



HAL
open science

Unsupervised learning of seismic wavefield features: clustering continuous array seismic data during the 2009 L'Aquila earthquake

Peidong Shi, Léonard Seydoux, Piero Poli

► **To cite this version:**

Peidong Shi, Léonard Seydoux, Piero Poli. Unsupervised learning of seismic wavefield features: clustering continuous array seismic data during the 2009 L'Aquila earthquake. *Journal of Geophysical Research: Solid Earth*, 2020, 10.1029/2020JB020506 . insu-03094514

HAL Id: insu-03094514

<https://insu.hal.science/insu-03094514v1>

Submitted on 4 Jan 2021

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.

1 **Unsupervised learning of seismic wavefield features:**
2 **clustering continuous array seismic data during the**
3 **2009 L'Aquila earthquake**

4 **Peidong Shi¹, Léonard Seydoux¹, and Piero Poli¹**

5 ¹Institut de Sciences de la Terre, Université Grenoble Alpes, CNRS (UMR5275), Grenoble, France.

6 **Key Points:**

- 7 • Identification of frequency-dependent wavefield features from array analysis of con-
8 tinuous seismic data
9 • Wavefield features reveal a time and frequency evolution of seismic wavefield re-
10 lated to the seismic source properties and fault states
11 • Unsupervised learning of wavefield features identifies distinct clusters well corre-
12 lated with different periods of seismic cycle without explicit physical modeling

Corresponding author: Peidong Shi, peidong.shi@univ-grenoble-alpes.fr

Abstract

We apply unsupervised machine learning to three years of continuous seismic data to unravel the evolution of seismic wavefield properties in the period of the 2009 L'Aquila earthquake. To obtain sensible representations of the wavefield properties variations, we extract wavefield features (i.e. entropy, coherency, eigenvalue variance and first eigenvalue) from the covariance matrix analysis of the continuous wavefield data. The defined wavefield features are insensitive to site-dependent local noise, and inform the spatiotemporal properties of seismic waves generated by sources inside the array. We perform a sensitivity analysis of these wavefield features, and track the evolution of source properties from the unsupervised learning of the uncorrelated features. By clustering the wavefield features, our unsupervised analysis avoids explicit physical modeling (e.g. no requirement for event location and magnitude estimation) and can naturally separate peculiar patterns solely from continuous seismic data. Our model-free unsupervised learning of wavefield features reveals distinct clusters well correlated with different periods of the seismic cycle, which are consistent with previous model-dependent studies.

1 Introduction

Seismological observations are a primary source of information about fault physics and its evolution in time and space (Gutenberg & Richter, 1956; Scholz, 2002; Aki & Richards, 2002). Seismic catalogs are nowadays the main way of labeling seismic data, by association of waveforms with earthquakes occurring in a given position and at a certain time (Gutenberg & Richter, 1956; Scholz, 2002; Aki & Richards, 2002). While earthquake catalogs are among the main source of information to study faults, the continuous stream of seismic data is likely to hide important additional information about fault physics, which cannot be easily summarized into discrete observables. For example, the slow earthquakes and tremors show very different wavefield properties compared to that of regular earthquakes, requiring alternative approaches to derive information about their physics (Ide et al., 2007; Beroza & Ide, 2011). Therefore, it is worthwhile to explore the potential to assess physical properties of faults from direct analysis of continuous seismic wavefields.

The latter idea has been recently explored in laboratory-scale fracture experiments. Indeed, recent studies based on laboratory observations, show that continuous acoustic emission (AE) contain essential information about the physical state of the rock (Rouet-Leduc et al., 2017; Bolton et al., 2019; Hulbert et al., 2019). In these studies, statistical features of the continuous AE signals (e.g. amplitudes, variance etc.), are used for supervised or unsupervised machine learning (ML) and classification, to characterize the wavefield variations and study the evolution of the (laboratory) seismic cycle (Bolton et al., 2019), including the estimation of failure time (Rouet-Leduc et al., 2017). The experiments carried at a laboratory scale already involve complex, nonlinear relationships between the continuous signal properties and the fault states, suggesting that systems of higher complexity such as the real geological settings should also be investigated with machine learning tools, as in the present study.

In addition to the laboratory studies, unsupervised machine learning has been applied to real continuous seismic data in volcanic settings to classify volcanic tremors and monitor volcanic activities (Langer et al., 2009; Esposito et al., 2008; Köhler et al., 2010; Langer et al., 2011; Carniel et al., 2013; Unglert et al., 2016). Unsupervised machine learning can distinguish seismic wavefield of distinct characteristics (e.g. spectral content) generated by different volcanic activities, such as pre-, co- and post-eruption, thus permits the recognition of different types of volcanic activities directly from continuous seismic data.

In summary, both laboratory experiments (Rouet-Leduc et al., 2017; Bolton et al., 2019; Hulbert et al., 2019; Shreedharan et al., 2020) and real volcanic seismic data analysis (Esposito et al., 2008; Köhler et al., 2010; Langer et al., 2011; Carniel et al., 2013;

Unglert et al., 2016) show promising potential to utilize real continuous seismic wavefield and ML algorithms to understand physical processes occurring inside the Earth. However, to our knowledge, no studies have been performed so far on clustering of long-term real continuous array seismic data to establish the space-time evolution of the physical state of the faults where significant earthquakes occur.

We here present an unsupervised class-membership identification (clustering) of ensemble wavefield features, which capture the nature of the seismic wavefield as seen by an array of stations. The choice of array features is aimed at reducing the sensitivity of single-station statistical features to noise intensity (e.g. daily/weekly variation of human activity and variation of meteorological conditions, Cara et al., 2003; Poli et al., 2020) and enhancing the identification of spatio-temporal properties of (possibly mixed) seismic sources (Seydoux et al., 2016a; Soubestre et al., 2019). We can thus recognize patterns within seismic signals and track their temporal evolution, which can be related to particular fault states occurring at different stages of the seismic cycle (e.g. earthquake nucleation, afterslip etc.). Differently from laboratory experiments (Rouet-Leduc et al., 2017; Hulbert et al., 2019; Shreedharan et al., 2020), we have no independent information about the fault state (e.g. stress, friction). That is why we use unsupervised analysis and self-learn from the continuous data.

To test our approach, we used three years of vertical-component seismic data recorded in the region of L’Aquila, Italy (Figure 1). We use this region as a test case, as it hosting a magnitude 6 earthquake (6 of April, 2009, Chiarabba et al., 2009; Di Luccio et al., 2010) preceded by a long-lasting preparatory phase (Sugan et al., 2014; Vuan et al., 2018). Previous studies also reported that the fault properties may have changed dramatically in the preparatory phase of the main event due to fluid movement (Di Luccio et al., 2010; Chiarabba et al., 2020), velocity change (Baccheschi et al., 2020), and variation of elastic and anisotropic parameters (Lucente et al., 2010). In addition, this region is well instrumented with permanent seismic stations (Figure 1a), allowing an array-based analysis. The complex faulting processes and high quality continuous seismic data make the L’Aquila earthquake a perfect test case to investigate the feasibility of tracking fault states directly from continuous seismic wavefield.

We explore spatial wavefield features of long time windows (60 days) and their temporal evolution with respect to the main earthquake in the area using cluster analysis. We highlight different patterns in the wavefield and relate them to the physical processes of the fault (e.g. the preparation, afterslip etc.). Our results show the feasibility of using array-based wavefield properties to directly assess the fault state and characterize different stages of the seismic cycle.

2 Data and Processing

We focus on a time period of about 3 years (2008-2010, included) around the Mw 6.1 L’Aquila earthquake (6 April 2009, Chiarabba et al., 2009; Di Luccio et al., 2010). This event has been chosen because it presented a prominent and long-lasting preparation period, starting 3-4 months before the mainshock, and including several dozens of foreshocks and possible significant changes in the fault rock properties (Di Luccio et al., 2010; Lucente et al., 2010; Di Stefano et al., 2011; Herrmann et al., 2011; Sugan et al., 2014; Vuan et al., 2018; Chiarabba et al., 2020).

The three years of continuous vertical-component seismic data (from 2007-11-03 to 2010-08-23) recorded by the six nearest stations (Figure 1a) at a 50 Hz sampling rate were downloaded from the *Istituto Nazionale Geofisica e Vulcanologia* (INGV) data center (INGV Seismological Data Centre, 2006). Data have been transformed into velocity using the instrument response and processed to remove gaps and glitches. Data gaps and glitches are filled or replaced with white random noise of minimal amplitudes (~ 10

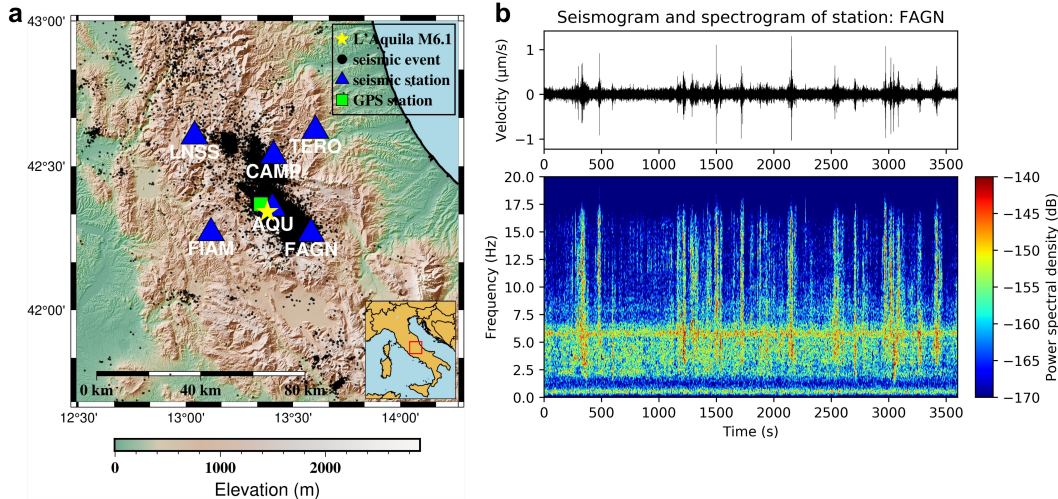


Figure 1. (a) Location of the 2009 L'Aquila earthquake and the nearby permanent seismic array. Yellow star indicates the epicenter of the mainshock. Blue triangles represent the seismic stations. Green square denotes the GPS (Global Positioning System) station. Black dots show the locations of earthquakes including the foreshocks and aftershocks of the 2009 L'Aquila earthquake from 2008-2010 in this region (seismic catalog from INGV). Red rectangular in the bottom-right inserted regional map highlights the current study area. (b) Hour-long example of vertical ground velocity and corresponding spectrogram recorded at the station FAGN. The records start at 2009-04-05 10:00:00 (UTC).

114 orders of magnitude lower than the average signal amplitudes) to allow a continuous analysis of seismic data and eliminate data anomalies. We have tested that this random noise
 115 is not affecting our subsequent analysis. Spectral analysis of the continuous data shows
 116 the dominant frequency range of the local earthquakes is around 0.5-18 Hz (Figure 1b),
 117 while below 0.5 Hz, micro-seismic noise dominates. We thus focus on the frequency range
 118 of 0.5-18 Hz to reduce the effects of ambient noise and also influence of regional or re-
 119 mote earthquakes.
 120

121 3 Decomposition of the Wavefield and Features Extraction

122 3.1 Covariance Matrix Analysis of Continuous Seismic Data

123 We define a set of features relevant for characterizing the propagation of seismic
 124 waves beneath the seismic array, over a broad frequency range (0.5-18 Hz). Following
 125 Seydoux et al. (2016a), we extract these features from the factorization of the covariance
 126 matrix of continuous array seismic data. Such analyses were successfully used for detect-
 127 ing and classifying seismovolcanic tremors (Soubestre et al., 2018), teleseismic earthquakes
 128 (Seydoux et al., 2016a), and for analyzing ambient noise wavefield (Seydoux et al., 2016b).

129 The covariance matrix is built from the time average of the Fourier cross-spectra
 130 matrices calculated over a set of half-overlapping sub-windows (Seydoux et al., 2016a,
 131 and Figure 2a). Two types of time window are involved in the calculation of covariance
 132 matrix. A short one (sub-window, W_1) where the Fourier cross-spectra matrix is calcu-
 133 lated, and a longer one (averaging-window, W_2) used to average the cross-spectra ma-
 134 trices, which in turn defines the covariance matrix of a particular time scale (Figure 2a).
 135 We here use a backward-looking approach to time stamp the results: the end time of each
 136 averaging time window (W_2) is assigned as the time stamp associated with the covari-

137 ance matrix of the time window, hence the obtained results are causal. The size of W1
 138 depends on the size of seismic array and the frequency range of interest (Seydoux et al.,
 139 2016a). In this study, we use a W1 of 80 seconds to ensure the slowest waves to fully travel
 140 the aperture of the seismic array.

141 The size of W2 is crucial to define the time resolution of our analysis. We here aim
 142 at classifying long-lasting patterns in the seismic signals, and thus we average the co-
 143 variance matrix over 60 days and shift W2 by one day. The use of long averaging win-
 144 dow would probably increase the influence of external wavefield properties originated out-
 145 side the seismic array. However, as we perform our analysis at relatively high frequen-
 146 cies (0.5-18 Hz), the analysis inherently focuses on a local area (i.e. inside the array) due
 147 to the attenuation of high-frequency waves generated from distant sources. Because we
 148 want to analyze seismic sources seen by the ensemble of seismic stations, we apply spec-
 149 tral whitening to the daily seismograms before computing the covariance matrix (Seydoux
 150 et al., 2016a, 2016b). In this way, the spectral energy is not taken into account and the
 151 analysis mostly relies on the phase coherence between the seismic stations, thus cancelling
 152 non-propagative signals (e.g. local noise, traffic, wind).

153 3.2 Wavefield Features

154 From the eigendecomposition of the covariance matrix, the eigenvalues $\lambda(f, t)$ and
 155 corresponding eigenvectors $\mathbf{v}(f, t)$ are obtained for each frequency f and time t (Figure
 156 2a). Note that the covariance matrix is inherently Hermitian and positive semi-definite;
 157 the matrix is therefore always diagonalizable and the eigenvalues are positive and real.
 158 From the eigenvalues, we define four features: (1) the Shannon entropy, (2) the coherency,
 159 (3) the eigenvalue variance, and (4) the first eigenvalue.

160 1. The Shannon entropy, initially developed in the frame of information theory (Shannon,
 161 1948) and applied to the case of discrete operators by Von Neumann (1986), provides
 162 a measurement of the quantity of information present in a multivariate dataset. If we
 163 consider the normalized covariance matrix eigenvalues $p_i(f, t) = \lambda_i(f, t) / \sum_{i=1}^N \lambda_i(f, t)$
 164 such as $\sum_{i=1}^N p_i(f, t) = 1$ (where N is the total number of stations in the array and p_i
 165 represents the normalized i -th eigenvalue of the covariance matrix at a given time and
 166 frequency), we can consider each normalized eigenvalue (p_i) to represent the probabili-
 167 ty of each source (identified by each corresponding eigenvector) to be observed in the
 168 studied time period. The Shannon entropy σ_e is then defined as:

$$169 \quad \sigma_e(f, t) = - \sum_{i=1}^N p_i(f, t) \ln(p_i(f, t)). \quad (1)$$

170 Following Shannon (1948), the higher the entropy, the more chaotic the wavefield and
 171 the lower the wavefield spatial coherence. A coherent wavefield generated by only one
 172 source or many co-located sources in the analyzed time window is likely to be spanned
 173 by a single dominating eigenvalue (Figure 2b). Therefore, low values of the entropy will
 174 be observed when the wavefield is dominated by the coherent sources localized in space.
 175

176 2. The coherency function, commonly used in exploration geophysics (Gersztenkorn
 177 & Marfurt, 1999), is defined as the ratio between dominating wavefield component (first
 178 eigenvalue) and the full wavefield (sum of all eigenvalues), and reports the wavefield co-
 179 herence σ_c :

$$180 \quad \sigma_c(f, t) = \frac{\lambda_1(f, t)}{\sum_{i=1}^N \lambda_i(f, t)}. \quad (2)$$

181

182 3. To estimate the flatness of covariance matrix eigenvalues distribution, we define
 183 the eigenvalue variance σ_v as:

$$184 \quad \sigma_v(f, t) = \frac{\sum_{i=1}^N (\lambda_i(f, t) - \mu)^2}{N}, \quad (3)$$

185 where $\mu = \sum_{i=1}^N \lambda_i(f, t)/N$ is the mean eigenvalue at a given time and frequency. The
 186 eigenvalue variance is related to both wavefield coherence and source energy (Figures 2b-
 187 2d). For example, for one dominating source in the studied time window (W2), the cor-
 188 responding eigenvalue variance will be large and the wavefield is coherent as well.

189 4. Finally, we use the first eigenvalue σ_f :

$$190 \quad \sigma_f(f, t) = \lambda_1(f, t). \quad (4)$$

191 Theoretically this value defines the coherence of a single source over the time window
 192 W2. As it is resulting from phase multiplication, this value can be affected by noise, for
 193 example biasing the estimation of the phase correlation. There is thus an imprint of the
 194 frequency dependent signal-to-noise level in this measure. For example, stronger source
 195 and/or a large number of co-located coherent sources in the studied time window (W2)
 196 will result in larger phase correlations (because of higher signal-to-noise ratio after av-
 197 eraging) and thus lead to a larger eigenvalue. Therefore, the first eigenvalue provides a
 198 measurement of the strength of the dominating source in the wavefield.

199 These four features are obtained at each time step (1 day) and frequency (from 0.5
 200 to 18 Hz). We thus have a time-frequency representation of the wavefield (Figure 2a),
 201 which can be used to track its evolution. As mentioned above, the features contain in-
 202 sights about the wavefield spatio-temporal properties, and thus provide insights on the
 203 seismic signals generated inside the array. Since the wavefield features are calculated us-
 204 ing a long window (60 days), many seismic sources can exist in the same time window
 205 of analysis. Among the different potential scenarios, we can distinguish the following ex-
 206 treme cases.

207 If many seismic sources occur in a small region with respect to the wavelength and
 208 the array aperture (e.g. an earthquake swarm or co-located sources), the average covari-
 209 ance matrix will exhibit a dominant eigenvalue while the other eigenvalues will be small
 210 (scenario illustrated in Figure 2b), giving small values for the entropy and high values
 211 for the coherency, eigenvalue variance and first eigenvalue.

212 If many independent seismic sources are acting in the same time window (W2) and
 213 scattered in a vast area with respect to the array aperture, the eigenvalue distribution
 214 will follow a steadily decaying distribution (scenario illustrated in Figure 2c) specific to
 215 the array geometry, the structure of the underlying medium and the duration of the av-
 216 eraging window W2 (Seydoux et al., 2016a). In this situation, the entropy and first eigen-
 217 value will be high and the coherency and eigenvalue variance will be small, indicating
 218 an incoherent ensemble wavefield with many incoherent seismic sources in the analyzed
 219 time scale (W2).

220 Finally, if the records only contain electronic noise or spatially distributed incoher-
 221 ent perturbations (e.g. rain, wind, road traffic etc.), the covariance matrix eigenvalues
 222 will be approximately equal and small (scenario illustrated in Figure 2d) depending on
 223 the estimation parameters (Menon et al., 2014). In this situation, the entropy will be
 224 high and the coherency, eigenvalue variance and first eigenvalue will be small, indicat-
 225 ing an incoherent ensemble wavefield with no sources in the analyzed time scale (W2).

226 In summary, the defined wavefield features permit to discern the behavior of the
 227 wavefield over different frequencies and as a function of time. We use these features to

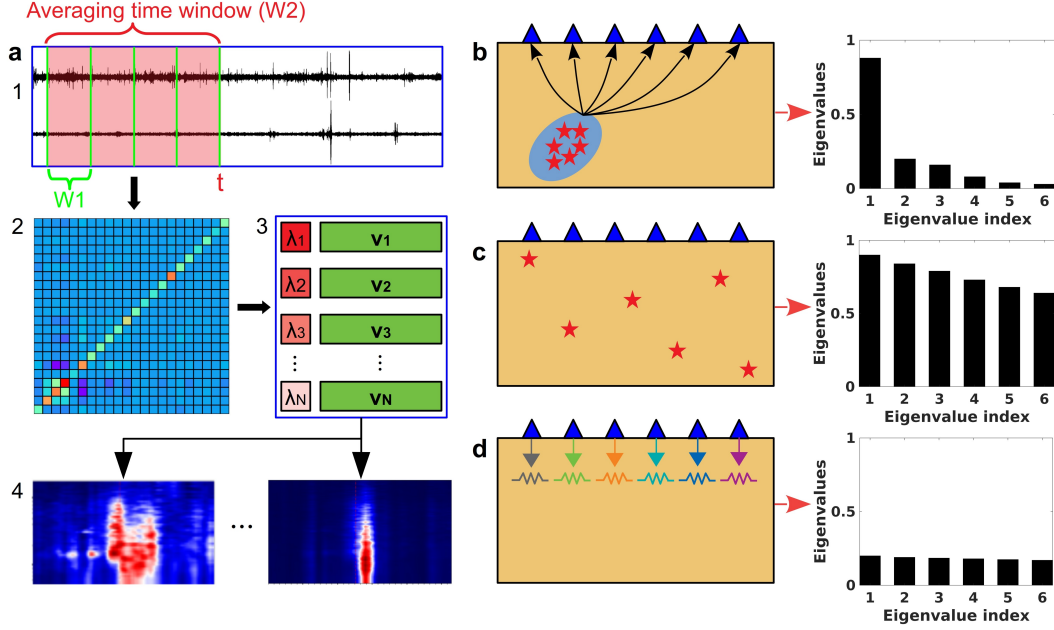


Figure 2. (a) Workflow of wavefield features extraction and analysis, which includes: 1. seismic data processing (e.g. filtering etc.) and time window determination, 2. covariance matrix calculation, 3. eigendecomposition of covariance matrix, and 4. wavefield feature extraction. Right panel shows three representative scenarios of source distribution and the corresponding eigenvalue distribution of covariance matrix in the time window of analysis, which are (b) many co-located seismic sources, (c) many independent (spatially scattered) seismic sources, and (d) electronic or local non-seismic sources. Blue triangles indicate seismic stations and red stars indicate seismic sources.

228 track the evolution of the wavefield during the seismic cycle (short term in this case, 3
 229 years), and to assess if seismic signals contain information about the evolution of the fault
 230 state.

231 4 Feature Analysis and Clustering

232 4.1 Feature Relationship and Sensitivity Analysis

233 The extracted wavefield features over the full dataset are shown in Figure 3 as a
 234 function of time and frequency over the 3 years centered on the L’Aquila earthquake.
 235 In particular, the coherency and entropy features (Figures 3a and 3b) are increasing and
 236 dropping respectively before the mainshock, suggesting the activation of localized sources
 237 in the 3 months before the mainshock at 1-10 Hz. After the strike of the mainshock, during
 238 the aftershock sequence, the frequency content of the coherent wavefield moves to
 239 a lower frequency range (below 5 Hz). Yet, depending on the ratio between the wave-
 240 length of the seismic wavefield and the seismic array aperture, multiple sources distributed
 241 in space are likely to induce a low coherence value (as depicted in Figure 2c). As observed
 242 in the wavefield features (Figures 3a and 3b), at high frequencies, the coherence almost
 243 vanishes, whereas at lower frequencies (below 5 Hz), a high coherence is still observed
 244 (due to larger wavelengths). This is in agreement with the spread of aftershocks near the
 245 rupture zone due to the stress redistribution after the mainshock. The eigenvalue vari-
 246 ance and first eigenvalue features (Figures 3c and 3d) indicate that the fault is most active
 247 during the aftershock periods. In addition, the eigenvalue variance tends to increase

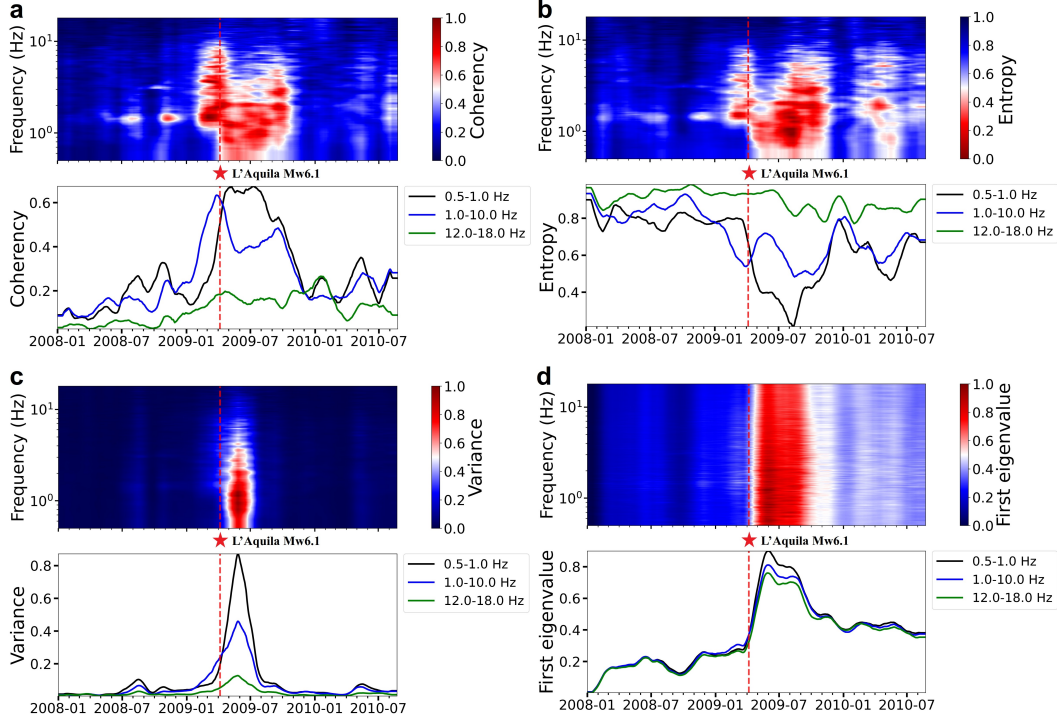


Figure 3. The extracted wavefield features using an averaging window of 60 days. For each sub-figure, the top panel shows the feature with respect to time and frequency (frequency axis in log scale, ranges from 0.5 to 18 Hz), and the bottom panel shows the features averaged in three different frequency bands. The horizontal axis shows time (ranges from 2008-01-01 to 2010-08-23). The red dashed line and star highlight the origin time of the 2009 L'Aquila earthquake. (a) Coherency; (b) Entropy; (c) Eigenvalue variance; (d) First eigenvalue.

248 as the mainshock is approaching, especially in the frequency range of 1-10 Hz, suggest-
 249 ing an activation of relatively strong sources in the area (Figure 3c). The overall time-
 250 frequency evolution of the wavefield features in the studied region visually suggests that
 251 different physical processes are acting during the pre- and post-seismic stages.

252 To quantitatively assess if the observed features can isolate different stages of the
 253 seismic cycle (e.g. pre- and post-seismic) we apply an unsupervised class-membership
 254 identification (clustering). Our approach is similar to the clustering of laboratory data
 255 of Bolton et al. (2019). Our scope is to naturally separate periods with potential differ-
 256 ent physical processes in the fault region, solely from data. We thus avoid any explicit
 257 physical modeling (e.g. location of events, magnitude estimation) and time constrain (e.g.
 258 before and after the earthquake), and learn relevant characteristics with implicit mod-
 259 els from the data itself.

260 Visually, some of the proposed features (e.g. entropy and coherency, Figure 3) show
 261 some similarities, and will be likely redundant in the identification of classes. To quan-
 262 tify any redundancy in our dataset, we analyze the relationship between wavefield fea-
 263 tures and select the uncorrelated ones (i.e. features that are independent and respon-
 264 sible for different source properties) for clustering.

265 To that scope, we calculate the correlation coefficients between different features
 266 in different frequency ranges. Results of this analysis are reported in Figure 4. The en-
 267 tropy and coherency, which provide an estimate of the wavefield coherence, are well cor-

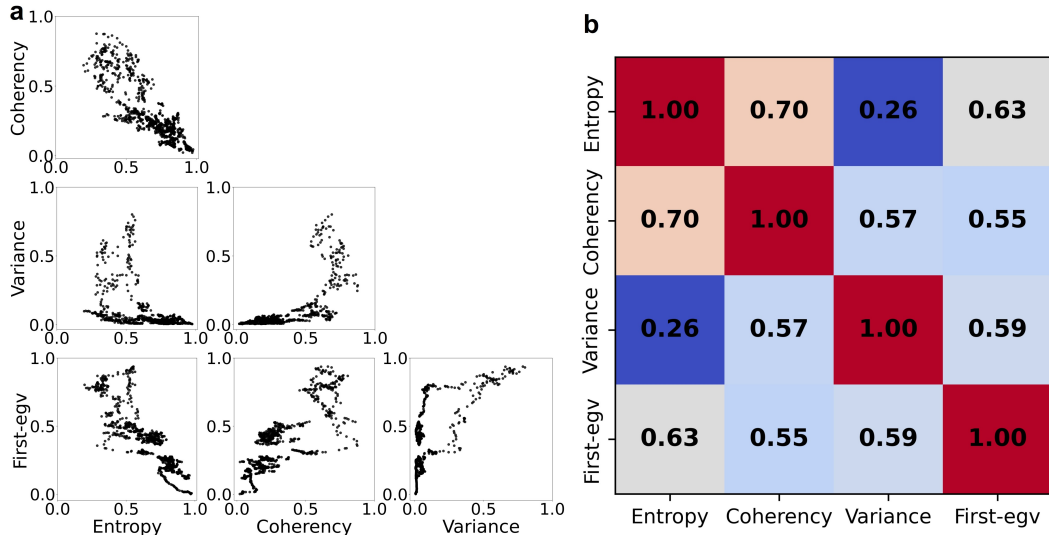


Figure 4. Correlation analysis between different features. (a) Cross-plot of different features at frequency band: 2-2.1 Hz (corresponds to the lower triangular part of the correlation coefficient matrix in (b)). (b) Average correlation coefficients between different features over the full frequency range (0.5-18 Hz). Correlation coefficients are first calculated based on the average features of 0.1 Hz frequency bins, and then are averaged to obtain the final correlation coefficients over the whole frequency band.

268 related with each other over a broad frequency range (0.5-18 Hz) with an average cor-
 269 relation coefficient of 0.7. The eigenvalue variance shows an average correlation coeffi-
 270 cient of 0.57 and 0.59 with the coherency and first eigenvalue respectively, which indi-
 271 cates it contains information about the wavefield coherence and the source energy at the
 272 same time.

273 4.2 Cluster Analysis

274 According to the sensitivity analysis of all features (Section 4.1), the coherency, eigen-
 275 value variance and first eigenvalue are poorly correlated (Figure 4) indicating a sensi-
 276 tivity to different properties of the wavefield (Figures 2 and 3). These three features are
 277 thus selected for the unsupervised analysis. For each time window, the number of fre-
 278 quency points is large (2800 points), therefore defining a very large feature space of 3
 279 x 2800 dimensions. In order to reduce the dimension of the feature space, we focus on
 280 the sensitive frequency range (0.5-10 Hz) and average each feature in frequency bins of
 281 0.1 Hz from 0.5 to 10 Hz. We end up with 95 frequency bins for each of the three fea-
 282 tures. In addition, we linearly normalize the feature magnitude in the interval [0, 1]
 283 with the feature maximum over all the frequencies in order to balance the information pro-
 284 vided by each feature (e.g. Bolton et al., 2019). In this way, the relative amplitude of
 285 the features in different frequency bins is kept. Finally, the three normalized features
 286 are combined together, forming a feature space of 285 dimensions (3 x 95) for cluster
 287 analysis.

288 We extract 966 samples (time segments of W2) in total over the dataset for clus-
 289 tering analysis. Clusters found in seismic data are likely to be unbalanced, because dif-
 290 ferent physical processes may occur at different timescales (e.g. seismic data are mostly
 291 composed by noise). Yet, many clustering approaches are essentially based on the clus-
 292 ter size balance in order to evaluate the clustering quality (for instance K-Means). More

generally, class imbalance is a general issue in clustering, and only few algorithms allow to overcome this problem. Hierarchical clustering (Maimon & Rokach, 2005) is recognized as one of the most powerful approach to cluster unevenly distributed class of data. This is done by building a hierarchy of nested clusters by successively merging or splitting data samples based on any pairwise distance between the data points. In this study, we use an agglomerative strategy which treats each data sample as a cluster and successively merges the two clusters with the smallest distance until all clusters are gathered by a root cluster (Pedregosa et al., 2011). We use L1 distance to measure the distances between data samples.

The hierarchy of our clustering can be represented by a dendrogram, which indicates the distance and splitting between clusters (Figure 5a). We then use a silhouette analysis (Rousseeuw, 1987) to determine the optimal number of clusters (Figures 5b and 5c). The silhouette score is a measure of the average distance between a sample in one cluster to the samples in the neighboring clusters and thus provides a way to assess cluster separation. It is calculated from the normalized difference between the mean nearest inter-cluster distance and the mean intra-cluster distance. Therefore, a large average silhouette score generally indicates large separating distances between the resulting clusters, and hence better clustering results. We vary the number of clusters between 3 to 15, and found that 6 clusters allow to achieve the best separation (Figures 5b and 5c).

4.3 Clustering Results

Because the dimension of the feature space is large, we propose to visualize the clustering results from the two main principal features components. We extract these components with principal component analysis (PCA) as shown in Figure 6. PCA projects data from the original feature space into a principal component (PC) space. Each PC is a linear combination of all the original features, scaled by a corresponding correlation coefficient. PCA also allows to observe the data variance explained by each component. In our case, we see that the first three PCs (PC1-PC3) respectively explain about 80%, 10% and 6% of the total data variance, while all other PCs account for less than 1% of the total data variance each (Figure 6a). Since the first two PCs account for almost 90% of the data variance together, we can thus effectively represent and visualize our data in a 2D PC space.

We use PCA to identify the most relevant wavefield features and frequency ranges to each PC by looking at the linear combination coefficients of the original features, which is useful to interpret the clustering results in a more physical way (Figure 6b). The PCA results indicate that the first PC is highly correlated with the first eigenvalue (Figure 6b), while the second PC is highly related to the wavefield coherence (Figure 6b).

The clustering results are presented in the space formed by the first two principal components in Figure 7. Six clusters are presented along with other independent measurements, i.e. GPS displacement and seismic catalog (Figure 7c). As shown in Figure 7a, the six clusters are well separated in the PC space indicating there are clear and well recognizable patterns in the continuous seismic wavefield. The distribution of different clusters in the original feature space also demonstrates the clustering results are a natural partition according to the wavefield property variations (Figures 8 and 9). The temporal evolution of the clustered data points is shown in the PC space (Figure 7b) and corresponding to each measurement (i.e., PC1-PC3, GPS and seismic catalog, Figure 7c). In Figure 7c, the different PCs, GPS measurements and seismic catalog are color-coded according to the identified clusters to better observe differences among the different clusters.

Before discussing the properties for each cluster, it is worth to remind that the features are extracted from 60 days of data, and each point in Figure 7 is at the end of the time window. Thus, each point has seen data for the preceding 60 days (see Figure 7c),

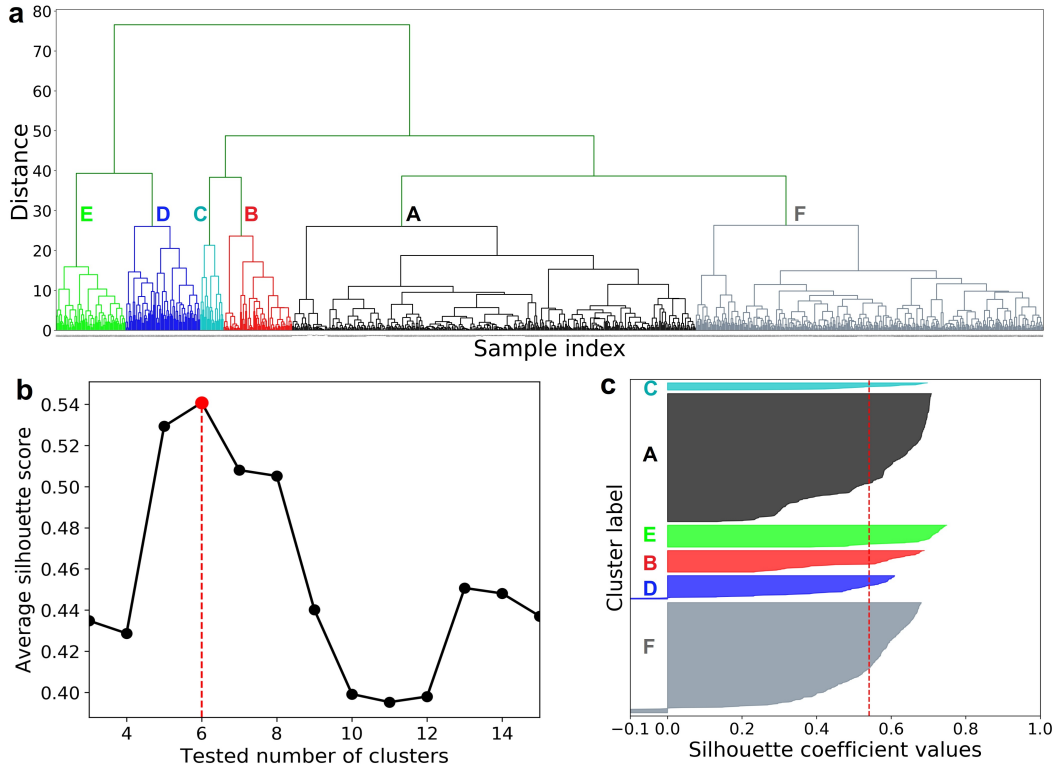


Figure 5. (a) Dendrogram of hierarchical clustering. Different clusters are marked by different colors and annotated using the cluster labels: A-F. The color-code and label of different clusters are consistent with that in Figure 7. The sample index correspond to the date index. (b) Variation of average silhouette score with the number of clusters. Red dashed line indicates that when the number of clusters is 6, the silhouette score reaches to a maximum of about 0.54. (c) Silhouette scores of the data points in each cluster when the number of cluster is 6. Different colors correspond to different clusters. Red dashed line shows the average silhouette score. Most data points in the six clusters have a silhouette score larger than the average score, which indicates a favorable clustering result.

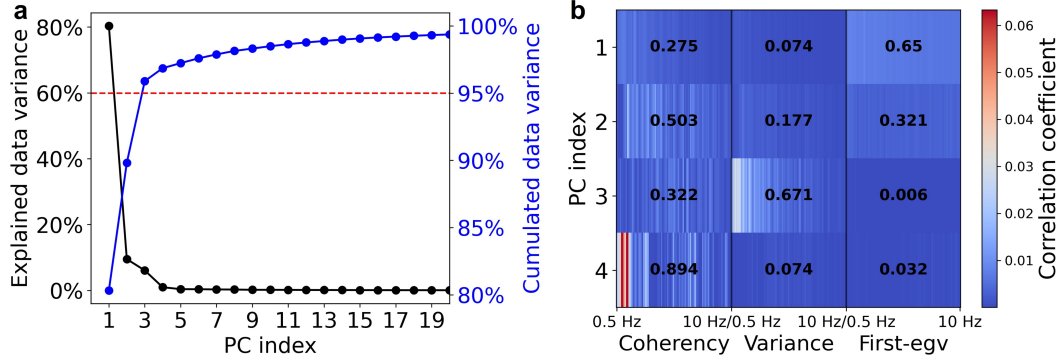


Figure 6. PCA of all the input features. (a) Black line shows the explained data variance (in percentage) of the first 20 principal components (correspond to the left axis). Blue line shows the cumulative explained data variance for the number of principal components used (correspond to the right axis). Red dashed line highlights 95% cumulative percentage. (b) The correlation coefficients between the first four principal components and the features in different frequency bins. The number marked on each section shows the cumulative correlation coefficient over the whole adopted frequency range (0.5-10 Hz) for the coherency, eigenvalue variance and first eigenvalue feature, respectively.

344 and for example, cluster C contains a mixture of signals from times prior and after the
345 mainshock.

346 Cluster A identified with a low wavefield coherency (Figures 3a, 3b, and 9) and small
347 first eigenvalues (Figures 3d and 9), corresponds to a quiet period (low seismicity). Clus-
348 ter B exhibits increased wavefield coherency (especially in the frequency band of 1-10
349 Hz, see in Figures 3a, 3b, and 9), eigenvalue variance (Figure 3c and 9), and first eigen-
350 values (Figures 3d and 9). It corresponds to the increment of seismic activity prior the
351 2009 L'Aquila earthquake. During this period, the earthquake rate increased in this re-
352 gion (Sugan et al., 2014; Vuan et al., 2018) and the earthquakes also tend to localize around
353 the fault (Figures 7c and 9).

354 Clusters C is likely resolving the last period before the main event, but is also af-
355 fected by the mainshock and some aftershocks. It is showing clear differences respect to
356 A and B, in particular an increment of first eigenvalue and a reduction of coherency at
357 1-10 Hz (Figures 3a, 3d, 7c and 9). The group D, which shows strong wavefield coherency
358 in the low frequency range (0.5-1 Hz) and large first eigenvalues (Figures 3a, 3d, and 7c),
359 corresponds to a short period of aftershock sequences immediately after the 2009 L'Aquila
360 earthquake.

361 Compared with cluster D, cluster E shows increasing wavefield coherency (at 1-10
362 Hz) and decreasing first eigenvalues (Figures 3, 7c and 9). It is worth to note that al-
363 though both clusters D and E fall into aftershock sequences, they exhibit distinct coherency
364 variations in different frequency ranges (0.5-1 and 1-10 Hz, see in Figures 3 and 9). More-
365 over, there is a jump in PC2 from cluster D to E (Figure 7c). According to the PCA
366 analysis (figure 6), PC2 is mainly related to wavefield coherency. Therefore, the jump in PC2
367 from cluster D to E is mainly due to a change in the wavefield coherency, which can be
368 confirmed in the extracted coherency feature (Figures 3 and 9). Compared with clus-
369 ter D, the wavefield coherency of cluster E increases at 1-10 Hz. This behavior can be
370 interpreted (see in Figure 2) as an activation of localized seismic sources of low magni-
371 tudes (especially the event cluster in 30 km away to the main event, Figure 9). These
372 observations suggest an evolution of the aftershock behavior. During the period of clus-

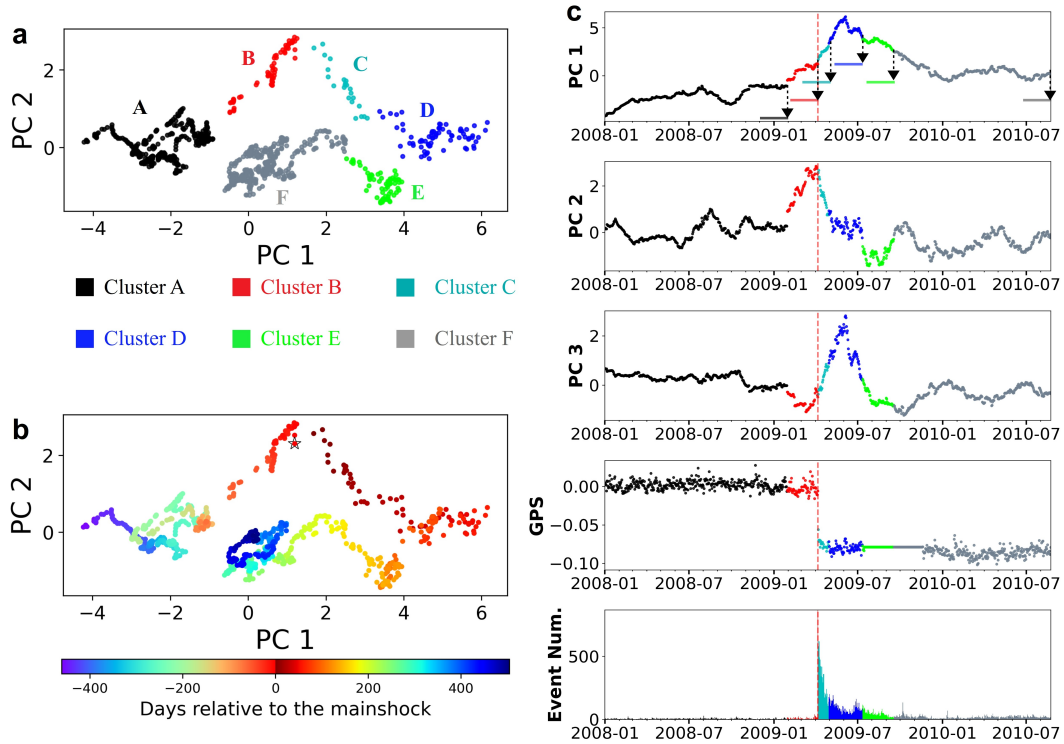


Figure 7. Clustering results shown in the 2D PC space with horizontal axis showing the first PC and vertical axis showing the second PC. Six clusters are color-coded and marked with labels A to F (consistent with Figure 5). (b) Temporal variation of the clustered points in the 2D PC space. The data points are color-coded according to the days relative to the mainshock as indicated by the colorbar at the bottom. The star highlights the day when the 2009 L'Aquila earthquake occurred. (c) Temporal variation of the principal components, GPS measurements and number of seismic events per day. The red dashed lines exhibits the origin time of the 2009 L'Aquila earthquake. The first to third rows show the variation of the first three PCs with time. The fourth row shows the ground displacements in the vertical direction measured by a GPS station in L'Aquila (location shown in Figure 1a). The fifth row shows the detected number of seismic events per day in the INGV catalog. The different measurements are color-coded according to the corresponding cluster. The time window (60 days) for extracting wavefield features at the last data sample in each cluster is highlighted by the black arrow and the corresponding color-coded bar in the top panel.

373 ter E, the earthquake rate is much lower than the previous aftershock stages (C and D)
 374 and localized swarm-like seismicity of low magnitudes emerges (Figures 7c and 9).

375 The last cluster (F) shows low wavefield coherency and steady decreasing first eigen-
 376 values (Figures 3 and 7). During this period, the aftershocks sequence reduces and the
 377 earthquake rate in the region starts to recover to a background level (Figures 7c and 9).
 378 As shown in the dendrogram (Figure 5a) and in Figure 7a, the A and F clusters are close
 379 to each other and belong to the same root cluster. Compared to the other clusters which
 380 are more seismically active, they correspond to quieter periods.

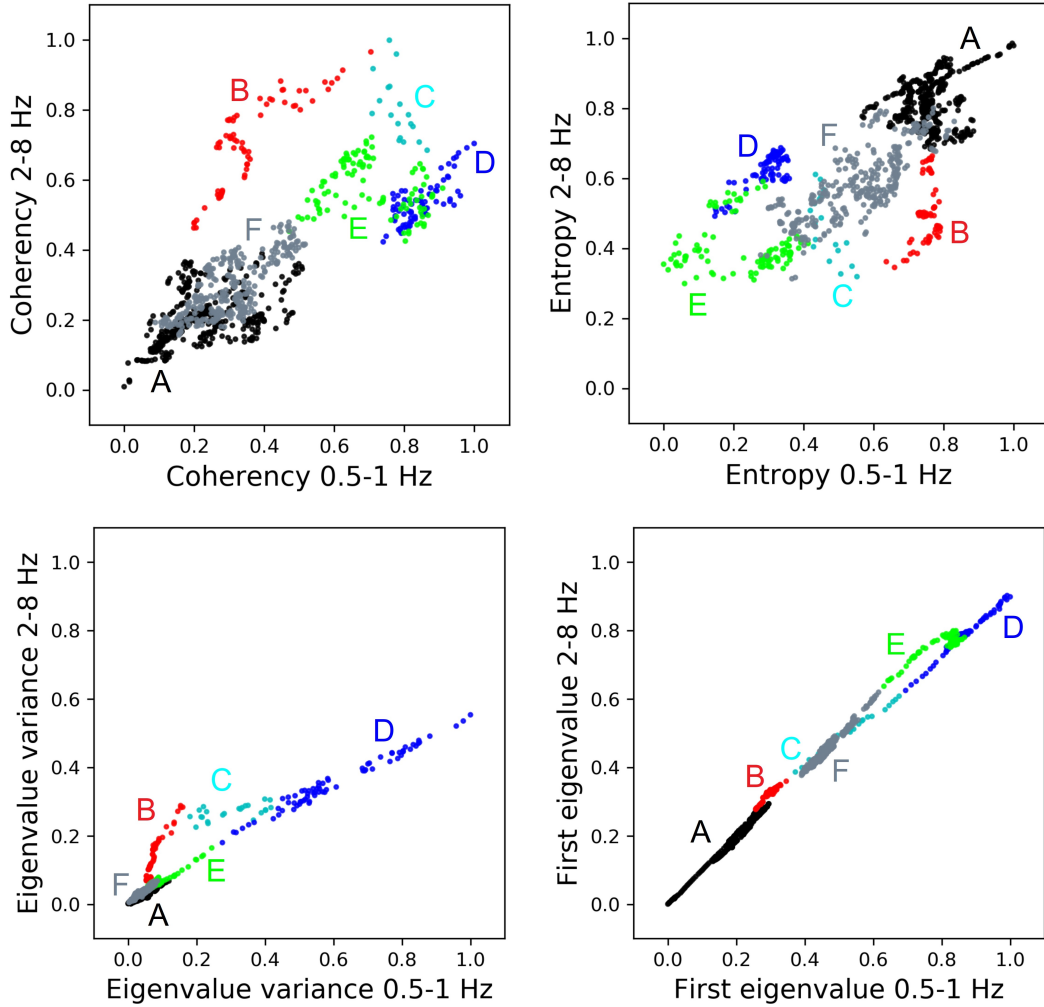


Figure 8. The distribution of the identified clusters in the wavefield feature space. The wavefield features are extracted at each frequency point from 0.5-18 Hz. Here for better visualizing the clustering results in a 2D feature space, the features are averaged and shown in two frequency bands, which are (1) low frequency band: 0.5-1 Hz and (2) higher frequency band: 2-8 Hz.

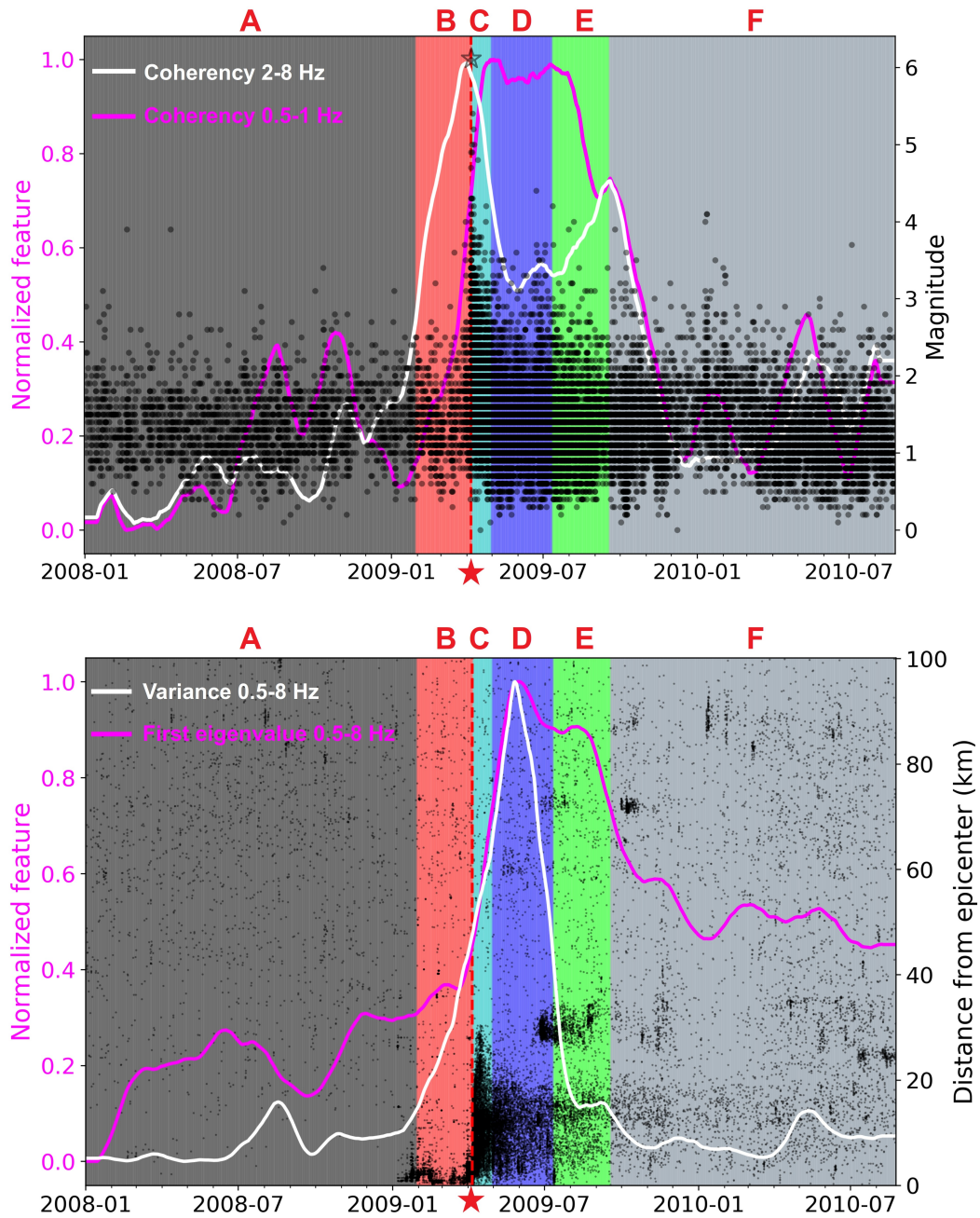


Figure 9. Clustering results with different colors representing different clusters (clusters: A-F). Seismic events from INGV catalog and the wavefield features are also displayed for comparison. The origin time of the 2009 L'Aquila earthquake is marked by the red star and red dash line. In the top panel, clustering results are shown together with the average coherency feature in different frequency bands (correspond to the left axis) and local magnitudes of seismic events (correspond the right axis). In the bottom panel, clustering results are shown together with the average eigenvalue variance feature (left axis), the first eigenvalue feature (left axis) and the event distances from the epicenter of the 2009 L'Aquila earthquake (right axis).

5 Discussion

We show the existence of the time and frequency evolution of wavefield features derived from continuous seismic records, and how the analysis of these features reveals distinct clusters well correlated with different periods of the seismic cycle (Figure 7). As the analysis is performed over long-term scale, the time-frequency evolution reflects the statistical wavefield properties and is related to the evolution in position, size, number and distribution pattern of the seismic sources inside the array (Figure 2). Hence, by analyzing the statistical properties of continuous wavefield, we draw conclusions about physical processes occurring in the fault region without the need of any modeling.

As the used features have physical meanings, they can provide important information about the processes occurring in each cluster. For example, cluster B, which is characterized by increasing coherence and first eigenvalue, suggests the activation of localized sources prior to the main event (see Figures 2, 3 and 7). This behavior agrees with previous studies on L'Aquila earthquake, suggesting the occurrence of localized foreshocks and increased earthquake rate before the mainshock (Sugan et al., 2014). However, differently from previous studies, no explicit modeling is involved in our analysis, and we show how this information emerges naturally from our chosen representation of the continuous seismic wavefield. Clusters D and E show a reduction of the coherency compared to clusters B and C especially in the high frequency range (1-10 Hz). This behavior suggests that seismicity is spread around the fault, as stress is redistributed after the mainshock (Marsan, 2005). Previous study based on earthquake catalog (Marsan et al., 2014) shows that earthquakes preceded by accelerating seismicity rate produce more aftershocks on average and exhibit more spatial spreading aftershock sequences, which agrees with our model free analysis here. A similar phenomenon has also been recently observed for the Ridgecrest earthquake (Trugman et al., 2020; Ross et al., 2019), where the temporal and spatial variations of the earthquake waveform similarity before and after the 2019 Ridgecrest earthquake are compared. Significant reduction of the earthquake similarity in the aftershock sequences (compared to the pre-event seismicity) is observed and interpreted as a result of small scale heterogeneities in the residual stress field initiated by the main event. Their observations of the temporal variation of coherence using earthquake waveforms of catalogued events correspond well with our results derived from continuous seismic data. But again, our observations and analysis are model free and do not require additional seismological dataset such as earthquake catalogs and velocity models.

More complex is the interpretation of cluster C, which partially covers the last pre-seismic period and a portion of time after the event. This issue comes from the limitation of our methodology to a given timescale. In fact, the use of a long-term window (60 days) with daily step, reduces the possibility of resolving short-lasting clusters and focuses on long-lasting processes. Attempting to reduce the time window will be the subject of future research. However, despite this limitation the method is clearly highlighting different parts of the seismic cycles (including the quiet period, clusters A and F), without the need of modeling.

As in stick-slip rock failure experiments in laboratory (Bolton et al., 2019), our study highlights that fault state can be tracked from continuous seismic data. The ability of unrevealing peculiar patterns in seismic data, extend the laboratory-based idea that continuous data are rich enough to inform us about evolution of physical properties of the fault (e.g. Rouet-Leduc et al., 2017; Bolton et al., 2019). In contrast with the laboratory setting, real data cannot be associated with other boundary conditions (e.g. absolute stress level), and only part of the seismic cycle can be resolved. It is thus unlikely that our features-based approach will permit any kind of machine learning based prediction of the rupture (e.g. Rouet-Leduc et al., 2017). It will however permit to rapidly parse large amount of data and extract peculiar patterns, which can be related to other estimates (e.g. geodetic data) to better characterize different stages of the seismic cy-

434 cle. Our features can also be used to regress seismic data into other information (e.g. GPS
435 displacement, Frank et al., 2015) to explore slip rate during aseismic slip episodes.

436 Finally, in the present study, we defined spatial features for exploring spatially dis-
437 tributed sensors. One of the main advantage is the ability to easily identify propagative
438 signals and to disregard any site-dependent patterns that may bias the analysis (e.g. lo-
439 cal noise level). Given the large number of seismic arrays deployed worldwide, the de-
440 velopments of features that account for spatial properties of the wavefield is of great in-
441 terest and will be in the scope of future studies.

442 6 Conclusions

443 We analyze the wavefield properties with unsupervised machine learning to directly
444 assess fault state and its temporal evolution from continuous seismic data. Unlike tra-
445 ditional statistical features calculated from single station, we extract frequency-dependent
446 wavefield features from the array covariance matrix analysis, which provide the inter-
447 pretation of the physical properties of the seismic sources. The array-based wavefield fea-
448 tures enable to analyze the overall source properties and its temporal evolution for un-
449 derstanding the fault activities in the study region. Our study shows the value of advanced
450 array processing and machine learning analysis to reveal information embedded in the
451 continuous seismic data. Our study builds a bridge between the laboratory experiments
452 and the real earthquake observations and is a step towards understanding the fault physics.
453 Our future work involves further unraveling hidden signals in continuous seismic data
454 for studying fault physics.

455 Acknowledgments

456 This research received funding from the European Research Council (ERC) under the
457 European Union Horizon 2020 Research and Innovation Programme (grant agreements,
458 802777-MONIFaults). L.S. acknowledges support from the European Research Coun-
459 cil under the European Union Horizon 2020 research and innovation program (grant agree-
460 ment no. 742335, F-IMAGE). Computations were performed using the UGA High-Performance
461 Computing infrastructures CIMENT. The continuous seismic data and seismic catalog
462 used in this study are available at the Istituto Nazionale Geofisica e Vulcanologia (INGV)
463 seismological data center (http://cnt.rm.ingv.it/webservices_and_software/, last
464 accessed March 2020). The GPS data used in this study are available at the Nevada Geode-
465 tic Laboratory website (<http://geodesy.unr.edu/>, last accessed March 2020). Figure
466 1a was made using the Generic Mapping Tools v.6.0.0 (<http://gmt.soest.hawaii.edu/>,
467 last accessed March 2020; Wessel et al., 2019). The unsupervised machine learning al-
468 gorithm used in this study were based on scikit-learn v.0.22.2 ([https://scikit-learn](https://scikit-learn.org/)
469 [.org/](https://scikit-learn.org/), last accessed March 2020). We thank editor Rachel Abercrombie, associate ed-
470 itor Nori Nakata and two anonymous reviewers for their efforts and constructive com-
471 ments.

472 References

- 473 Aki, K., & Richards, P. G. (2002). *Quantitative seismology*. University Science
474 Books.
- 475 Baccheschi, P., De Gori, P., Villani, F., Trippetta, F., & Chiarabba, C. (2020).
476 The preparatory phase of the Mw 6.1 2009 L’Aquila (Italy) normal faulting
477 earthquake traced by foreshock time-lapse tomography. *Geology*, *48*(1), 49–55.
- 478 Beroza, G. C., & Ide, S. (2011). Slow earthquakes and nonvolcanic tremor. *Annual*
479 *Review of Earth and Planetary Sciences*, *39*, 271–296.
- 480 Bolton, D. C., Shokouhi, P., Rouet-Leduc, B., Hulbert, C., Rivière, J., Marone, C.,
481 & Johnson, P. A. (2019). Characterizing acoustic signals and searching for

- precursors during the laboratory seismic cycle using unsupervised machine learning. *Seismological Research Letters*, *90*(3), 1088–1098.
- Cara, F., Di Giulio, G., & Rovelli, A. (2003). A study on seismic noise variations at Colfiorito, central Italy: implications for the use of H/V spectral ratios. *Geophysical Research Letters*, *30*(18).
- Carniel, R., Jolly, A. D., & Barbui, L. (2013). Analysis of phreatic events at Ruapehu volcano, New Zealand using a new SOM approach. *Journal of Volcanology and Geothermal Research*, *254*, 69–79.
- Chiarabba, C., Amato, A., Anselmi, M., Baccheschi, P., Bianchi, I., Cattaneo, M., ... others (2009). The 2009 L'Aquila (central Italy) MW6.3 earthquake: Main shock and aftershocks. *Geophysical Research Letters*, *36*(18).
- Chiarabba, C., Buttinelli, M., Cattaneo, M., & De Gori, P. (2020). Large earthquakes driven by fluid overpressure: The Apennines normal faulting system case. *Tectonics*, *39*(4), e2019TC006014.
- Di Luccio, F., Ventura, G., Di Giovambattista, R., Piscini, A., & Cinti, F. (2010). Normal faults and thrusts reactivated by deep fluids: The 6 April 2009 Mw 6.3 L'Aquila earthquake, central Italy. *Journal of Geophysical Research: Solid Earth*, *115*(B6).
- Di Stefano, R., Chiarabba, C., Chiaraluce, L., Cocco, M., De Gori, P., Piccinini, D., & Valoroso, L. (2011). Fault zone properties affecting the rupture evolution of the 2009 (Mw 6.1) L'Aquila earthquake (central Italy): Insights from seismic tomography. *Geophysical Research Letters*, *38*(10).
- Esposito, A. M., Giudicepietro, F., D'Auria, L., Scarpetta, S., Martini, M. G., Coltelli, M., & Marinaro, M. (2008). Unsupervised neural analysis of very-long-period events at Stromboli volcano using the self-organizing maps. *Bulletin of the Seismological Society of America*, *98*(5), 2449–2459.
- Frank, W. B., Radiguet, M., Rousset, B., Shapiro, N. M., Husker, A. L., Kostoglodov, V., ... Campillo, M. (2015). Uncovering the geodetic signature of silent slip through repeating earthquakes. *Geophysical Research Letters*, *42*(8), 2774–2779.
- Gersztenkorn, A., & Marfurt, K. J. (1999). Eigenstructure-based coherence computations as an aid to 3-D structural and stratigraphic mapping. *Geophysics*, *64*(5), 1468–1479.
- Gutenberg, B., & Richter, C. F. (1956). Magnitude and energy of earthquakes. *Annali di Geofisica*, *9*(1), 1–15.
- Herrmann, R. B., Malagnini, L., & Munafò, I. (2011). Regional Moment Tensors of the 2009 L'Aquila Earthquake Sequence. *Bulletin of the Seismological Society of America*, *101*(3), 975–993.
- Hulbert, C., Rouet-Leduc, B., Johnson, P. A., Ren, C. X., Rivière, J., Bolton, D. C., & Marone, C. (2019). Similarity of fast and slow earthquakes illuminated by machine learning. *Nature Geoscience*, *12*(1), 69–74.
- Ide, S., Beroza, G. C., Shelly, D. R., & Uchide, T. (2007). A scaling law for slow earthquakes. *Nature*, *447*(7140), 76–79.
- INGV Seismological Data Centre. (2006). *Rete Sismica Nazionale (RSN). Istituto Nazionale di Geofisica e Vulcanologia (INGV), Italy.* <https://doi.org/10.13127/SD/X0FXNH7QFY>.
- Köhler, A., Ohrnberger, M., & Scherbaum, F. (2010). Unsupervised pattern recognition in continuous seismic wavefield records using self-organizing maps. *Geophysical Journal International*, *182*(3), 1619–1630.
- Langer, H., Falsaperla, S., Masotti, M., Campanini, R., Spampinato, S., & Messina, A. (2009). Synopsis of supervised and unsupervised pattern classification techniques applied to volcanic tremor data at Mt Etna, Italy. *Geophysical Journal International*, *178*(2), 1132–1144.
- Langer, H., Falsaperla, S., Messina, A., Spampinato, S., & Behncke, B. (2011). Detecting imminent eruptive activity at Mt Etna, Italy, in 2007–2008 through

- 537 pattern classification of volcanic tremor data. *Journal of Volcanology and*
538 *Geothermal Research*, 200(1-2), 1–17.
- 539 Lucente, F. P., De Gori, P., Margheriti, L., Piccinini, D., Di Bona, M., Chiarabba,
540 C., & Agostinetti, N. P. (2010). Temporal variation of seismic velocity and
541 anisotropy before the 2009 MW 6.3 L’Aquila earthquake, Italy. *Geology*,
542 38(11), 1015–1018.
- 543 Maimon, O., & Rokach, L. (2005). *Data mining and knowledge discovery handbook*.
544 Springer.
- 545 Marsan, D. (2005). The role of small earthquakes in redistributing crustal elastic
546 stress. *Geophysical Journal International*, 163(1), 141–151.
- 547 Marsan, D., Helmstetter, A., Bouchon, M., & Dublanchet, P. (2014). Foreshock ac-
548 tivity related to enhanced aftershock production. *Geophysical Research Letters*,
549 41(19), 6652–6658.
- 550 Menon, R., Gerstoft, P., & Hodgkiss, W. S. (2014). On the apparent attenuation
551 in the spatial coherence estimated from seismic arrays. *Journal of Geophysical*
552 *Research: Solid Earth*, 119(4), 3115–3132.
- 553 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ...
554 others (2011). Scikit-learn: Machine learning in Python. *Journal of Machine*
555 *Learning Research*, 12, 2825–2830.
- 556 Poli, P., Boaga, J., Molinari, I., Cascone, V., & Boschi, L. (2020). The 2020 coro-
557 navirus lockdown and seismic monitoring of anthropic activities in Northern
558 Italy. *Scientific Reports*, 10(1), 1–8.
- 559 Ross, Z. E., Idini, B., Jia, Z., Stephenson, O. L., Zhong, M., Wang, X., ... others
560 (2019). Hierarchical interlocked orthogonal faulting in the 2019 ridgecrest
561 earthquake sequence. *Science*, 366(6463), 346–351.
- 562 Rouet-Leduc, B., Hulbert, C., Lubbers, N., Barros, K., Humphreys, C. J., & John-
563 son, P. A. (2017). Machine learning predicts laboratory earthquakes. *Geophys-
564 ical Research Letters*, 44(18), 9276–9282.
- 565 Rousseeuw, P. J. (1987). Silhouettes: a graphical aid to the interpretation and vali-
566 dation of cluster analysis. *Journal of Computational and Applied Mathematics*,
567 20, 53–65.
- 568 Scholz, C. H. (2002). *The mechanics of earthquakes and faulting*. Cambridge univer-
569 sity press.
- 570 Seydoux, L., Shapiro, N. M., De Rosny, J., Brenguier, F., & Landès, M. (2016a).
571 Detecting seismic activity with a covariance matrix analysis of data recorded
572 on seismic arrays. *Geophysical Journal International*, 204(3), 1430–1442.
- 573 Seydoux, L., Shapiro, N. M., De Rosny, J., & Landès, M. (2016b). Spatial coher-
574 ence of the seismic wavefield continuously recorded by the USArray. *Geophys-
575 ical Research Letters*, 43(18), 9644–9652.
- 576 Shannon, C. E. (1948). A mathematical theory of communication. *The Bell System*
577 *Technical Journal*, 27(3), 379–423.
- 578 Shreedharan, S., Bolton, D. C., Rivière, J., & Marone, C. (2020). Preseismic fault
579 creep and elastic wave amplitude precursors scale with lab earthquake mag-
580 nitude for the continuum of tectonic failure modes. *Geophysical Research*
581 *Letters*, 47(8), e2020GL086986.
- 582 Soubestre, J., Seydoux, L., Shapiro, N. M., De Rosny, J., Droznin, D., Droznina,
583 S. Y., ... Gordeev, E. (2019). Depth migration of seismovolcanic tremor
584 sources below the Klyuchevskoy volcanic group (Kamchatka) determined from
585 a network-based analysis. *Geophysical Research Letters*, 46(14), 8018–8030.
- 586 Soubestre, J., Shapiro, N. M., Seydoux, L., De Rosny, J., Droznin, D. V., Droznina,
587 S. Y., ... Gordeev, E. I. (2018). Network-based detection and classification
588 of seismovolcanic tremors: Example from the Klyuchevskoy volcanic group in
589 Kamchatka. *Journal of Geophysical Research: Solid Earth*, 123(1), 564–582.
- 590 Sukan, M., Kato, A., Miyake, H., Nakagawa, S., & Vuan, A. (2014). The prepara-
591 tory phase of the 2009 Mw 6.3 L’Aquila earthquake by improving the detection

- 592 capability of low-magnitude foreshocks. *Geophysical Research Letters*, *41*(17),
593 6137–6144.
- 594 Trugman, D. T., Ross, Z. E., & Johnson, P. A. (2020). Imaging stress and faulting
595 complexity through earthquake waveform similarity. *Geophysical Research Let-*
596 *ters*, *47*(1), e2019GL085888.
- 597 Unglert, K., Radić, V., & Jellinek, A. M. (2016). Principal component analysis vs.
598 self-organizing maps combined with hierarchical clustering for pattern recog-
599 nition in volcano seismic spectra. *Journal of Volcanology and Geothermal*
600 *Research*, *320*, 58–74.
- 601 Von Neumann, J. (1986). *Papers of John von Neumann on computers and computer*
602 *theory*. The MIT Press, Cambridge, MA.
- 603 Vuan, A., Sukan, M., Amati, G., & Kato, A. (2018). Improving the Detection of
604 Low-Magnitude Seismicity Preceding the Mw 6.3 L’Aquila Earthquake: De-
605 velopment of a Scalable Code Based on the Cross Correlation of Template
606 Earthquakes. *Bulletin of the Seismological Society of America*, *108*(1), 471–
607 480.
- 608 Wessel, P., Luis, J., Uieda, L., Scharroo, R., Wobbe, F., Smith, W., & Tian, D.
609 (2019). The generic mapping tools version 6. *Geochemistry, Geophysics,*
610 *Geosystems*, *20*(11), 5556–5564.