1

# The Assumption of Poisson Seismic-Rate Variability in

# CSEP/RELM Experiments

4

**A.M. Lombardi & W. Marzocchi**

**Istituto Nazionale di Geofisica e Vulcanologia, Via di Vigna Murata 605, 00143 Rome, Italy.**

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

**Abstract**

Evaluating the performances of earthquake forecasting/prediction models is the main rationale behind some recent international efforts like the Regional Earthquake Likelihood Model (RELM) and the Collaboratory for the Study of Earthquake Predictability (CSEP). Basically, the evaluation process consists of two steps: 1) to run simultaneously all codes to forecast future seismicity in well-defined testing regions; 2) to compare the forecasts through a suite of statistical tests. The tests are based on the likelihood score and they check both the time and space performances. All these tests rely on some basic assumptions that have never been deeply discussed and analyzed. In particular, models are required to specify a rate in space-time-magnitude bins, and it is assumed that these rates are independent and characterized by Poisson uncertainty. In this work we have explored in detail these assumptions and their impact on CSEP testing procedures when applied to a widely used class of models, i.e., the Epidemic-Type Aftershock Sequence (ETAS) models. Our results show that, if an ETAS model is an accurate representation of seismicity, the same "right" model is rejected by the current CSEP testing procedures a number of times significantly higher than expected. We show that this deficiency is due to the fact that the ETAS models produce forecasts with a variability significantly higher than that of a Poisson process, invalidating one of the main assumption that stands behind the CSEP/RELM evaluation process. Certainly, this shortcoming does not negate the paramount importance of the CSEP experiments as a whole, but it does call for a specific revision of the testing procedures to allow a better understanding of the results of such experiments.

## 1.    Introduction

The success of operational forecast indispensably depends on the use of reliable and skillful models (ICEF, 2009). In a nutshell, a model has to produce forecasts/predictions compatible with the future seismicity, and the forecasts/predictions have to be precise enough to be usable for practical

2

50    purposes (i.e., they need a good skill). Moreover, if a set of reliable models is available, it is

51    important to know what is the "best" one(s), i.e., the one(s) with the highest skill.

52    The evaluation of these pivotal features characterizing each forecasting/prediction model is the

53    primary goal of the Collaboratory for the Study of Earthquake Predictability (CSEP hereinafter;

54    Jordan 2006; http://www.cseptesting.org).

55    CSEP provides a rigorous framework for an empirical evaluation of any forecasting and prediction

56    model. CSEP can be considered the successor of the Regional Earthquake Likelihood Model

57    (RELM) experiment (Schorlemmer and Gerstenberger, 2007). While RELM was focusing on

58    California, CSEP extends this focus to many other regions (New Zealand, Italy, Japan, North- and

59    South-Western Pacific, and the whole World) as well as global testing centers (New Zealand,

60    Europe, Japan). The coordinated international experiment has two main advantages: the evaluation

61    process is supervised by an international scientific committee, not only by the modelers themselves,

62    and the cross-evaluation of a model performances in different regions of the world can facilitate its

63    evaluation in a much shorter period of time (see also Zechar et al., 2009).

64    All CSEP experiments performed in each testing region are truly prospective tests. In other words,

65    each experiment compares forecasts produced by several models under testing with real data

66    observed in the corresponding testing region after the forecasts have been produced. The forecasts

67    are generated in the testing center independent of the modelers. The testing procedure adopted can

68    be summarized in two subsequent steps: 1) to measure the *reliability* of each model; 2) to quantify

69    the relative *skill* among the set of reliable models. In the first step, the forecasts/predictions made by

70    each model are compared to the real seismicity through one or more goodness-of-fit tests. If the

71    seismicity observed is compatible with the output of the model and the model-based variability,

72    then the performance of the models can be contrasted with other models in the second step of the

73    analysis.    Specifically,    the    second    step    of    the    analysis    compares    quantitatively    the

74    forecasting/prediction capabilities of the models in order to establish a hierarchy of best performing

75    models.

76   In this paper, we explore the performances of the CSEP/RELM testing procedure for two classes of

77   forecasting models, Poisson and ETAS, that are largely represented in CSEP/RELM experiments

78   (for the reliability of the prediction models see, e.g., Marzocchi et al., 2003; Zechar & Jordan, 2008,

79   and references therein).

80

81   **2.      The CSEP/RELM suite of tests**

82   The CSEP/RELM suite of tests is originally composed of three different tests (Schorlemmer et al.,

83   2007; see also Kagan and Jackson 1994; 1995). The *L*-test (*Data-consistency test*) and *N*-test

84   (*Number of events test*) are intended to check the goodness-of-fit of the model, while the *R*-test

85   (*Hypotheses comparison*) compares the forecasting performances of different models.

86   The *L*-test and *R*-test are based on the well-known concept of conditional likelihood that is one of

87   most used statistical tools to check and compare the performance of one or more models on data.

88   The formulation of these tests requires the definition of bins that are specified intervals in space,

89   magnitude and time. Using the same symbols of *Schorlemmer et al. (*2007*)*, we define:

$\omega_i$                                number of earthquakes occurred in the *i*-th bin

$\lambda_i^j$                                rate of earthquake occurrence for the *i*-th bin and *j*-th model.

$L_i^j = L(\omega_i \mid \lambda_i^j)$          log-likelihood calculated for the *i*-th bin and *j*-th model

90

91   The joint log-likelihood for the *j*-th model is calculated as

$$L^j = \sum_{i=1}^{n} L(\omega_i \mid \lambda_i^j) \tag{1}$$

93   where *n* is the number of bins.

94   In order to get numbers from equation (1) $L(\omega_i / \lambda_i^j)$ must be defined. The basic assumption that

95   stands behind the CSEP/RELM testing procedure is that earthquakes are assumed to occur in each

96   bin according to a Poisson process with the rate specified by the model (Schorlemmer et al., 2007).

97    Note that this assumption is associated with the CSEP/RELM testing procedure not with the

98    loglikelihood tests that can manage any kind of arbitrary distributions. Therefore, equation (1)

99    becomes

100
$$L^j = \sum_{i=1}^{n} L(\omega_i \mid \lambda_i^j) = \sum_{i=1}^{n} \left( -\lambda_i^j + \omega_i \ln \lambda_i^j - \ln \omega_i! \right)$$
(2)

101   This assumption is crucial and a careful evaluation of its validity is mandatory to fully understand

102   the CSEP/RELM tests. This assumption means that the bins are spatially and temporally

103   independent, and the number of earthquakes in time has a variance equal to the average. Although

104   some authors have already categorized such assumptions as "unlikely" and foresee possible

105   inconsistencies of the tests (e.g., Werner and Sornette, 2008), the consequences have never been

106   explored in detail. Moreover, we argue that the current use of this testing procedure in CSEP

107   experiments may lead to think that the departures from this hypothesis could be considered as

108   negligible.

109   The log-likelihood obtained by equation (2) is used to get the significance level of the tests through

110   simulations. The *L*-test compares the observed log-likelihood value (see equation (2)) with a

111   prefixed number of synthetic values obtained under the Poisson assumption for each bin, i.e.,

112   simulating records where each bin has a number of earthquakes generated according to a Poisson

113   process with the rate given by the model. The quantile score $\gamma^j$ for the *j*-th model is the fraction of

114   simulated likelihood values that are less or equal to the observed *L*. This quantile score can be

115   considered the p-value of the test. Note that, compared to the analyses performed by Schorlemmer

116   et al. (2007) and Werner and Sornette (2008), here we do not consider the inclusion of

117   *uncertainties*, because we aim to explore the tests in an optimal situation, i.e., with negligible

118   uncertainty in the observations.

119   Schorlemmer et al. (2007) discussed the case in which a model can pass the *L*-test even if it is

120   wrong. For this reason, the authors proposed a second test, the *N*-test, that checks if the total

121   number of forecasted events is compatible with the observed number. In this case the quantile score,

122    $\delta^j$, is the probability to have no more than the observed number of events by a Poisson process

123    with a rate given by the model. In this case the test is two-sided, checking both possible over-

124    prediction and under-prediction. To summarize, a model is "good" (*reliable*) if it is not rejected by

125    both $L$ and $N$ tests. Only if the model passes these tests, then it is considered in the $R$-test, where it

126    is compared to other reliable models. In the next section, we explore the performances of the L- and

127    N-tests applied to synthetic catalogs. The goal is to check, in a controlled experiment, if the

128    proportion of rejections of the "right" model is comparable to the significance level of the test. We

129    anticipate that possible departures may point to inconsistencies of the Poisson variability for each

130    bin assumed in the CSEP/RELM testing procedure.

131

132    **3.      Application of the CSEP/RELM testing procedure to synthetic catalogs.**

133    In order to evaluate quantitatively the performances of CSEP/RELM testing procedure, we use

134    these tests in a controlled experiment where we know exactly the model that generates earthquakes.

135    The experiment can be described in three steps:

136    1. We generate 100 synthetic catalogues that we call "pseudo-real catalogs". Specifically we

137    simulate two sets of 100 pseudo-real catalogs: one is consistent with a stationary non-homogeneous

138    Poisson process, and another that is consistent with the well-known Epidemic-Type Aftershocks

139    Sequence (ETAS; e.g., Ogata 1998) model. The generation of the ETAS pseudo-real catalogs is

140    described in Appendix A and mimics the 1992 Landers sequence.

141    2. We generate one-day forecasts for a period of 10 days after the mainshock using exactly the same

142    models and relative parameters that generate the pseudo-real catalogs. After each one-day forecast,

143    the history is updated to take into account all events that occurred before the starting time of the

144    next forecast. The forecasts are computed and evaluated in terms of expected number of events with

145    magnitude above $M_l$ 3.0 in each cell $C_i$ of a grid, with a spacing of 0.1°x0.1° and covering the target

146 region [-117.5°W/33.25°N − -115.5°W/35.5°N].  Specifically for each cell $C_i$ and for each time

147 window $T_j$ we compute the relative forecast rate $\lambda_i^j$ by the formula

148
$$\lambda_i^j = \int\limits_{T_j} \iint\limits_{C_i} \int\limits_{M_c}^{M_{max}} \lambda(t,x,y,m/H_t)dtdxdydm \qquad (3)$$

149 where $\lambda(t,x,y|\mathcal{H}_t)$ is the space-time conditional intensity defined by Poisson and ETAS models (see

150 Appendix A), $M_c$ and $M_{max}$ are the minimum and maximum magnitude considered. The seismic

151 history $H_t$, i.e. the information coming from the events that occurred before the time $t$, is crucial for

152 time-dependent models, such as the ETAS model. On the other hand, the Poisson rate is

153 independent of $H_t$ and the time $t$.  For the ETAS model we include in seismic history $H_t$ the

154 parameters of earthquakes that occurred before the time window $T_j$. To take into account the

155 expected triggering effect of events that occurred during $T_j$, we simulate 1000 different stochastic

156 realizations of the model inside the time window $T_j$ and then we calculate for each bin the mean and

157 the variance of predictions $\lambda_i^j$ coming from each of these synthetic realizations.

158 3. We compare each one-day forecast with each pseudo-real catalog for both classes of models

159 (Poisson and ETAS). For each of 100 pseudo-real catalogs we apply the $N$ and $L$ tests in order to

160 verify the agreement between observations and forecasts. In this case, the model is certainly right;

161 therefore we expect to see a number of rejections by both tests comparable to the significance level

162 used.

163 In Figure 1 we show the fraction of rejections of both $L$ (one tail test) and $N$-tests (two tails test) on

164 100 ETAS pseudo-real catalogs at significance level 0.05, for daily and cumulative tests, and for

165 each time window $T_j$. The plots show that the proportions of rejections of $N$-test are above 30% (see

166 Figure 1a), much larger than the theoretical fraction (i.e., 5%). Similar results are found for the $L$-

167 test (see Figure 1b), computed on whole region, for which the fraction of rejections is above 20%.

168 In order to verify the spatial distribution of $L$-test failures we show in Figure 2 the maps of quantile

169 scores $\gamma_j^i$ for each time window. The figure shows that the failures are mainly near Landers and Big

170 Bear locations, where the number of events is larger and the spatial clustering is more evident.

171    The same analyses on Poissonian catalogues show that the fractions of rejections for both tests are

172    in perfect agreement with the significance level (0.05) adopted (see Figure 3).

173    To explain one of the possible reasons for this discrepancy, we report in Figure 4, the ratio between

174    mean and variance of the number of events recorded into 1000 synthetic ETAS catalogues,

175    simulated following the same rules used for the 100 pseudo-real catalogs (see Appendix A). This

176    ratio is much smaller than the unity, the value that characterizes the Poisson distribution (see Figure

177    4). This proves that the variability of the number of events is much larger than that expected in the

178    case of a Poisson process. By performing a Chi-squared test, the Poisson distribution is rejected for

179    all time-windows at a significance level of 0.01, and this independent of how the data are regrouped

180    to compare expected and observed distributions.

181    To quantify the differences between the variability of the seismic rate due to Poisson and ETAS

182    distributions, we plot in Figure 5 the differences of their 95% confidence bounds. Specifically, for

183    each pseudo-real ETAS catalog and for each day, we compute the variability of the seismic rate

184    $\Delta\alpha_{95\%}^{POISSON}$ expected by the Poisson distribution and assumed by CSEP tests; this value is compared

185    with the empirical variability $\Delta\alpha_{95\%}^{ETAS}$ of the ETAS distribution that has been calculated numerically

186    by the 1000 synthetic rates used for producing forecasts. Figure 5 shows the average of the

187    differences $\Delta\alpha_{95\%}^{ETAS} - \Delta\alpha_{95\%}^{POISSON}$ calculated for 100 pseudo-real catalogs. The positive differences

188    mean that the variability for the ETAS model is much larger than the variability of the Poisson

189    distribution. Interestingly, this difference decreases with time, implying that this difference becomes

190    less serious when the seismic rate tends to decrease.

191

192    **3.    Discussion and conclusions.**

193    In this paper we show that part of the CSEP/RELM testing procedure does not perform correctly for

194    a widely used class of models, i.e., the ETAS models. Specifically, by reproducing the CSEP

195    experiment on "pseudo-real" ETAS catalogs – for which we know the right model – we find that

196 the rejections are much more than expected. We identify one main reason for this deficiency: the

197 assumption that the number of earthquakes per bin has a Poisson distribution does not hold for

198 ETAS models. The latter have a variability of occurrences much larger than what predicted by a

199 Poisson distribution. The underestimation of the variability made under the Poisson hypothesis

200 unavoidably leads to a high rejection frequency during the CSEP experiments, at least for the ETAS

201 class of models. It is worth noting that a higher variability compared to what assumed by the

202 Poisson hypothesis is also observed on real catalogs (e.g., Saichev and Sornette, 2007; Kagan, 2009

203 and references therein) possibly (but not necessarily) leading to a wide generalization of the

204 conclusions reported in this paper (see also Schorlemmer et al., 2010). These results may be

205 generalized in this way: forecasting models that produce a higher variability of the seismic rates

206 compared to the Poisson process may be rejected too often also when they represent an accurate

207 representation of the observed spatio-temporal evolution of the seismicity. On the other hand, we

208 also foresee that forecasting models producing a variability of the seismic rates smaller than that

209 expected in the case of a Poisson process may be not rejected often enough even in case they do not

210 represent an accurate representation of the seismicity. Figure 2 shows also another interesting

211 departure from the Poisson distribution. Rejected bins appear clustered in space. The Poisson

212 distribution assumes that the seismic variability per bin is conditioned only by the seismic rate of

213 the model. Actually, the observed rate in a bin is also conditioned by the seismicity occurred in the

214 adjacent bins during the forecasting time window; this component is neglected in the testing phase

215 and it may play an important role on the results of the L-test.

216 In this paper we have investigated a strongly clustered sequence (pseudo real catalogs mimicking an

217 aftershock sequence) that is characterized by bins with a large number of events. In other cases,

218 such as the one-day forecasts during a quiet period or the forecast of large events (M≥5.0) in a 5-

219 year time period, the expected number of events is probably much smaller. In these cases, the bias

220 may be less serious as showed by Figure 1 (cf. the rejection rates for M3+ and M4+ events) and

221 Figure 5, and also as expected by the theory of hypothesis testing (basically, the fewer the data, the

222 more difficult is to reject an hypothesis).

223 Although these results indicate a bias of the current testing procedures of the CSEP experiments, we

224 stress that these experiment remain of paramount importance and they are unavoidable if we wish to

225 maintain earthquake forecasting in a scientific domain that requires formulation of hypothesis and

226 testing. The lesson to be learned is that some of the CSEP/RELM testing procedures should be

227 improved and/or implemented. Specifically, in order to get reliable results, we argue that the

228 CSEP/RELM suite of tests needs a significant revision. We identify three possible strategies that

229 could be implemented for current and future experiments:

230 1. Each forecasting model has to provide the likelihood function. This allows the likelihood tests to

231 be applied correctly because the Poisson assumption of the seismic rate variability is no longer

232 necessary, and other goodness-of-fit tests and skill measures may be applied, like the residuals

233 analysis (Ogata 1998; Marzocchi and Lombardi, 2009) and the Information Gain (e.g., Daley and

234 Vere-Jones, 2003). Notably, this approach would also avoid potential biases in the testing phase due

235 to the spatial correlation of the rejected bins (see figure 2). This is maybe the optimal choice from a

236 statistical point of view, but it is not applicable to models that do not have a likelihood function,

237 such as many pattern recognition algorithms.

238 2. The forecasts have to be described by a distribution of the expected number of earthquakes (see

239 also Werner and Sornette 2008), not by a single value as now. For example, the forecasts may be

240 composed by 1000 expected number of events, from which a central value and the dispersion can be

241 easily retrieved (see Marzocchi and Lombardi, 2009). This strategy is in principle applicable to

242 every model, but it would require a change in the CSEP procedures. In our mind, this option is

243 probably the easiest to implement for future experiments, but it is inapplicable to the present

244 forecasts that are composed just by one single expected number of earthquakes. Moreover, being

245 still based on binning forecasts, we remark that this strategy would not avoid possible biases

246 induced by the spatial correlation of the bins; a careful analysis of such potential bias is required.

247 3. The model-based variability of the number of earthquakes in each bin may be set by some

248 empirical rules that take into account the higher variability that characterize many models. This is

249 widely applicable for all models and all experiments so far completed or running, but certainly it

250 raises important technical problems. The first one is the introduction of a key parameter (i.e. the

251 dispersion) *after* the forecasts have been made. This would corrupt the prospective philosophy of

252 the experiments. Second, the choice of the empirical adjustment rule becomes critical for the

253 evaluation process. Unavoidably, this choice would raise a lot of debate about what is the best

254 adjustment rule, and if different rules should be applied to different models. In any case, it may be

255 difficult to establish these rules objectively and independently from the modelers.

256

257 **Data and Resources**

258 The Landers earthquake data were obtained from Southern California Earthquake Data Center,

259 website (http://data.scec.org/research/altcatalogs.html). The maps were made using the Generic

260 Mapping Tools (www.soest.hawaii.edu/gmt). The MATLAB GNU codes used in the present work

261 to run the N and L tests have been provided by the Southern California Earthquake Center CSEP

262 software development team.

263

269

270 **References**

271 Daley D.J. and D. Vere-Jones (2003), *An Introduction to the Theory of Point Processes*, Springer-
272   Verlag, New York, 2-nd ed., Vol. 1, pp. 469.

273    ICEF (International Commission on Eathquake Forecasting for Civil Protection) (2009).

274          Operational Earthquake Forecasting: State of Knowledge and Guidelines for Utilization.

275          Report for Italian Civil Protection .

276    Jordan, T.H. (2006). Earthquake Predictability:Brick by brick, *Seismol.Res.Lett.,* 77(1)*,* 3-6.

277    Kagan Y.Y. (2009). Statistical distributions of earthquake numbers: consequence of branching

278          process, *Geophys. J. Int*., submitted.

279    Kagan Y.Y., D.D. Jackson (1994). Long-term probabilistic forecasting of earthquakes, J. Geophys.

280          Res., 99, 13685-13700.

281    Kagan Y. Y., D. D. Jackson (1995). New seismic gap hypothesis: Five years after. J. Geophys. Res.,

282          100, 3943-3959.

283    Marzocchi W., L. Sandri, E. Boschi (2003). On the validation of earthquake-forecasting models: the

284          case of pattern recognition algorithms. *Bull. Seismol. Soc. Am*., 93, 1994-2004.

285    Marzocchi W., A.M. Lombardi (2009). Real-time forecasting following a damaging earthquake.

286          *Geophys. Res. Lett*., 36, L21302, doi:10.1029/2009GL040233.

287    Ogata, Y. (1998). Space-Time Point-Process Models for Earthquake Occurrences, *Ann. Inst. Statist.*

288          *Math.* 50(2), 379-402.

289    Saichev A., D. Sornette (2007). Power law distribution of seismic rates, *Tectonophysics,* 431*,* 7-13

290    Schorlemmer, D., and M. C. Gerstenberger (2007). RELM testing center. *Seismol.Res.Lett.,* 78(1)*,*

291          30-36.

292    Schorlemmer, D., M.C. Gerstenberger, S. Wiemer, D.D. Jackson and D.A. Rhoades (2007).

293          Earthquake Likelihood Model Testing, *Seismol.Res.Lett.,* 78(1)*,* 17-29.

294    Schorlemmer D., J.D. Zechar, M.J. Werner, E.H. Field, D.D. Jackson, T.H. Jordan (2010). First

295          results of the Regional Earthquake Likelihood Models experiment, *Pure Applied Geophys*.,

296          in press.

297    Werner M. J., and D. Sornette (2008), Magnitude uncertainties impact seismic rate estimates,

298        forecasts, and predictability experiments, *J. Geophys. Res.* 113, B08302,

299        doi:10.1029/2007JB005427.

300    Zechar J.D., T.H. Jordan (2008), Testing alarm-based earthquake predictions, *Geophys. J. Int.*, 172,

301        715-724.

302    Zechar J.D., D. Schorlemmer, M. Liukis, J. Yu, F. Euchner, P.J. Macheling, and T.H. Jordan

303        (2009), The Collaboratory for the Study of Earthquake Predictability perspective on

304        computational earthquake science, Concurrency and Computation: Practice and Experience,

305        doi: 10.1002/cpe.1519

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

**Figure Captions**

**Figure 1:** Fractions of rejections of the daily $L$ and $N$-test on 100 pseudo-real ETAS catalogs.

**Figure 2:** Spatial distribution of fractions of rejections on 100 pseudo-real ETAS catalogs for $L$-test conducted on 10 time windows.

**Figure 3:** The same of Figure 1 but for pseudo-real Poisson catalogs.

**Figure 4:** Ratio between mean and variance of events recorded in 1000 ETAS pseudo-real catalogs for 10 time windows.

**Figure 5:** Difference between the 95% confidence intervals of the ETAS and Poisson distributions as a function of the forecasting time window; each point represents the average of the differences calculated for the 100 ETAS pseudo-real catalogs used for $L$ and $N$-tests.

# APPENDIX A. Generating the pseudo-real synthetic catalogs

In this appendix, we report the strategy adopted to generate ETAS and Poissonian pseudo-real catalogs.

The total space-time conditional intensity $\lambda(t,x,y/\mathcal{H}_t)$ of the ETAS model (i.e. the probability of an earthquake occurring in the infinitesimal space-time volume conditioned to all past history) is defined by equation:

$$\lambda(t,x,y,m/\mathcal{H}_t) = \left[ \nu u(x,y) + \sum_{t_i<t} \frac{K}{(t-t_i+c)^p} e^{\alpha(M_i-M_c)} \frac{c^i_{d,q,\gamma}}{\left[r_i^2 + \left(de^{\gamma(M_i-M_c)}\right)^2\right]^q} \right] \beta e^{\beta(m-M_c)} \qquad (A1)$$

where $\mathcal{H}_t = \{(t_i,x_i,y_i,M_i); \; t_i < t\}$ is the observation history up time $t$, $M_c$ is the completeness magnitude of the catalog, $u(x,y)$ is the spatial probability density function (PDF) of background events, $c^i_{d,q,\gamma} = \frac{q-1}{\pi} [(de^{\gamma(M_i-M_c)})^2]^{q-1}$ is the normalization constant of the spatial PDF for triggered events, and $r_i$ is the distance between location $(x,y)$ and the epicenter of $i$-th event $(x_i,y_i)$ (Lombardi et al., 2009). Finally $\beta = b \cdot ln(10)$ is the parameter of the well-known Gutenberg-Richer Law (Gutenberg and Richter, 1954), assumed as distribution for magnitude of all events.

The set of parameters $\Theta = (\nu,K,c,p,\alpha,d,q,\gamma,\beta)$ of the model, for the events occurred within a time interval $[T_1,T_2]$ and a region $R$, can be estimated by maximizing the log-likelihood function (Daley and Vere-Jones, 2003), given by

$$logL(\Theta) = \sum_{i=1}^{N} log\,\lambda(t_i,x_i,y_i,m_i/H_{t_i}) - \int_{T_1}^{T_2}\int_R\int_{M_c}^{M_{max}} \lambda(t,x,y,m/H_t)\,dtdxdydm \qquad (A2)$$

A careful method to obtain the best parameters of the model is the iteration algorithm developed by Zhuang et al. (2002), providing also an estimation of the PDF $u(x,y)$ for background events.

Our pseudo-real ETAS catalogs are simulated in agreement with the ETAS model estimated for the region hit by the Landers earthquake. Specifically we use the relocated data set (Hauksson and Shearer, 2005) recorded by the California Institute of Technology/U.S. Geological Survey (CIT / USGS) Southern California Seismic Network and available at the SCEDC (Southern California Earthquake Data Center) website (http://data.scec.org/research/altcatalogs.html). We consider earthquakes with a depth less than 30 km and a magnitude above 3.0, occurred from Jan 1 1984 to Dec 31 2004 and located in the region [-119.0°W/32.5°N − -115.0°W/36.5°N] (5757 events). The parameters estimated by using the procedure proposed by Zhuang et al. (2002) are listed in Table

378   A1. We perform simulations by including in the past history the real observed seismicity above

379   magnitude 3.0, occurred before July 1 1992, 3 days after the $M_L7.3$ Landers mainshock. In this way

380   we take into account knowledge coming from the initial phase of the sequence, including also the

381   $M_L6.4$ Big Bear aftershock.

382   We simulate the Poisson pseudo-real catalogs by imposing a rate of 60 day$^{-1}$ and adopting the PDF

383   $u(x,y)$, estimated for the ETAS model, for the spatial distribution of events. All pseudo-real catalogs

384   recover a time period of 10 days. We remark that we intend to perform simulations by reproducing

385   the type of forecasts usually tested in CSEP laboratories, no matter the specific region or time

386   period we consider.

387   In order to verify the reliability of our pseudo-real catalogs, we analyze their residuals. The residual

388   analysis is a common diagnostic technique for stochastic point processes based on transformation of

389   the time axis $t$ into a new scale $\tau$ by the increasing function

390
$$\tau = \Lambda(t) = \int_{T_{start}}^{t} dt \int_{R} dxdy \int_{M_c}^{M_{max}} dm\ \lambda(t,x,y,m/\mathcal{H}_t) \qquad (A3)$$

391   where $T_{start}$ is the starting time of the observation history $H_t$ (Ogata, 1998). The random variable $\tau$

392   represents the expected number of occurrences in time period $[T_{start}, t]$. If a model with conditional

393   intensity $\lambda(t,x,y,m/\mathcal{H}_t)$ describes the temporal evolution of the process, the transformed data $\tau_i$

394   $=\Lambda(t_i)$, known in statistical seismology with the name of *residuals,* are expected to behave like a

395   stationary Poisson process with the unit rate (Ogata, 1998); i.e. the values $\Delta\tau_i = \tau_{i+1}-\tau_i$ are

396   independent and exponentially distributed (with mean equal to 1) random variables. We check this

397   hypothesis for residuals by means of two nonparametric tests: the Runs test, to verify the reliability

398   of the independence property, and the one-sample Kolmogorov-Smirnov (KS1) test, to check the

399   standard exponential distribution (Gibbons and Chakraborti, 2003; Lombardi and Marzocchi, 2007).

400   Specifically the Runs-test can be used to test if a process is not auto-correlated and consists in

401   testing the randomness of runs, i.e. of uninterrupted subsequences of values above or below the

402   mean (see Gibbons and Chakraborti, 2003; Lombardi and Marzocchi, 2007 for details). We use

403   both tests because all goodness-of-fit tests (as KS1) are ineffective to check the presence of a

404   memory in the time series. Hence, any discrepancy of residuals by Poisson hypothesis, identified by

405   just one or both tests, is a sign of inadequacy of ETAS model to explain all basic features of

406   analyzed seismicity. We stress that this check analysis is similar to the RELM/CSEP $N$-test. As the

407   $N$-test, it consists in a comparison between the observed and the expected total number of events

408   and it is directed to highlight under or over-prediction. On the other side the residual analysis does

409   not need the discretization of the temporal scale in time bins. As explained along the text, this is a

410   crucial point of RELM/CSEP tests. In Figure A1 we show the empirical cumulative function of p-

values of KS1 and Runs tests, for the 100 pseudo-real ETAS catalogs, together with the 99% confidence bounds. The confidence level is calculated assuming that for each point of the curve the expected fraction of rejection is given by the p-value reported on the x-axis, and the variability (1 sigma) is given by $\sqrt{p(1-p)/N}$. Note that, for both tests the cumulative distribution is inside the 99% confidence interval.

**References**

Daley, D.J. and D. Vere-Jones (2003). *An Introduction to the Theory of Point Processes*, Springer-Verlag, New York, 2-nd ed., Vol. 1, pp. 469.

Gibbons, J.D. and S. Chakraborti (2003). Non-parametric Statistical Inference, 4th ed., rev. and expanded, New York: Marcel Dekker, 645 pp.

Gutenberg, B. and C.F. Richter (1954). *Seismicity of the Earth and Associated Phenomena*, Princeton, pp. 273.

Hauksson, E. and P. Shearer (2005). Southern California Hypocenter Relocation withWaveform Cross-Correlation, Part 1: Results Using the Double-Difference Method, *Bull. Seismol. Soc. Am.*, **95 (3)**, 896–903, doi:10.1785/0120040167.

Lombardi, A. M., and W. Marzocchi (2007). Evidence of clustering and nonstationarity in the time distribution of large worldwide earthquakes, *J. Geophys. Res.*, **112**, B02303, doi:10.1029/2006JB004568.

Lombardi, A.M., M. Cocco and W. Marzocchi (2009). On the increase of background seismicity rate during the 1997-1998 Umbria-Marche (central Italy) sequence: apparent variation or fluid-driven triggering? Submitted to *Bull. Seismol. Soc. Am.*

Ogata, Y. (1998). Space-Time Point-Process Models for Earthquake Occurrences, *Ann. Inst. Statist. Math.* **50(2)**, 379-402.

Zhuang, J., Y. Ogata and D. Vere-Jones (2002). Stochastic declustering of space-time earthquake occurrence, *J. Am. Stat. Assoc.,* **97**, 369-380.

438 **Figure Captions**

439 **Figure A1:** Cumulative function of the empirical p-values (solid black lines) for KS1 (panel a) and

440 RUNS (panel b) Test applied to Residuals of 100 simulated ETAS catalogues. Dashed gray lines

441 mark the 99% confidence bounds.

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

| Parameter | Value |
|---|---|
| $\nu$ | $0.10 \pm 0.004$ (day$^{-1}$) |
| K | $0.043 \pm 0.002$ (day$^{p-1}$) |
| p | $1.20 \pm 0.01$ |
| c | $0.030 \pm 0.004$ (day) |
| $\alpha$ | $1.20 \pm 0.03$ (mag$^{-1}$) |
| d | $0.30 \pm 0.01$ (km) |
| q | $\equiv 1.5$ |
| $\gamma$ | $0.60 \pm 0.03$ (mag$^{-1}$) |
| Log-likelihood | -21277.5 |

472

473 **TableA1:** Maximum Likelihood parameters (with relative errors) and log-likelihood of ETAS
474 model for Landers region seismicity [ -119.0° W/32.5° N − -115.0° W/36.5° N]
475 ($M_c$ = 3.0; Jan 1 1984 – Dec 31 2004; 5757 events)
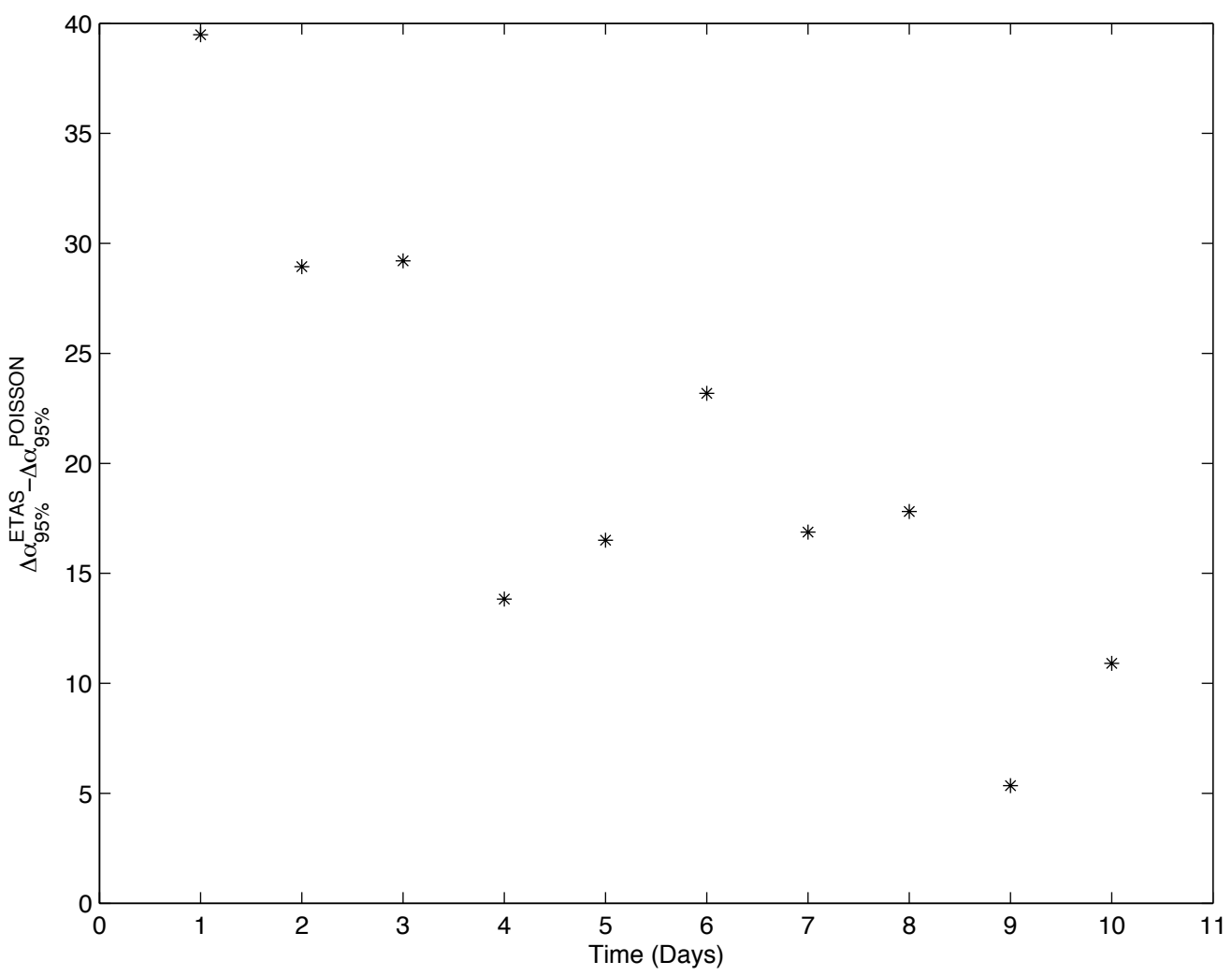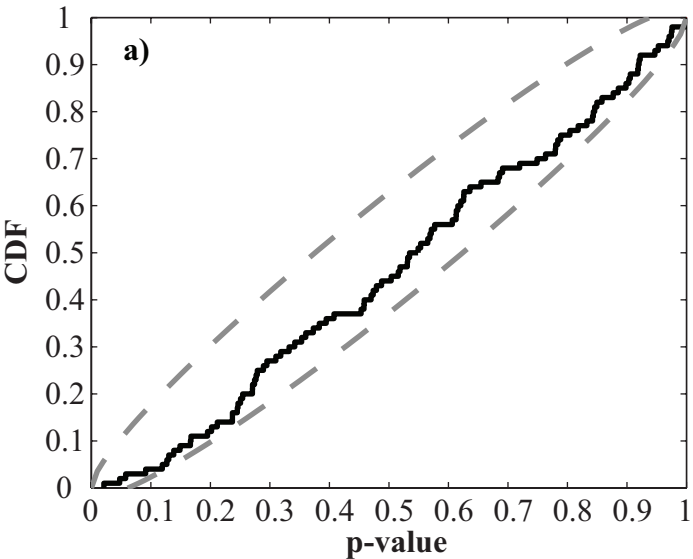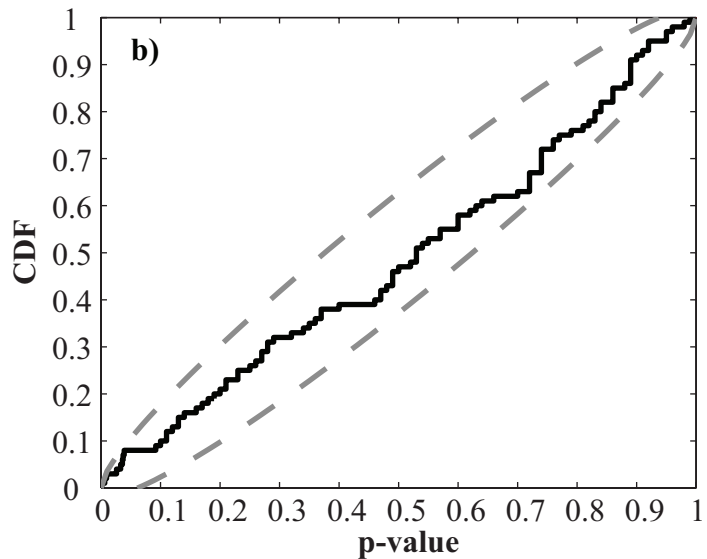
476

477

478

479

480

481

482

Figure 1

Figure 2

Figure 3

Figure 4

**KS1 TEST**

**RUNS TEST**

Figure A1