.03	2.09	1.99	3.35	3.47	3.28

Fortunately, computing ratios of the mean, median and rms values of the  $\alpha_{95}$  pdfs to the mean, median and rms values of the MAD and aMAD pdfs, brings interesting additional insight (Tables 5 and 6). It now shows that for a given pair of  $(d, n)$  values, the mean, median and rms of the MAD pdf all scale in the same way (to within a maximum of 10 per cent in relative value, reflecting the differences in the shape of the pdfs) with respect to the mean, median and rms of the  $\alpha_{95}$  pdf. The corresponding ratios hardly depend on the value of  $d$ . In contrast, they are a function of  $n$ . But as  $n$  increases, this dependence becomes much weaker, with all ratios showing a tendency to converge to values of order 2. A very similar behaviour can be found when considering ratios with respect to aMAD mean, median and rms values, except for the fact that ratios now tend to converge to values slightly larger than 3.

This now suggests that reconstructing an estimate of the  $\alpha_{95}$  associated with a direction recovered from a PC analysis (resp. aPC analysis) might be possible, at least in an approximate way, by simply multiplying the corresponding MAD (resp. aMAD) by some appropriate factor with no dependence on  $d$  and vanishing dependence on  $n$  when  $n$  becomes large enough.

## 5 ASYMPTOTIC EXPRESSIONS FOR THE MAD, aMAD AND $\alpha_{95}$ RMS

To confirm the asymptotic behaviours tentatively identified in Tables 5 and 6, we now focus on the MAD, aMAD and  $\alpha_{95}$  rms and search for asymptotic expressions of these quantities as  $n$  becomes large.

### 5.1 MAD root mean square

First consider the case of the MAD recovered using the standard PC analysis. Given that we describe the successive points in the Zijderveld plot to be the result of a random walk with successive steps  $\mathbf{S}_i$  as defined by eq. (1), and that we ignore measurement errors, these successive points  $k$  will be defined by eq. (2) with  $\sigma_\beta = 0$ , that is, by

$$\mathcal{R}_k = \sum_{i=1}^k \mathbf{S}_i = k\mathbf{D} + \sum_{i=1}^k \mathcal{A}_i, \quad (12)$$

**Table 6.** Same as Table 5, but for  $d = 10$ .

$n$	$\alpha_{95}$ to MAD ratios			$\alpha_{95}$ to aMAD ratios		
	$\frac{\text{mean}(\alpha_{95})}{\text{mean(MAD)}}$	$\frac{\text{median}(\alpha_{95})}{\text{median(MAD)}}$	$\frac{\text{rms}(\alpha_{95})}{\text{rms(MAD)}}$	$\frac{\text{mean}(\alpha_{95})}{\text{mean(aMAD)}}$	$\frac{\text{median}(\alpha_{95})}{\text{median(aMAD)}}$	$\frac{\text{rms}(\alpha_{95})}{\text{rms(aMAD)}}$
3	3.98	4.14	3.76	5.77	5.84	5.73
4	2.92	2.97	2.85	4.71	4.77	4.67
5	2.58	2.61	2.54	4.26	4.36	4.21
6	2.42	2.45	2.38	4.03	4.11	3.96
7	2.32	2.36	2.29	3.88	3.98	3.82
8	2.26	2.30	2.23	3.77	3.88	3.71
9	2.21	2.26	2.18	3.72	3.80	3.64
10	2.19	2.22	2.15	3.66	3.78	3.57
11	2.16	2.20	2.13	3.62	3.73	3.53
12	2.13	2.19	2.11	3.56	3.71	3.49
13	2.12	2.17	2.10	3.55	3.66	3.47
14	2.12	2.15	2.08	3.54	3.65	3.45
15	2.10	2.15	2.07	3.50	3.62	3.44
16	2.09	2.14	2.06	3.50	3.59	3.40
100	2.00	2.06	1.97	3.30	3.46	3.23

the coordinates of which will be denoted  $(x_{1k}, x_{2k}, x_{3k})$ , for consistency with Section 3.

Because the contribution of the systematic term in eq. (12) (first term in the right-hand side) is proportional to  $k$ , it will prove useful to rescale this term into

$$\mathcal{L}_k = \frac{k}{n} \mathbf{D}, \tag{13}$$

with coordinates  $(l_k, 0, 0)$ , where  $l_k = (k/n)d$ . Likewise, since the contribution of the random walk term in eq. (12) (second term in the right-hand side) has coordinates with zero expectation and  $\sqrt{k}$  standard deviation (recall our assumptions for  $\mathcal{A}_i$  in Section 2), it will also prove useful to rescale this term into

$$\mathcal{W}_k = \frac{1}{\sqrt{n}} \sum_{i=1}^k \mathcal{A}_i, \tag{14}$$

with coordinates  $(w_{1k}, w_{2k}, w_{3k})$ . Then  $\mathcal{R}_k$  can be written as

$$\mathcal{R}_k = n\mathcal{L}_k + \sqrt{n}\mathcal{W}_k, \tag{15}$$

which highlights the leading dependence of each term with the total number of points  $n$ .

Ideally, our goal is now to start from eq. (15), compute analytical expressions for the eigenvalues  $\lambda_1, \lambda_2$  and  $\lambda_3$  of the  $T$  matrix defined by eq. (3), and infer the MAD from eq. (4). This is a very significant task. Fortunately this task can be by-passed, because all quantities of interest are expected to have asymptotic behaviours, as soon as  $n$  becomes large enough. In practice, we will thus just focus on expanding quantities of interest in series of  $\varepsilon = 1/\sqrt{n}$ , to recover leading order expressions for  $\lambda_1, \lambda_2 + \lambda_3$  and thus the MAD.

We first introduce the inner product:

$$\langle g, h \rangle = \frac{1}{n} \sum_{k=1}^n g_k h_k, \tag{16}$$

and the notations  $\bar{g} = \langle g, \mathbf{1} \rangle$  (inner product with the unit function, providing an average value  $\bar{g}$  of  $g$ ) and  $g' = g - \bar{g}$ . One can then write each element  $T_{ij}$  of the  $T$  matrix as

$$T_{ij} = n \langle x'_i, x'_j \rangle, \tag{17}$$

consistent with eq. (3). Each such element can next be expanded in the following forms, depending on which is considered:

$$\begin{aligned} T_{11} &= n^3 (\langle l', l' \rangle + 2 \langle w'_1, l' \rangle \varepsilon + \langle w'_1, w'_1 \rangle \varepsilon^2) \\ T_{1i} &= n^3 (\langle l', w'_i \rangle \varepsilon + \langle w'_1, w'_i \rangle \varepsilon^2) \quad i \neq 1 \\ T_{ij} &= n^3 \langle w'_i, w'_j \rangle \varepsilon^2 \quad i, j \neq 1. \end{aligned} \tag{18}$$

As far as  $l_k$  (the only non-zero coordinate of  $\mathcal{L}_k$ ) is concerned, one easily derives

$$\begin{aligned} \bar{l} &= \frac{d}{2} \left( 1 + \frac{1}{n} \right) = \frac{d}{2} + O(\varepsilon^2) \quad \text{and} \\ \langle l', l' \rangle &= \frac{d^2}{12} \left( 1 - \frac{1}{n^2} \right) = \frac{d^2}{12} + O(\varepsilon^4). \end{aligned} \tag{19}$$

Inferring values in eq. (18) for the inner products that involve the renormalized random term  $\mathcal{W}_k$  is less straightforward. However, precisely because  $\mathcal{W}_k$  has been renormalized, its random coordinates are known to take values within a finite range, independent of  $n$  (and thus  $\varepsilon$ ). These terms can thus be considered as order zero when compared to  $\varepsilon$ . This key property is what allows us to proceed.

We start by considering the trace of the  $T$  matrix, which provides the sum of all three eigenvalues  $\lambda_1 \geq \lambda_2 \geq \lambda_3$ :

$$\begin{aligned} \lambda_1 + \lambda_2 + \lambda_3 &= T_{11} + T_{22} + T_{33} \\ &= n^3 \left[ \frac{d^2}{12} + 2 \langle w'_1, l' \rangle \varepsilon + (\langle w'_1, w'_1 \rangle + \langle w'_2, w'_2 \rangle + \langle w'_3, w'_3 \rangle) \varepsilon^2 + o(\varepsilon^2) \right], \end{aligned} \tag{20}$$

where eq. (19) has been taken into account. We next consider the trace of the product  $T^T T$ , where  $T^T$  stands for the transpose of  $T$ . This now provides the sum of the squares of all three

eigenvalues:

$$\begin{aligned} \lambda_1^2 + \lambda_2^2 + \lambda_3^2 &= T_{11}^2 + T_{22}^2 + T_{33}^2 + 2(T_{12}^2 + T_{23}^2 + T_{13}^2) \\ &= n^6 \left[ \frac{d^4}{12^2} + 4 \frac{d^2}{12} \langle w'_1, l' \rangle \varepsilon + \left( 2 \frac{d^2}{12} \langle w'_1, w'_1 \rangle \right. \right. \\ &\quad \left. \left. + 4 \langle w'_1, l' \rangle^2 + 2 \langle w'_2, l' \rangle^2 + 2 \langle w'_3, l' \rangle^2 \right) \varepsilon^2 + o(\varepsilon^2) \right]. \end{aligned} \quad (21)$$

But each of these eigenvalues can itself be expanded in series of  $\varepsilon$  in the form of

$$\begin{aligned} \lambda_1 &= n^3 \left( \frac{d^2}{12} + a_1 \varepsilon + b_1 \varepsilon^2 + o(\varepsilon^2) \right), \\ \lambda_2 &= n^3 (a_2 \varepsilon + b_2 \varepsilon^2 + o(\varepsilon^2)), \\ \lambda_3 &= n^3 (a_3 \varepsilon + b_3 \varepsilon^2 + o(\varepsilon^2)), \end{aligned} \quad (22)$$

where  $a_2 \geq 0$  and  $a_3 \geq 0$  because all eigenvalues are non-negative; the factor  $n^3$  can be inferred from the same factor in eq. (18); and the lack of constant values in the bracket for  $\lambda_2$  and  $\lambda_3$  can be inferred from considering the limit case when  $n$  goes to infinity (and  $\varepsilon$  goes to zero).

Using these expansions into eqs (20) and (21) then leads to the following identities

$$\begin{cases} a_1 + a_2 + a_3 = 2 \langle w'_1, l' \rangle \\ b_1 + b_2 + b_3 = \langle w'_1, w'_1 \rangle + \langle w'_2, w'_2 \rangle + \langle w'_3, w'_3 \rangle \\ 2 \frac{d^2}{12} a_1 = 4 \frac{d^2}{12} \langle w'_1, l' \rangle \\ 2 \frac{d^2}{12} b_1 + a_1^2 + a_2^2 + a_3^2 = 2 \frac{d^2}{12} \langle w'_1, w'_1 \rangle + 4 \langle w'_1, l' \rangle^2 \\ \quad + 2 \langle w'_2, l' \rangle^2 + 2 \langle w'_3, l' \rangle^2. \end{cases} \quad (23)$$

From the first and third identities, we derive that  $a_1 = 2 \langle w'_1, l' \rangle$  and  $a_2 = a_3 = 0$ . This can next be used in the last identity to get  $b_1 = \langle w'_1, w'_1 \rangle + \frac{12}{d^2} \langle w'_2, l' \rangle^2 + \frac{12}{d^2} \langle w'_3, l' \rangle^2$ , which can also be used in the second identity to get  $b_2 + b_3 = \langle w'_2, w'_2 \rangle + \langle w'_3, w'_3 \rangle - \frac{12}{d^2} \langle w'_2, l' \rangle^2 - \frac{12}{d^2} \langle w'_3, l' \rangle^2$ . Coefficients  $a_1$  and  $b_1$  being known together with  $b_2 + b_3$  (and since  $a_2 = a_3 = 0$ ) expansion (22) for  $\lambda_1$  and  $\lambda_2 + \lambda_3$  can then be used to infer

$$\frac{\lambda_2 + \lambda_3}{\lambda_1} = \frac{12}{d^2} \sum_{i=2}^3 \left( \langle w'_i, w'_i \rangle - \frac{12}{d^2} \langle w'_i, l' \rangle^2 \right) \varepsilon^2 + o(\varepsilon^2). \quad (24)$$

This can next be used in eq. (4) to infer

$$\text{MAD} = \varepsilon \frac{\sqrt{12}}{d} \sqrt{\sum_{i=2}^3 \left( \langle w'_i, w'_i \rangle - \frac{12}{d^2} \langle w'_i, l' \rangle^2 \right) + o(\varepsilon)}, \quad (25)$$

which provides the MAD to leading order in  $\varepsilon$  as a function of the coordinates of the random walk  $\mathcal{R}_k$ , for any specific Zijderfeld plot.

In principle, eq. (25) could then be used to infer any of the quantities of the MAD distribution we investigated in Section 4. From an analytical point of view, however, the simplest quantity to derive is the rms of the MAD, that is  $\sqrt{E(\text{MAD}^2)}$  (where we recall that  $E()$  stands for the mathematical expectancy). Deriving this rms requires lengthy, but tractable, calculations, which we just summarize. We first note that

$$\begin{aligned} \langle w'_i, w'_i \rangle &= \langle w_i, w_i \rangle - 2 \langle w_i, \bar{w}_i \rangle + \langle \bar{w}_i, \bar{w}_i \rangle \\ \langle w'_i, l' \rangle^2 &= \langle w_i - \bar{w}_i, l \rangle^2 = \langle w_i, l \rangle^2 + \langle \bar{w}_i, \bar{l} \rangle^2 - 2 \langle w_i, l \rangle \langle \bar{w}_i, \bar{l} \rangle, \end{aligned} \quad (26)$$

and next derive the following intermediate results

$$\begin{aligned} E(\langle w_i, w_i \rangle) &= \frac{n+1}{2n} = \frac{1}{2} + O(\varepsilon^2); \\ E(\langle \bar{w}_i, \bar{w}_i \rangle) &= E(\langle w_i, \bar{w}_i \rangle) = \frac{(n+1)(n+\frac{1}{2})}{3n^2} = \frac{1}{3} + O(\varepsilon^2) \\ E(\langle \bar{w}_i, \bar{l} \rangle^2) &= \frac{d^2 (n+1)^3}{12 n^4} \left( n + \frac{1}{2} \right) = d^2 \left( \frac{1}{12} + O(\varepsilon^2) \right) \\ E(\langle w_i, l \rangle^2) &= \frac{d^2 n+1}{30 n^4} (4n^3 + 6n^2 + 4n + 1) = d^2 \left( \frac{2}{15} + O(\varepsilon^2) \right) \\ E(\langle w_i, l \rangle \langle \bar{w}_i, \bar{l} \rangle) &= \frac{d^2 (n+1)^2}{48 n^4} (5n^2 + 5n + 2) = d^2 \left( \frac{5}{48} + O(\varepsilon^2) \right), \end{aligned} \quad (27)$$

where it can be seen that, as anticipated, all quantities remain finite when  $n$  becomes large and  $\varepsilon$  becomes small. We may then finally conclude from eq. (25) that

$$\begin{aligned} \text{rms}(\text{MAD}) &= \sqrt{E(\text{MAD}^2)} = \frac{\varepsilon}{d} \sqrt{\frac{8}{5}} + o(\varepsilon) \\ &= \frac{1}{d\sqrt{n}} \sqrt{\frac{8}{5}} + o\left(\frac{1}{\sqrt{n}}\right). \end{aligned} \quad (28)$$

Thus, to leading order in  $n$ ,

$$\text{rms}(\text{MAD}) \sim \frac{1}{d\sqrt{n}} \sqrt{\frac{8}{5}} \quad (29)$$

where rms(MAD) is expressed in radians.

## 5.2 aMAD root mean square

A similar calculation can next be done for the aMAD recovered from the aPC analysis. We now start from definition (5) of the  $T$  matrix, the  $T_{ij}$  elements of which can be written as

$$T_{ij} = n \langle x_i, x_j \rangle. \quad (30)$$

Each such element can next be expanded in the following forms, depending on which is considered:

$$\begin{aligned} T_{11} &= n^3 (\langle l, l \rangle + 2 \langle w_1, l \rangle \varepsilon + \langle w_1, w_1 \rangle \varepsilon^2) \\ T_{ii} &= n^3 (\langle l, w_i \rangle \varepsilon + \langle w_1, w_i \rangle \varepsilon^2) \quad i \neq 1 \\ T_{ij} &= n^3 \langle w_i, w_j \rangle \varepsilon^2 \quad i, j \neq 1, \end{aligned} \quad (31)$$

and one can easily derive

$$\langle l, l \rangle = \frac{d^2}{3} \left( 1 + \frac{3}{2n} + \frac{1}{2n^2} \right) = \frac{d^2}{3} + \frac{d^2}{2n} + O(\varepsilon^4). \quad (32)$$

Following the same line of reasoning as in the previous case, one can next infer the sum and quadratic sum of all three eigenvalues  $\lambda_1 \geq \lambda_2 \geq \lambda_3$ :

$$\begin{aligned} \lambda_1 + \lambda_2 + \lambda_3 &= n^3 \left[ \frac{d^2}{3} + 2 \langle w_1, l \rangle \varepsilon + \left( \langle w_1, w_1 \rangle + \langle w_2, w_2 \rangle \right. \right. \\ &\quad \left. \left. + \langle w_3, w_3 \rangle + \frac{d^2}{2} \right) \varepsilon^2 + o(\varepsilon^2) \right], \end{aligned}$$

$$\lambda_1^2 + \lambda_2^2 + \lambda_3^2 = n^6 \left[ \frac{d^4}{3^2} + 4 \frac{d^2}{3} \langle w_1, l \rangle \varepsilon + \left( 2 \frac{d^2}{3} \left( \langle w_1, w_1 \rangle + \frac{d^2}{2} \right) + 4 \langle w_1, l \rangle^2 + 2 \langle w_2, l \rangle^2 + 2 \langle w_3, l \rangle^2 \right) \varepsilon^2 + o(\varepsilon^2) \right], \quad (33)$$

from which one can derive

$$\lambda_1 = \frac{d^2}{3} + 2 \langle w_1, l \rangle \varepsilon + o(\varepsilon)$$

$$\lambda_2 + \lambda_3 = \left[ \langle w_2, w_2 \rangle + \langle w_3, w_3 \rangle - \frac{3}{d^2} (\langle w_2, l \rangle^2 + \langle w_3, l \rangle^2) \right] \varepsilon^2 + o(\varepsilon^2), \quad (34)$$

and thus

$$\frac{\lambda_2 + \lambda_3}{\lambda_1} = \frac{3}{d^2} \sum_{i=2}^3 \left( \langle w_i, w_i \rangle - \frac{3}{d^2} \langle w_i, l \rangle^2 \right) \varepsilon^2 + o(\varepsilon^2). \quad (35)$$

Eq. (35) can then be used with results from (27) to finally derive the rms of the aMAD:

$$\begin{aligned} \text{rms(aMAD)} &= \sqrt{E(\text{aMAD}^2)} = \frac{\varepsilon}{d} \sqrt{\frac{3}{5}} + o(\varepsilon) \\ &= \frac{1}{d\sqrt{n}} \sqrt{\frac{3}{5}} + o\left(\frac{1}{\sqrt{n}}\right). \end{aligned} \quad (36)$$

Thus, to leading order in  $n$ ,

$$\text{rms(aMAD)} \sim \frac{1}{d\sqrt{n}} \sqrt{\frac{3}{5}} \quad (37)$$

where rms(aMAD) is expressed in radians.

### 5.3 $\alpha_{95}$ root mean square

Having analytical expressions for the asymptotic behaviour of the MAD and aMAD rms, we now wish to derive an analogous expression for the asymptotic behaviour of the  $\alpha_{95}$  rms. This can be done in a much more straightforward way, starting from eq. (11). Developing the arccos and exponential functions to leading order first leads to

$$\alpha_{95}^2 \sim 2 \ln 20 \frac{n - R}{(n - 1)R}. \quad (38)$$

Introducing the following  $k$  parameter

$$k = \frac{n - 1}{n - R} \quad (39)$$

next allows to transform the previous equation into (to leading order in  $n$ ):

$$\alpha_{95}^2 \sim \frac{2 \ln 20}{kR} \sim \frac{2 \ln 20}{(k - 1)n}. \quad (40)$$

Finally, considering that the statistical model we use for the Zijdeveld plots (recall eq. 1) implies that the individual unit vectors  $\mathbf{s}_i$  in eq. (6) are very close to satisfy a Fisher distribution with precision parameter  $\kappa = d^2$  (Khokhlov *et al.* 2006), that relevant values of  $k$  are typically larger than 25 so that  $(k - 1) \sim k$  to within a few percent at most, and that for such a Fisher distribution  $k^{-1}$  is an unbiased estimate for  $\kappa^{-1}$  (McFadden 1980), we may finally conclude that to

leading order in  $n$  (and to within a relative error commensurate with  $\kappa^{-1}$ )

$$\text{rms}(\alpha_{95}) = \sqrt{E(\alpha_{95}^2)} \sim \frac{1}{d\sqrt{n}} \sqrt{2 \ln 20}, \quad (41)$$

where rms( $\alpha_{95}$ ) is expressed in radians.

## 6 CONVERTING MAD AND aMAD INTO $\alpha_{95}$ ESTIMATES

As can be checked, eqs (29), (37) and (41) all reasonably predict the computed empirical MAD, aMAD and  $\alpha_{95}$  rms values as soon as  $n$  becomes large enough (see Tables 3 and 4, which also provide the corresponding predicted values). This is particularly true for the empirical MAD and aMAD rms which quickly converge towards the values predicted by eqs (29) and (37). Convergence is slower in the case of the  $\alpha_{95}$  rms, but the trend is clear. Eqs (29), (37) and (41) can also be used to predict the asymptotic values of the two rms( $\alpha_{95}$ )/rms(MAD) and rms( $\alpha_{95}$ )/rms(aMAD) ratios. This leads to

$$\frac{\text{rms}(\alpha_{95})}{\text{rms(MAD)}} \sim C_{\text{MAD}}^* = \sqrt{\frac{5 \ln 20}{4}} = 1.935... \quad (42)$$

and

$$\frac{\text{rms}(\alpha_{95})}{\text{rms(aMAD)}} \sim C_{\text{aMAD}}^* = \sqrt{\frac{10 \ln 20}{3}} = 3.160... \quad (43)$$

both consistent with the trends seen in the ratios of the corresponding empirical rms (Tables 5 and 6). Note, however, that the convergence is again not very quick, reflecting the slow convergence of the  $\alpha_{95}$  rms towards the values predicted by eq. (41).

If we now further note that the asymptotic rms( $\alpha_{95}$ ) of eq. (41) is identical to the  $Q_{95}$  predicted by eq. (10), and that this  $Q_{95}$  provides a measure converging towards the empirical 0.95 quantile angles  $Q_{95}(\mathcal{AG})$  and  $Q_{95}(\text{F})$  for the directions inferred from both the Angular Gaussian and Fisher analysis (recall Section 3 and Tables 1 and 2), this also implies that, as  $n$  gets larger, the rms( $\alpha_{95}$ ) value more and more accurately predicts these 0.95 quantile angles. This then suggests that multiplying the MAD (resp. aMAD) angle recovered from a PC analysis (resp. aPC analysis) by  $C_{\text{MAD}}^*$  (resp.  $C_{\text{aMAD}}^*$ ), as provided by eq. (42) (resp. eq. 43), could provide a reasonable estimate of the  $\alpha_{95}$  error affecting directions recovered from such a PC analysis (resp. aPC analysis). In particular, one would then expect the rescaled  $C_{\text{MAD}}^* \times \text{rms(MAD)}$  (resp.  $C_{\text{aMAD}}^* \times \text{rms(aMAD)}$ ) rms to properly predict the  $Q_{95}(\text{PC})$  (resp.  $Q_{95}(\text{aPC})$ ) quantiles.

Unfortunately, this is not the case. Multiplying rms(MAD) for  $d = 10$  and  $n = 16$  (Table 4) by  $C_{\text{MAD}}^* \sim 1.94$ , for instance, leads to  $3.49^\circ$ , whereas  $Q_{95}(\text{PC}) = 3.86^\circ$  (Table 2). This corresponds to a 10 per cent underestimation. For smaller values of  $n$ , and as one can easily check, discrepancies become even larger. The small size of  $n$  (hence the lack of convergence), however, is not the only cause of this mismatch. Indeed, multiplying rms(MAD) for  $d = 10$  and  $n = 100$  (Table 4) by  $C_{\text{MAD}}^* \sim 1.94$ , still reveals a comparable underestimation of order 10 per cent (it leads to  $1.40^\circ$ , whereas  $Q_{95}(\text{PC}) = 1.54^\circ$ , Table 2). This contrasts with the fact that for the same  $d = 10$  and  $n = 100$ , the rms( $\alpha_{95}$ ) value exactly matches both  $Q_{95}(\mathcal{AG})$  and  $Q_{95}(\text{F})$  (to within numerical accuracy see Tables 2 and 4). Similar underestimations are also revealed when testing rms(aMAD) values multiplied by  $C_{\text{aMAD}}^* \sim 3.16$ .

One cause of these discrepancies lies in the fact, observed in Tables 1 and 2, that even for large values of  $n$ ,  $Q_{95}(\text{PC})$  and  $Q_{95}(\text{aPC})$  never converge towards  $Q_{95}(\mathcal{AG})$  and  $Q_{95}(\text{F})$  (note that, in contrast,

**Table 7.**  $C_{\text{MAD}}$  and  $C_{\text{aMAD}}$  scaling factors as inferred by the method described in the text (Section 6) to convert MAD and aMAD angles into  $\alpha_{95}$  estimates, as a function of  $n$  for  $d = 5$  (left columns) and  $d = 10$  (right columns). Recall that these empirical results are expected to be accurate only to within a few unit changes in the last digit.

$n$	$d = 5$		$d = 10$	
	$C_{\text{MAD}}$	$C_{\text{aMAD}}$	$C_{\text{MAD}}$	$C_{\text{aMAD}}$
3	7.65	5.96	7.73	6.05
4	3.89	4.95	3.90	5.04
5	3.15	4.61	3.20	4.64
6	2.86	4.41	2.90	4.45
7	2.69	4.29	2.73	4.32
8	2.62	4.23	2.63	4.24
9	2.56	4.16	2.58	4.20
10	2.53	4.14	2.54	4.14
11	2.50	4.10	2.51	4.14
12	2.47	4.10	2.48	4.11
13	2.45	4.08	2.46	4.08
14	2.43	4.08	2.45	4.08
15	2.42	4.05	2.44	4.07
16	2.42	4.05	2.43	4.05
100	2.37	3.99	2.37	3.99

$Q_{95}(\text{PC})$  and  $Q_{95}(\text{aPC})$ , on one hand, and  $Q_{95}(\mathfrak{A}\mathfrak{G})$  and  $Q_{95}(\text{F})$ , on the other hand, converge towards each other). As a result, whereas rescaling MAD and aMAD values by  $C_{\text{MAD}}^*$  and  $C_{\text{aMAD}}^*$  leads to rms values perfectly converging towards  $Q_{95}(\mathfrak{A}\mathfrak{G})$  and  $Q_{95}(\text{F})$  (as can indeed be checked), these values fail to converge towards  $Q_{95}(\text{PC})$  and  $Q_{95}(\text{aPC})$ . These circumstances, together with the fact that  $n$  is never large in practice, implies that factors more appropriate than  $C_{\text{MAD}}^*$  and  $C_{\text{aMAD}}^*$  should be looked for to convert MAD and aMAD values into  $\alpha_{95}$  values in a more self-consistent way.

Focussing on the most intrinsic property an  $\alpha_{95}$  estimate is expected to have allows this to be done in an elegant way, whatever the values of  $d$  and  $n$ . This property states that when considering a set of directions known to satisfy a common Fisher distribution, each provided with its individual  $\alpha_{95}$  estimate, these directions should lie within their respective  $\alpha_{95}$  estimates of the true direction, 95 per cent of the time. Indeed, this property is very accurately satisfied by the directions and  $\alpha_{95}$  estimates recovered from the more than 100 000 Zijderfeld plots we numerically produced for each considered values of  $d$  and  $n$ , when using a Fisher analysis (i.e. eqs 6 and 11).

Since a PC analysis of the same Zijderfeld plots also produces estimates of the direction that satisfy a Fisher distribution (even for low values of  $n$ , recall Section 3), the following procedure can be applied: for each Zijderfeld plot, compute the angle  $\theta$  between the recovered direction and the known true direction, divide this angle by the MAD angle recovered from the same sample, that is, compute the  $\theta_{\text{MAD}} = \theta/\text{MAD}$  rescaled (non-dimensional) angle, build the  $\theta_{\text{MAD}}$  empirical pdf from the more than 100 000 Zijderfeld plots, and finally look for the 0.95 quantile associated with this pdf. If one denotes  $C_{\text{MAD}}$  this 0.95 quantile, then, as can readily be checked, multiplying each individual MAD value by this new  $C_{\text{MAD}}$  factor will produce  $\alpha_{95}$  estimates exactly satisfying the required property. Obviously, exactly the same procedure can also be applied to infer a new  $C_{\text{aMAD}}$  factor to convert aMAD values into  $\alpha_{95}$  estimates again satisfying the required property.

Table 7 provides the corresponding  $C_{\text{MAD}}$  and  $C_{\text{aMAD}}$  factors for  $d = 5$  and  $d = 10$ . As one would have hoped,  $C_{\text{MAD}}$  and  $C_{\text{aMAD}}$  val-

**Table 8.** Recommended  $C_{\text{MAD}}$  and  $C_{\text{aMAD}}$  scaling factors to convert MAD and aMAD angles into  $\alpha_{95}$  estimates, as a function of number of steps  $n$  analysed in a Zijderfeld plot.

$n$	$C_{\text{MAD}}$	$C_{\text{aMAD}}$
3	7.69	6.00
4	3.90	5.00
5	3.18	4.63
6	2.88	4.43
7	2.71	4.31
8	2.63	4.24
9	2.57	4.18
10	2.54	4.14
11	2.51	4.12
12	2.48	4.11
13	2.46	4.08
14	2.44	4.08
15	2.43	4.06
16	2.43	4.05
100	2.37	3.99

ues only very slightly depend on the value of  $d$ . But as anticipated, all  $C_{\text{MAD}}$  and  $C_{\text{aMAD}}$  values significantly differ from respectively  $C_{\text{MAD}}^*$  and  $C_{\text{aMAD}}^*$ , even for  $n = 16$  and  $n = 100$ . They nevertheless converge towards limit values, which are already very well approximated by the  $C_{\text{MAD}} \sim 2.37$  and  $C_{\text{aMAD}} \sim 3.99$  values recovered for  $n = 100$ . One less expected, but important, property of these new scaling factors is that they, too, fail to rescale the rms(MAD) and rms(aMAD) values into correct predictions of the  $Q_{95}(\text{PC})$  and  $Q_{95}(\text{aPC})$  quantiles. If one again considers the case of  $d = 10$  and  $n = 16$ , for instance, and multiply rms(MAD) =  $1.80^\circ$  by the corresponding  $C_{\text{MAD}} = 2.43$ , one then gets  $4.37^\circ$ , which overestimates  $Q_{95}(\text{PC}) = 3.86^\circ$  by roughly 10 per cent. A similar overestimation is observed when considering  $d = 10$  and  $n = 100$ , and also when considering rms(aMAD) multiplied by  $C_{\text{aMAD}}$ . In other words, whereas rescaling the MAD and aMAD angles using the  $C_{\text{MAD}}^*$  and  $C_{\text{aMAD}}^*$  factors led to rms underestimating the actual  $Q_{95}(\text{PC})$  and  $Q_{95}(\text{aPC})$  quantiles, using the new  $C_{\text{MAD}}$  and  $C_{\text{aMAD}}$  factors now leads to rms overestimating these quantiles. This undesired property highlights the important fact that contrary to true  $\alpha_{95}$  estimates, rescaled MAD and aMAD angles cannot provide unbiased rms estimates of the 0.95 quantile and simultaneously satisfy the intrinsic property we used above to recover the  $C_{\text{MAD}}$  and  $C_{\text{aMAD}}$  factors. Unfortunately, only one property of  $\alpha_{95}$  estimates can be guaranteed at a time, reflecting the fact that a single scaling factor cannot exactly recast MAD and aMAD pdfs into an  $\alpha_{95}$  pdf (recall Fig. 3). Since for all practical purposes, the one property one most expects from  $\alpha_{95}$  estimates is that they define the angular limit within which the true direction should lie 95 per cent of the time, it is our recommendation that these new  $C_{\text{MAD}}$  and  $C_{\text{aMAD}}$  factors be used to rescale all MAD and aMAD angles into  $\alpha_{95}$  estimates. Fortunately, this can readily be done, thanks to the fact already pointed out that these factors have very little dependence with  $d$ , so that averaging factors recovered from simulations with  $d = 5$  and  $d = 10$  can be used to produce  $d$ -independent values accurate within better than a few percent. Table 8 provides the corresponding recommended  $C_{\text{MAD}}$  and  $C_{\text{aMAD}}$  factors as a function of the number  $n$  of demagnetization steps involved in the principal component and anchored principal component analysis of a Zijderfeld plot.



## 7 TESTING THE ZIJDERVELD PLOT RANDOM WALK MODEL AGAINST REAL DATA

Before concluding, we finally turn to the important question of the relevance of the random walk based statistical model of Zijderfeld plots we assumed to derive all our results so far. We argued in Section 2 that this model was producing plausible Zijderfeld plots, as illustrated by, for example, Fig. 1. But visual analogy does not constitute solid evidence. In addition, and as already noted, the freedom palaeomagneticians may take in choosing the successive demagnetization steps can introduce variations in the statistical properties of Zijderfeld plots, not to mention the nature of the rocks analysed and the techniques used to demagnetize the sample, all of which may also have some influence. We therefore decided to carry out tests on a variety of Zijderfeld plots extracted from published studies expected to be representative of plots used in practice.

Ideally, testing our model would entail collecting large sets of real Zijderfeld plots, and for each plot, checking that the distribution of successive points is consistent with the distribution of  $\mathcal{R}_k$  points as predicted by eqs (1) and (2), for some known value  $n$  of steps but unknown values of  $\delta$ ,  $\sigma_\alpha$  and  $\sigma_\beta$  and an unknown drift direction, amounting to a total of five unknown parameters. But the small value of  $n$  (compared to the five unknown parameters) would clearly make such tests difficult to carry out in a conclusive way. Besides the key issue of interest is not so much whether the Zijderfeld plots exactly follow the proposed statistical model, but rather whether the recovered MAD and aMAD angles reasonably satisfy the statistical properties needed to convert them into  $\alpha_{95}$  as proposed in Table 8.

To address this issue we rely on the following procedure. We first assemble several collections of Zijderfeld plots, each collection consisting of plots associated with samples having been analysed in a common way (either by thermal or by alternative field demagnetization) and coming from a common geological source (sediments or lava flows). Samples from a given collection may then be expected to share the same statistical properties, except for the direction of the field they have recorded, which may vary from one sample to the next (reflecting, for instance, the secular variation of the geomagnetic field). From a modelling point of view, we may thus expect all the corresponding Zijderfeld plots to individually be described by statistical models sharing the same  $\delta$ ,  $\sigma_\alpha$  and  $\sigma_\beta$  parameters but not necessarily the same drift direction. If this indeed is the case, we may next proceed and for each Zijderfeld plot build an analogue of the  $\theta_{\text{MAD}}$  rescaled angle investigated in the previous section. Unfortunately, exactly the same  $\theta_{\text{MAD}} = \theta/\text{MAD}$  cannot be computed, since even though both the MAD angle and the recovered direction can be computed, the true field direction remains unknown. We may, however, take advantage of the fact that each Zijderfeld plot can also be analysed using an aPC analysis (being it understood that, for the purpose of the present tests, only primary magnetization components are being considered in each Zijderfeld plot). Then, rather than computing the angle  $\theta$  between the direction recovered from the PC analysis and the (unknown) true direction, we may compute the angle  $\theta'$  between the direction recovered from the PC analysis and that recovered from the aPC analysis, thus using the latter as a proxy for the true direction. Much like the  $\theta$  angle was rescaled into  $\theta_{\text{MAD}}$ ,  $\theta'$  can then further be rescaled into  $\theta'_{\text{MAD}} = \theta'/\text{MAD}$ . Finally, and to make the most of all the information available in each Zijderfeld plot, this procedure can be applied not only to the full plot (i.e. all the points usually analysed to recover a PC direction), but also to subset of points within the plot.

More specifically, the way we proceeded is the following. For each Zijderfeld plot numbered with index  $i$  within a given collection of  $N$  Zijderfeld plots (so that  $1 \leq i \leq N$ ):

(i) We first compute the direction recovered from the aPC analysis using all the  $n_i$  points selected as relevant in the plot; this direction is next used as a proxy for the true field direction in all subsequent steps.

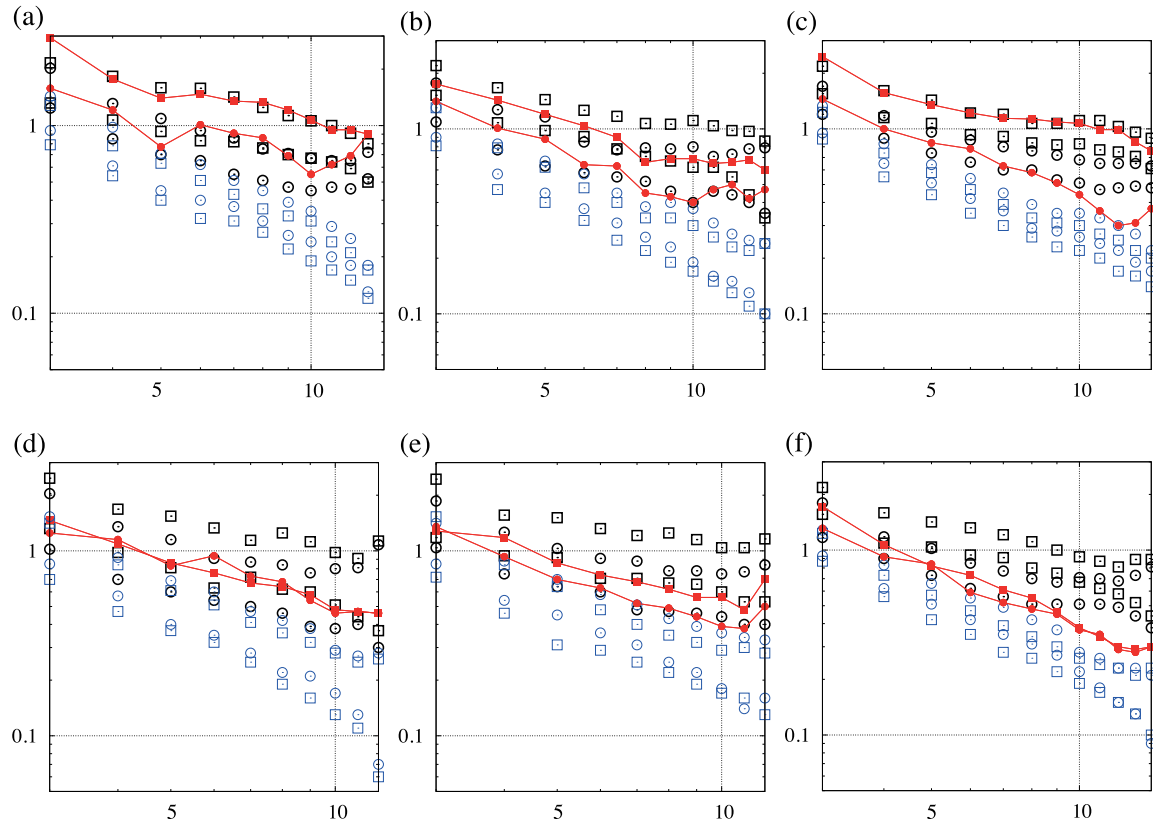
(ii) Next, numbering selected points in the Zijderfeld plot from 1 (closest to the origin) to  $n_i$ , we use the first  $m$  points to recover a MAD angle and a direction from the PC analysis, compute the  $\theta'$  angle between this direction and the proxy aPC direction computed in the previous step, and finally compute the corresponding  $\theta'_{\text{MAD}} = \theta'/\text{MAD}$  rescaled (non-dimensional) angle; this is repeated for all  $m$  points such that  $3 \leq m \leq n_i$ ; for each  $m$  the corresponding  $\theta'_{\text{MAD}}$  angle is then denoted  $\theta'_{\text{MAD,ini}}(m, i)$ .

(iii) Finally, and still using the same numbering of points in the Zijderfeld plot, we repeat the same procedure but use the final (rather than initial)  $m$  points to carry on the PC analysis; for each  $m$  such that  $3 \leq m \leq n_i$ , the corresponding  $\theta'_{\text{MAD}}$  angle is then denoted  $\theta'_{\text{MAD,fin}}(m, i)$ ; note that by construction, for  $m = n_i$ ,  $\theta'_{\text{MAD,fin}}(m, i) = \theta'_{\text{MAD,ini}}(m, i)$ .

Relying on all of the  $N$  plots from a given collection, and for each value of  $m$  (i.e. common number of points used in each Zijderfeld plot to recover the MAD angles), we may then collect all the  $\theta'_{\text{MAD,ini}}(m, i)$  values. Note that the exact number of values collected in this way will be  $N$  for small values of  $m$ , but may be smaller for larger values of  $m$  since samples within a collection do not always have the same number of reliable points in their Zijderfeld plots. For each value of  $m$ , we may then use this population of  $\theta'_{\text{MAD,ini}}(m, i)$  and infer its 0.5 quantile, which we denote  $C'_{\text{ini}}(m)$ . Likewise, we may use the population of  $\theta'_{\text{MAD,fin}}(m, i)$  values and infer its 0.5 quantile, which we denote  $C'_{\text{fin}}(m)$ .

The  $C'_{\text{ini}}(m)$  and  $C'_{\text{fin}}(m)$  0.5 quantile values, as just defined, are then similar to the  $C_{\text{MAD}}$  0.95 quantile values computed in the previous section. They differ, however, in several important ways. First, rather than considering 0.95 quantiles, we only consider 0.5 quantiles. This is for practical reasons. As we shall later see,  $N$ , the total number of samples available in each of the collections we assembled, ranges from 52 to 229. This makes it impossible for us to derive stable 0.95 quantile estimates. But it is enough to derive robust 0.5 quantile estimates. Second,  $C'_{\text{ini}}(m)$  and  $C'_{\text{fin}}(m)$  are derived from a population of  $\theta'_{\text{MAD}}$  normalized angles, whereas  $C_{\text{MAD}}$  is derived from a population of  $\theta_{\text{MAD}}$  normalized angles. This is a subtle difference. In particular, the fact that  $\theta'_{\text{MAD}}$  normalized angles are computed with respect to the proxy aPC direction (and not the true, unknown, direction) introduces a potential dissymmetry between the two  $\theta'_{\text{MAD,ini}}(m, i)$  and  $\theta'_{\text{MAD,fin}}(m, i)$  populations, itself leading to potential systematic differences between the  $C'_{\text{ini}}(m)$  and  $C'_{\text{fin}}(m)$  quantiles. This property can usefully be taken advantage of, as we shall shortly see. Finally, and most importantly, contrary to the  $C_{\text{MAD}}$  0.95 quantile, the two  $C'_{\text{ini}}(m)$  and  $C'_{\text{fin}}(m)$  0.5 quantiles can be computed from both true and synthetic Zijderfeld plots, thus providing an indirect means to assess the relevance of the Zijderfeld plot random walk model we used to compute the  $C_{\text{MAD}}$  conversion factors of Table 8.

The data we relied on consisted in six collections of Zijderfeld plots extracted from palaeomagnetic studies differing in both the nature of the rocks that were analysed and the way these rocks had been analysed. We did not carry out any selection ourselves and only requested the authors of the studies to provide us with samples they considered reliable by the quality criteria used in their



**Figure 4.** Plots of the  $C'_{\text{ini}}(m)$  (red circles) and  $C'_{\text{fin}}(m)$  (red squares) quantiles (see Section 7 for definition and details) as a function of the number  $m$  of points analysed in the Zijdeveld plots, for each of the six collections of Zijdeveld plots used for testing the statistical assumptions used in this study: (a), (b) and (c) correspond to sediment collections A, B and C; (d), (e) and (f) correspond to volcanic collections D, E and F. Shown in black symbols are the envelopes within which the  $C'_{\text{ini}}(m)$  (open circles) and  $C'_{\text{fin}}(m)$  (open squares) quantiles are expected to lie if the Zijdeveld plots were satisfying the random walk model as defined by eqs (1) and (2) with no measurement errors (i.e.  $\sigma_{\beta} = 0$ ). Also shown, in blue symbols, the envelopes within which these quantiles would be expected to lie if the Zijdeveld plots did not follow a random walk but rather a systematic drift with pure measurement errors, as defined by eqs (1) and (2) with  $\sigma_{\alpha} = 0$  (Note the log–log scales).

study. The first collection (to which we will refer to as collection A) corresponds to recent (Holocene) sediment samples from the Laguna Potrok Aike in Argentina, extracted from the study of Lisé-Pronovost *et al.* (2013) (courtesy of A. Lisé-Pronovost). The selected 101 plots were all taken from a core that had been analysed using a 2G cryogenic magnetometer with U-channels, using 13 alternative field demagnetization steps. The second collection (collection B) was also assembled from sediment samples, but from the much older (roughly 1 Ga) Kartochka formation in the East Angara terrane of the Yenisey Ridge region (southwestern Siberian platform, Gallet *et al.* 2012). The 91 selected hand samples (no longer a core) had also been analysed using a 2G cryogenic magnetometer, but now relying on thermal demagnetization steps (between 13 and 18 steps). The third collection (collection C) fairly comparable in nature, could be made even larger (courtesy of V. Pavlov and Y. Gallet). It consisted of 229 hand samples selected from equally old (roughly 1 Ga) sediments coming from the Uchur-Maya region in eastern Siberia (Pavlov & Gallet 2010). Demagnetization had again been carried out in the same way (thermal demagnetization, using between 12 and 18 steps) still relying on the same 2G cryogenic magnetometer. The final three collections were assembled from volcanic samples. Two were extracted from the studies of Tanty *et al.* (2015) and Quidelleur *et al.* (2009) corresponding to volcanic samples of young age (between 2 Ma to recent times) coming from respectively Martinique Island and the back-arc volcanism east of the Andean Cordillera in Argentina (courtesy of J. Carlut). In both

cases, samples were measured using an Agico JR-5 spinner magnetometer and demagnetization was achieved either thermally (with 15–16 temperature steps) or with the help of alternating fields (up to 12 steps). Since we did not witness any significant differences in the behaviour of the samples from the two studies when analysed in the same way, collections were assembled based on the demagnetization method used rather than on the geographical provenance of the samples. Collection D thus consisted of 52 samples (21 from Martinique Island, 31 from Argentina) demagnetized with alternating fields, while collection E consisted of 59 samples (28 from Martinique Island, 31 from Argentina) thermally demagnetized. The final collection (collection F) consisted of substantially older ( $\sim 250$  Ma) volcanic samples selected within the Onkuchaskaya suite of the Pavlov *et al.* (2011) study of the Siberian Permian-Triassic traps (courtesy of V. Pavlov). The corresponding 178 samples were thermally demagnetized (with at least 15 steps) using a 2G cryogenic magnetometer.

Fig. 4 shows plots (in log-log scales) of the  $C'_{\text{ini}}(m)$  (red circles) and  $C'_{\text{fin}}(m)$  (red squares) quantiles computed in the way described above as a function of the number  $m$  of points analysed in the Zijdeveld plots. In each of these plots we also show envelopes within which the  $C'_{\text{ini}}(m)$  (open black circles) and  $C'_{\text{fin}}(m)$  (open black squares) quantiles would be expected to lie if these Zijdeveld plots were satisfying the random walk model we assumed so far (with no measurement errors), as well as the envelopes within which

$C'_{ini}(m)$  (open blue circles) and  $C'_{fin}(m)$  (open blue squares) quantiles would lie if the true Zijdeveld plots would rather follow a statistical model only involving a systematic drift with pure measurement errors, of the type proposed by Kent *et al.* (1983).

In practice, and for each given collection of true Zijdeveld plots (A, B, C, D, E and F), the random walk envelopes were computed by (i) producing 100 random collections of synthetic Zijdeveld plots, each synthetic collection consisting of the same amount  $N$  of Zijdeveld plots as the original collection, with exactly the same number  $n_i$  of points in the  $i$ th (with  $i = 1$  to  $N$ ) plot; (ii) computing the  $C'_{ini}(m)$  and  $C'_{fin}(m)$  quantiles for each of these 100 random collections in the same way as these quantities had been computed for the original collection; and (iii) retaining the smallest (lower envelope) and largest (upper envelope) values of the 100  $C'_{ini}(m)$  and  $C'_{fin}(m)$  computed in this way. The values chosen to produce the random walk, as defined by eqs (1) and (2), were  $\delta = 10$ ,  $\sigma_\alpha = 1$  (hence  $d = \delta/\sigma_\alpha = 10$ ) and  $\sigma_\beta = 0$ , consistent with the set of values already used to produce the  $C_{MAD}$  values in Table 7. As a matter of fact, we also repeated the calculation with  $\delta = 5$ ,  $\sigma_\alpha = 1$  (hence  $d = \delta/\sigma_\alpha = 5$ ) and  $\sigma_\beta = 0$  and also checked that, as was already the case for the  $C_{MAD}$  values, the resulting envelopes did not significantly depend on the value chosen for  $d$ . This is an encouraging intermediary result, as it also implies that the envelopes computed in this way are not expected to depend on the (*a priori* unknown) values of  $\delta$  and  $\sigma_\alpha$ , but only on  $m$  and the known number of samples used to build the  $C'_{ini}(m)$  and  $C'_{fin}(m)$  quantiles.

The procedure to produce the second set of envelopes, those predicted in the event Zijdeveld plots would rather follow a systematic drift with pure measurement errors as proposed by Kent *et al.* (1983), was similar, except for the fact that different  $\delta$ ,  $\sigma_\alpha$  and  $\sigma_\beta$  parameters were used. We set  $\sigma_\alpha = 0$  to remove the random walk component. In contrast, we assumed non-zero values for  $\sigma_\beta$  to introduce pure measurement errors. Kent *et al.* (1983) argued that appropriate values for  $\sigma_\beta$  in eq. (2) should depend on the intensity of the magnetization to be measured, that is, on the modulus  $\mu_k$  of  $\sum_{i=1}^k \mathbf{S}_i$ , which amounts to  $\mu_k = k\delta$  in the present case. More specifically, they argued for a dependency of the type  $\sigma_{\beta k}^2 = a\mu_k^2 + b$  (their eq. 3), where  $a$  and  $b$  are parameters depending on the sample analysed. In practice, three different sets of parameters were tested. In the first case, we assumed  $\delta = 10$ ,  $a = 0.01$  and  $b = 0$ . This amounted to assume the same fundamental step of length  $\delta = 10$  as in the random walk case previously considered (which led to the random walk envelopes plotted in Fig. 4), but a measurement noise of  $\sigma_{\beta k} = 0.1\mu_k = k$ , rather than a random walk component of  $\sigma_\alpha = 1$ . In short, whereas we previously assumed that each step in the Zijdeveld plot involved an uncertainty of relative size  $\sigma_\alpha/\delta = 0.1$ , thus producing a cumulative random walk type of error on the resulting vector  $\mathcal{R}_k$  at step  $k$  with the key implication that the error on  $\mathcal{R}_k$  would be partly correlated to the error in previous resulting vectors  $\mathcal{R}_{k'}$  (with  $k' < k$ ), we here followed Kent *et al.* (1983) and assumed independent measurement errors of relative size  $\sigma_{\beta k}/\mu_k = 0.1$  on each resulting vector  $\mathcal{R}_k$ , thus implying no correlations among errors produced on successive  $\mathcal{R}_k$  vectors. Quite remarkably again, the resulting set of envelopes (shown in Fig. 4) did not appear to significantly depend on the detailed choice made for the value of the relative size  $\sigma_{\beta k}/\mu_k$ . This could be checked by repeating the same calculation with  $\delta = 5$ ,  $a = 0.01$  and  $b = 0$ , resulting in  $\sigma_{\beta k}/\mu_k = 0.2$ . We also tested the possibility of having  $\sigma_{\beta k}$  independent of  $k$ , by setting  $\delta = 10$ ,  $a = 0$  and  $b = 0.01$ , amounting to assume  $\sigma_{\beta k} = 1$  (and  $\sigma_\beta/\delta = 0.1$ ). This again led to only very mild changes in the behaviour of the envelopes. All this led us to the second important conclusion that no matter the

(reasonably realistic) values assumed for  $\delta$ ,  $a$  and  $b$ , and therefore  $\sigma_\beta$ , the envelopes produced under the assumption proposed by Kent *et al.* (1983) always behave as shown in Fig. 4.

Bearing all these results in mind, we may now turn to the interpretation of Fig. 4. First consider the most important case of the sediment samples (Figs 4a–c). As can be seen, the behaviour of the  $C'_{ini}(m)$  and  $C'_{fin}(m)$  quantiles as a function of  $m$  very closely resemble that predicted by the random walk based statistical model of Zijdeveld plots we assumed to derive the results of Table 8. In particular, we note a clear trend of  $C'_{ini}(m)$  quantiles (red circles) to lie below the  $C'_{fin}(m)$  quantiles (red squares). This behaviour is also clearly seen in the envelopes predicted by the random walk model (black open circles and black open squares), and is in fact a direct consequence of the non-reversible nature of random walks. We also note an overall satisfying agreement in the trends of the  $C'_{ini}(m)$  and  $C'_{fin}(m)$  quantiles which most often lie within the predicted envelopes. We do recognize, however, that this is not always the case, testifying for the limits of the very simple random walk based statistical model we assumed. Generally, however, the disagreement is small. It can be quantified in terms of the model correctly predicting the  $C'_{ini}(m)$  and  $C'_{fin}(m)$  quantiles to within 30 per cent error (as can be inferred from the typical shift needed to bring observed quantiles safely within the corresponding predicted envelopes in these log-log diagrams). In contrast, it clearly appears that modelling Zijdeveld plots in terms of a systematic drift with pure measurement errors, as proposed by Kent *et al.* (1983), fails to produce adequate predictions. These predictions lead to no significant dissymmetry between  $C'_{ini}(m)$  and  $C'_{fin}(m)$  quantiles and all envelopes are always significantly too low. From all these considerations we thus conclude that sediment data, from U-channels or individual samples, young or old, demagnetized thermally or by using alternative fields, are all much better described by the random walk model we proposed in this study than by pure measurement errors. This, we note, also has the interesting implication that for sediment data, it thus is the noise in the magnetization acquisition mechanism (responsible for the random walk), rather than the measurement errors, which dominates the signal recovered in Zijdeveld plots. But most importantly, these results also validate our suggestion of relying on the values of Table 8 to rescale MAD and aMAD angles into  $\alpha_{95}$  estimates when dealing with sediment data, with the clear understanding that the resulting  $\alpha_{95}$  estimates may be in error by up to 30 per cent in relative value.

Turning to the case of volcanic samples (Figs 4d–f) reveals a slightly different story. The  $C'_{ini}(m)$  and  $C'_{fin}(m)$  quantiles now essentially behave in the same way, and both tend to lie on the low side of, or even below, the envelopes predicted by the random walk model. Modelling Zijdeveld plots in terms of a systematic drift with pure measurement errors, however, still fails to produce adequate predictions (with envelopes still significantly too low). What this now suggests is that errors in volcanic samples are not as strongly dominated by the noise in the magnetization acquisition mechanism. But it again also suggests that errors cannot entirely be interpreted in terms of measurement errors, either. As a matter of fact, adding adequate measurement errors to our random walk model (i.e. relying on a model with non-zero values for both  $\sigma_\alpha$  and  $\sigma_\beta$  in eqs 1 and 2) can easily lead to the observed behaviour for both  $C'_{ini}(m)$  and  $C'_{fin}(m)$  quantiles. This, however, was an investigation we did not pursue beyond this observation, as it would not have brought much further relevant information for the present study. It simply shows that relying on the values of Table 8 to rescale MAD and aMAD angles into  $\alpha_{95}$  estimates when dealing with volcanic data is not as secure as in the case of sediment data. From the systematic

shift seen in Figs 4(d)–(f) between the actual  $C'_{\text{ini}}(m)$  and  $C'_{\text{fin}}(m)$  quantile values and those predicted by our model, we may nevertheless infer the important conclusion that using the rescaling factors of Table 8 would usually lead to overestimating the predicted  $\alpha_{95}$  estimates by up to 50 per cent (as inferred from the systematic shift needed to bring observed quantiles safely within the corresponding predicted envelopes in these log-log diagrams). Still, given the absolute values of the  $C_{\text{MAD}}$  and  $C_{\text{aMAD}}$  parameters, systematically above 2, it is our contention that using these rescaling factors, rather than just relying on the MAD and aMAD angles as proxies for  $\alpha_{95}$  estimates, can also be of substantial use even for volcanic samples, with the clear understanding that the resulting  $\alpha_{95}$  estimates would then likely overestimate the true  $\alpha_{95}$  by up to 50 per cent in relative value.

## 8 SUMMARY AND ILLUSTRATION

In this paper, we introduced a simple statistical model of the way demagnetization techniques lead to Zijdeveld plots, when these display sequences of aligned points testifying for magnetic carriers having been magnetized within a single common magnetic field. This model statistically relates directions recovered from such plots to the direction of the field assumed to have magnetized the sample under investigation. It can be used to derive analytical results and generate synthetic Zijdeveld plots under controlled statistical conditions.

Analysing large numbers of Zijdeveld plots generated in this way, we empirically showed that classical principal component and anchored principal component analysis (as defined in Section 3) both provide estimates of the palaeomagnetic direction that very closely satisfy a Fisher distribution. This suggests that directions recovered from true palaeomagnetic Zijdeveld plots using either of these two analysis can be assumed to belong to a Fisher distribution and assigned a meaningful  $\alpha_{95}$  estimate.

Principal component and anchored principal component analysis, however, do not directly provide such  $\alpha_{95}$  estimates. Rather, they provide MAD and aMAD angles. Analysing MAD and aMAD angles recovered from large numbers of synthetic Zijdeveld plots, and comparing these with  $\alpha_{95}$  estimates associated with directions recovered from a Fisher analysis of the same synthetic Zijdeveld plots, showed that MAD, aMAD and  $\alpha_{95}$  estimates display different pdfs. These pdfs have slightly different shapes, and do not lead to identical mean, median, and rms values. This unfortunately implies that MAD and aMAD angles do not scale directional errors in the way  $\alpha_{95}$  estimates do. But it also shows that optimal scaling factors can be found to convert MAD and aMAD angles into appropriate  $\alpha_{95}$  estimates, provided one clearly states which property of  $\alpha_{95}$  estimates one wishes the rescaled MAD and aMAD angles to satisfy.

Using guidance from asymptotic expressions of the rms of the MAD, aMAD and  $\alpha_{95}$ , which we analytically derived, and numerous simulations, we showed that such factors could indeed be found. They essentially do not depend on the errors affecting the data, but are a function of the number  $n$  of demagnetization steps involved in the analysis of the Zijdeveld plot.

The scaling factors we recommend for practical use are the  $C_{\text{MAD}}$  and  $C_{\text{aMAD}}$  factors provided in Table 8. These have been optimized to ensure that  $\alpha_{95}$  estimates define the angular limit about the recovered direction within which the true direction should lie 95 per cent of the time. These factors accurately apply only to the extent that the Zijdeveld plots analysed strictly conform to the statistical model we assumed to derive our results. Fortunately, tests with real Zijdeveld

plots showed that for all practical purposes and when considering typical sediment data, these factors can also safely be used to convert MAD and aMAD angles into  $\alpha_{95}$  angles to within an acceptable relative uncertainty of 30 per cent. The situation is slightly different in the case of volcanic data. In that case, real data suggests that using  $C_{\text{MAD}}$  and  $C_{\text{aMAD}}$  factors provided in Table 8 leads to  $\alpha_{95}$  angles that can be overestimated by up to 50 per cent. This is still much better than directly using MAD and aMAD angles as (vastly underestimated, by factors larger than 2) proxies for  $\alpha_{95}$  estimates.

To illustrate how MAD and aMAD angles would translate into useful and mutually consistent  $\alpha_{95}$  estimates using Table 8, we now turn back to the real and synthetic Zijdeveld plots shown in Fig. 1. Each of these plots were analysed using a principal component, an anchored principal component, and a Fisher analysis. The KT68 plot and its synthetic SKT68 counterpart were analysed using  $n = 5$  steps (between  $T = 430^\circ$  and  $T = 560^\circ$ ). The KT75 and SKT75 plots were analysed using  $n = 12$  steps (between  $T = 390^\circ$  and  $T = 680^\circ$ ).

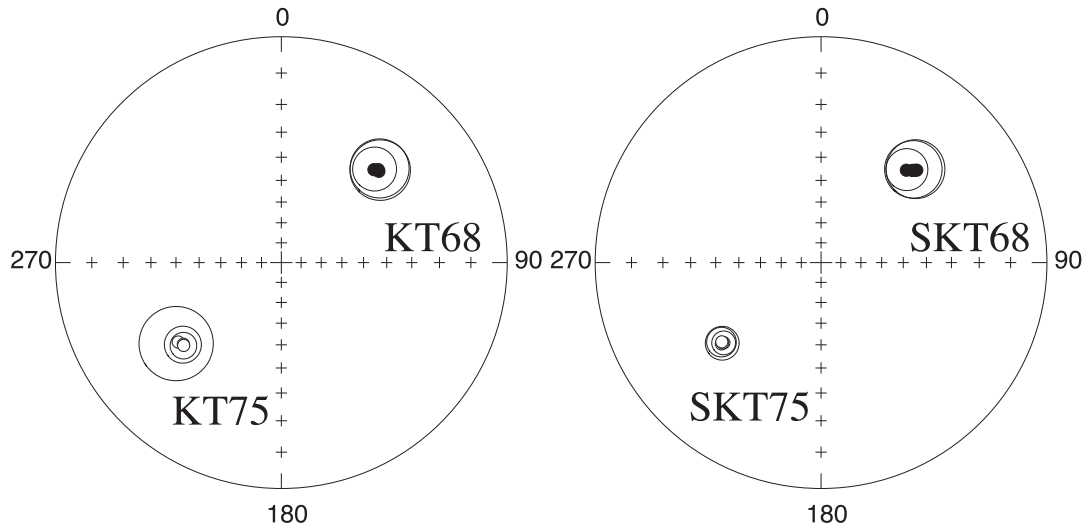
Considering synthetic case SKT68 (Fig. 5) reveals that all recovered directions agree with each other: any given direction is within the  $\alpha_{95}$  confidence circle of any other recovered direction, and all confidence circles include the reference direction. Note also that all  $\alpha_{95}$  estimates are of comparable sizes, with differences no larger than one would naturally expect from the width of the various pdfs shown in Fig. 3.

Considering synthetic case SKT75 illustrates a case with an identical level of noise in the Zijdeveld plot, but analysed using significantly more demagnetization steps ( $n = 12$  rather than  $n = 5$ ). Recovered directions are now more tightly grouped, as one would expect, and the  $\alpha_{95}$  confidence circles are of much smaller size. Again all directions fit within the  $\alpha_{95}$  confidence circle of any other recovered direction, and all confidence circles include the reference direction.

Turning to the real sediment sample KT68, with identical level of noise and analysed with the same number of demagnetization steps as SKT68, reveals a very similar situation (Fig. 5). Again, all directions and associated  $\alpha_{95}$  estimates are mutually consistent.

KT75 reveals a different but interesting situation. A quick comparison between Figs 1(b) and (d) reveals that demagnetization steps are more irregular in KT75 than in SKT75, despite the fact that SKT75 was generated using a random walk best matching KT75. Interestingly, however, all recovered directions still agree with each other in KT75, and both the principal component and anchored principal component analysis lead to  $\alpha_{95}$  estimates consistent with each other and comparable with those we had found for SKT75. The most significant difference turns out to arise from the Fisher analysis, which provides a direction remaining within the confidence circles of the two other directions, but leads to a much larger  $\alpha_{95}$  estimate. The reason for this behaviour is related to the smallest demagnetization steps in the KT75 Zijdeveld plot between  $T = 560^\circ$  and  $T = 660^\circ$  (Fig. 1b). Such steps are outliers with respect to the statistical model we assumed in this study. They are typical of small steps taken when approaching a Curie temperature (here that of Magnetite). In the present instance, they arrange themselves in a triangular shaped sequence in the Zijdeveld plot, which produces strong changes in the directions of the unit vectors  $\mathbf{s}_i$  used in the Fisher analysis of such plots (recall eq. 6). This affects both the recovered direction and its associated  $\alpha_{95}$  estimate. The impact, however, is weaker on the recovered direction, because two successive steps tend to compensate their effect on the recovered direction. It also is very weak on the directions recovered from principal component and anchored principal component analysis and on the values





**Figure 5.** Stereographic projections of the palaeomagnetic directions recovered from the Zijderveld plots shown in Fig. 1. Plots on the left-hand side correspond to the real KT68 and KT75 samples, and show the three (overlapping) directions recovered from a principal component, an anchored principal component and a Fisher analysis of the corresponding Zijderveld plot. They also show the  $\alpha_{95}$  confidence circles as inferred from the MAD and aMAD using the recommended conversion Table 8 for the principal component and anchored principal component directions, and the usual  $\alpha_{95}$ , as inferred from eq. (11), for the Fisher direction. Corresponding sets of  $(D, I, \alpha_{95})$  values are  $(46.9^\circ, 28.8^\circ, 11.3^\circ)$ ,  $(44.9^\circ, 29.7^\circ, 8.3^\circ)$ ,  $(46.2^\circ, 28.5^\circ, 11.0^\circ)$  for KT68 and  $(230.2^\circ, -31.1^\circ, 7.1^\circ)$ ,  $(229.7^\circ, -30.9^\circ, 5.1^\circ)$ ,  $(232.4^\circ, -30.1^\circ, 14.0^\circ)$  for KT75. The analogous plots on the right-hand side correspond to the synthetic SKT68 and SKT75 cases. The corresponding sets of  $(D, I, \alpha_{95})$  values are  $(44.7^\circ, 30.3^\circ, 10.9^\circ)$ ,  $(42.7^\circ, 31.8^\circ, 8.1^\circ)$ ,  $(45.3^\circ, 29.7^\circ, 11.1^\circ)$  for SKT68 and  $(230.8^\circ, -31.3^\circ, 6.5^\circ)$ ,  $(231.4^\circ, -31.3^\circ, 5.7^\circ)$ ,  $(230.5^\circ, -31.7^\circ, 4.6^\circ)$  for SKT75; also plotted (but overlapping) in that case are the reference directions used to generate the data  $(D, I) = (46.0^\circ, 29.0^\circ)$  for SKT68 and  $(230.8^\circ, -30.7^\circ)$  for SKT75.

of the MAD and aMAD angles, because of the small size of these steps. The robustness of principal component and anchored principal component analysis with respect to such small steps is consistent with the results of the tests carried out in Section 7 and the very reason for their wide success. What the present study now shows is that using the  $C_{\text{MAD}}$  and  $C_{\text{aMAD}}$  factors provided in Table 8 finally makes it possible to combine the robustness of principal component and anchored principal component analysis with the possibility of simultaneously recovering  $\alpha_{95}$  estimates.

## 9 PRACTICAL RECOMMENDATIONS

We finally conclude with some practical recommendations to take advantage of our findings and produce  $\alpha_{95}$  estimates and 95 per cent confidence intervals on the declination and inclination recovered from single samples displaying Zijderveld plots with linear behaviours (possibly in the form of two successive directions; note, however, that this study did not dwell with Zijderveld plots that show curved or even more complex behaviours). The recipe for this is very simple.

Consider the case when a principal component analysis has been used to analyse the sample, providing a direction (declination  $D$  and inclination  $I$ ) and a MAD value (as defined in Section 3, where notations are also defined). If  $n$  is the number of demagnetization steps that have been used, then an  $\alpha_{95}$  estimate can be computed just using

$$\alpha_{95} = C_{\text{MAD}}(n) \times \text{MAD}, \quad (44)$$

where  $C_{\text{MAD}}(n)$  is to be extracted from Table 8. This  $\alpha_{95}$  estimate can next also be used to infer the corresponding declination ( $\Delta D_{95}$ ) and inclination ( $\Delta I_{95}$ ) 95 per cent confidence intervals, using the well-known formulae (see eqs A.44 and A.45 in Butler 1992)

$$\sin \Delta D_{95} = \frac{\sin \alpha_{95}}{\cos I}, \quad \Delta I_{95} = \alpha_{95}. \quad (45)$$

If, rather than the regular principal component analysis just discussed, an anchored principal component analysis has been used, the procedure is exactly the same, except for the fact that the  $\alpha_{95}$  estimate must then be computed using the aMAD value and

$$\alpha_{95} = C_{\text{aMAD}}(n) \times \text{aMAD}, \quad (46)$$

where  $C_{\text{aMAD}}(n)$ , also to be extracted from Table 8, is a scaling parameter different from  $C_{\text{MAD}}(n)$ . Note again that it thus is very important to clearly distinguish whether the sample was analysed using a standard principal component analysis or an anchored standard principal component analysis.

Tests carried out in Section 7 showed that the above procedure applied to sediment data then leads to  $\alpha_{95}$  estimates and 95 per cent confidence intervals on declination and inclination accurate to within an estimated 30 per cent relative error. In the case of volcanic data, these values could be overestimated by up to 50 per cent. This, we note, is still much better than just using MAD and aMAD angles as proxies for  $\alpha_{95}$  estimates. Besides, and as already noted, volcanic directional data are almost never recovered from single samples.

Most interesting is the case of magnetostratigraphic studies, for which no other practical way currently exist to routinely produce  $\alpha_{95}$  estimates and 95 per cent confidence intervals on declination and inclination time series recovered from sedimentary long sequences. The above procedure can very straightforwardly be used for that purpose.

Such a recipe, we finally note, would be a trivial matter to implement as a subroutine in software which already rely on principal component analysis to reconstruct paleodirections from demagnetization measurements (e.g. Cogné 2003; Mazaud 2005; Tauxe *et al.* 2014, or the so-called PMGSC software developed by R. Enkin at the Geological Survey of Canada). We encourage authors of such software to look into this possibility.



## ACKNOWLEDGEMENTS

The authors are much indebted to J. Kirschvink and two anonymous reviewer for their extremely useful comments on the original version of this manuscript and to A. Lisé-Pronovost, V. Pavlov, Y. Gallet and J. Carlut for providing the data analysed in Section 7. They also wish to thank G.M. Molchan for his support and very useful suggestions with respect to the analytical derivation used in Section 4, as well as J.P. Cogné and J.P. Valet for fruitful discussions, and Y. Gallet for assistance in producing Figs 1 and 5 (using the PaleoMac software of Cogné 2003) and constructive comments on early versions of the manuscript. This work was partly financed by grant N 14.Z50.31.0017 of the Russian Ministry of Science and Education. This is IGP contribution N°3690.

## REFERENCES

- Bingham, C., 1983. A series expansion for angular Gaussian distribution, in *Statistics on Spheres*, pp. 226–231, ed. Watson, C., Wiley-Interscience.
- Butler, R.F., 1992. *Paleomagnetism: Magnetic Domains to Geological Terranes*, Blackwell Scientific Publications.
- Cogné, J.-P., 2003. PaleoMac: A Macintosh (TM) application for treating palaeomagnetic data and making plate reconstructions, *Geochem. Geophys. Geosyst.*, **4**, 1007, doi:10.1029/2001GC000227.
- Dunlop, D.J. & Özdemir, O., 2007. Magnetization in rocks and minerals, in *Geomagnetism*, Vol. 5 of Treatise on Geophysics, eds Kono, M. & Schubert, G., Elsevier.
- Fisher, R., 1953. Dispersion on a sphere, *Proc. R. Soc. A*, **217**, 295–305.
- Gallet, Y., Pavlov, V., Halverson, G. & Hulot, G., 2012. Toward constraining the long-term reversing behavior of the geodynamo: a new “Maya” superchron  $\sim 1$  billion years ago from the magnetostratigraphy of the Kartočka Formation (southwestern Siberia), *Earth planet. Sci. Lett.*, **339**, doi:10.1016/j.epsl.2012.04.049.
- Hill, M.J., Gratton, M.N. & Shaw J. 2002. A comparison of thermal and microwave palaeomagnetic techniques using lava containing laboratory induced remanence, *Geophys. J. Int.*, **151**, 157–163.
- Hongre, L., Hulot, G. & Khokhlov, A., 1998. An analysis of the geomagnetic field over the past 2000 years, *Phys. Earth planet. Inter.*, **106**, 311–335.
- Hulot, G. & Le Mouél, J.L., 2009. A statistical approach to the Earth’s main magnetic field, *Phys. Earth planet. Inter.*, **82**, 167–183.
- Johnson, C.L. *et al.*, 2008. Recent investigations of the 0–5 Ma geomagnetic field recorded by lava flows, *Geochem. Geophys. Geosyst.*, **9**(4), Q04032, doi:10.1029/2007GC001696.
- Kent, J.T., Briden, J.C. & Mardia, K.V., 1983. Linear and planar structure in ordered multivariate data as applied to progressive demagnetization of palaeomagnetic remanence, *Geophys. J. R. astr. Soc.*, **75**, 593–621.
- Khokhlov, A. & Hulot, G., 2013. Probability uniformization and application to statistical palaeomagnetic field models and directional data, *Geophys. J. Int.*, **193**, 110–121.
- Khokhlov, A., Hulot, G. & Bouligand, C., 2006. Testing statistical palaeomagnetic field models against directional data affected by measurement errors, *Geophys. J. Int.*, **167**(2), 635–648.
- Kirschvink, J.L., 1980. The least-squares line and plane and the analysis of palaeomagnetic data, *Geophys. J. R. astr. Soc.*, **62**, 699–718.
- Korte, M., Genevey, A., Constable, C.G., Frank, U. & Schnepf, E., 2005. Continuous geomagnetic field models for the past 7 millennia: 1. A new global data compilation, *Geochem. Geophys. Geosyst.*, **6**, Q02H15, doi:10.1029/2004GC000800.
- Korte, M., Donadini, F. & Constable, C.G., 2009. Geomagnetic field for 0–3 ka: 2. A new series of time-varying global models, *Geochem. Geophys. Geosyst.*, **10**, Q06008, doi:10.1029/2008GC002297.
- Laj, C., Kissel, C. & Roberts, A.P., 2006. Geomagnetic field behavior during the Iceland Basin and Laschamp geomagnetic excursions: a simple transitional field geometry?, *Geochem. Geophys. Geosyst.*, **7**, Q03004, doi:10.1029/2005GC001122.
- Leonhardt, R. & Fabian, K., 2007. Geomagnetic field for 0–3 ka: 2. A new series of time-varying global models, *Earth planet. Sci. Lett.*, **253**, 172–195.
- Lhuillier, F., Fournier, A., Hulot, G. & Aubert, J., 2011. The geomagnetic secular-variation timescale in observations and numerical dynamo models, *Geophys. Res. Lett.*, **38**, L09306, doi:10.1029/2011GL047356.
- Lisé-Pronovost, A. *et al.*, 2013. High-resolution palaeomagnetic secular variations and relative paleointensity since the Late Pleistocene in southern South America, *Quat. Sci. Rev.*, **71**, 91–108.
- Lund, S.P., Schwartz, M., Keigwin, L. & Johnson, T., 2005. Deep-sea sediment records of the Laschamp geomagnetic field excursion ( $\sim 41,000$  calendar years before present), *J. geophys. Res.*, **110**, B04101, doi:10.1029/2003JB002943.
- Mazaud, A., 2005. User-friendly software for vector analysis of the magnetization of long sediment cores, *Geochem. Geophys. Geosyst.*, **6**, Q12006, doi:10.1029/2005GC001036.
- McFadden, P.L., 1980. The best estimate of Fisher’s precision parameter  $\kappa$ , *Geophys. J. R. astr. Soc.*, **60**, 397–407.
- McFadden, P.L. & McElhinny, M.W., 1990. Classification of the reversal test in palaeomagnetism, *Geophys. J. R. astr. Soc.*, **103**, 725–729.
- Merrill, R.T., McElhinny, M.W. & McFadden, P., 1996. *The Magnetic Field of the Earth*, Academic Press.
- Nagy, E.N. & Valet, J.P., 1993. New advances for palaeomagnetic studies of sediment cores using U-channels, *Geophys. Res. Lett.*, **20**, 671–674.
- Pavlov, V. & Gallet, Y., 2010. Variations in geomagnetic reversal frequency during the Earths middle age, *Geochem. Geophys. Geosyst.*, **11**, Q01Z10, doi:10.1029/2009GC002583.
- Pavlov, V., Fluteau, F., Veselovskiy, R., Fetisova, A. & Latyshev, A., 2011. Secular geomagnetic variations and volcanic pulses in the Permian-Triassic traps of the Norilsk and Maimecha-Kotui Provinces, *Izv. Phys. Solid Earth*, **47**, 402–417.
- Quidelleur, X., Carlut, J., Tchilinguirian, P., Germa, A. & Gillot, P.Y., 2009. Palaeosecular directions from mid-latitude sites in the southern hemisphere (Argentina): contribution to time averaged field models, *Phys. Earth planet. Inter.*, **172**, 199–209.
- Scheidegger, A.E., 1965. On the statistics of the orientation of bedding planes, grain axes, and similar sedimentological data, *U.S. Geo. Surv. Prof. Pap.*, **525-C**, 164–167.
- Tantý, C., Carlut, J., Valet, J.P. & Germa, A., 2015. Palaeosecular variation recorded by 9 ka to 2.5-Ma-old lavas from Martinique Island: new evidence for the La Palma aborted reversal  $\sim 617$  ka ago, *Geophys. J. Int.*, **200**, 917–934.
- Tauxe, L., LaBreque, J.L., Dodson, R., Fuller, M. & Dematteo, J., 1983. “U” Channels—a new technique for palaeomagnetic analysis of hydraulic piston cores, *EOS, Trans. Am. geophys. Un.*, **64**, 219.
- Tauxe, L., Banerjee, S.K., Butler, R.F. & van der Voo, R., 2014. *Essentials of Paleomagnetism*, 3rd web edn, Available at: <http://earthref.org/MAGIC/books/Tauxe/Essentials/>, last accessed 3 November 2015.
- Zijderveld, J.D.A., 1967. A.C. demagnetization in rocks: analysis of results, in *Methods in Paleomagnetism*, pp. 254–286, eds Collinson, D.W., Creer, K.M. & Runcorn, S.K., Elsevier.