



národní
úložiště
šedé
literatury

Coenocline reconstruction using graph theory and Bayesian probability data generator

Čejchan, Petr
2007

Dostupný z <http://www.nusl.cz/ntk/nusl-55874>

Dílo je chráněno podle autorského zákona č. 121/2000 Sb.

Tento dokument byl stažen z Národního úložiště šedé literatury (NUŠL).

Datum stažení: 09.08.2024

Další dokumenty můžete najít prostřednictvím vyhledávacího rozhraní nusl.cz .

Coenocline reconstruction using graph theory and Bayesian probability data generator

Peter A. Cejchan, Institute of Geology, Academy of Sciences of the Czech Republic, Rozvojova 269, CZ 16502 Prague, <cej@gl.cas.cz>, fax +420-2-20922670

19th February 2007

Keywords: paleobiology, paleoecology, indirect gradient analysis, coenocline, continuum concept, gradient reconstruction, fossil assemblages, Bayesian inference, population density, abundance, species performance, Quadratic Assignment Problem.

CopyLeft 2007 by Peter Cejchan. All rights reserved. The right to make & distribute verbatim copies of this document is granted hereby, provided that this copyleft statement is included.

There are three kinds of lies: lies, damned lies, and statistics.
Mark Twain, 1924

1 The problem

The goal is to design a novel way to reconstruct the 'species response function (SRF)' of unitary organisms on an undimensional environmental/ (time gradient solely from abundance data on species communities (indirect gradient analysis [?]). The word 'indirect' means that the gradient itself is latent, is not accessible, or even is presumed. The study will focus on the topology of the problem, thus disregarding the quantitative measuring of 'sample scores'. Formally, this task can be termed 'seriation of fossil communities (assemblages)'.

Bayesian inference is used here to 'squeeze' the population density distributions out of observed species' abundances. A large number of generated hypotheses serve to derive probability density functions of sample ranks, which are the result of the analysis.

2 Seriation

We use the term 'seriation' here simply for topological (preserving just order, not distances) ordination in one dimension, usually, but not necessarily environment, or time. Seriation is founded upon the assumption of humped unimodal 'response functions'. The task might seem trivial at first glance: permute rows of an abundance matrix until the humped pattern appears in columns. Unfortunately, the number of permutations increases with $n!/2$, where n is the number of samples. Thus, having to seriate 20 samples by exhaustive search means to evaluate 121645100488320000 different permutations.

3 Need for an improved method

It appears that there is up to now no method in use, which would take the nonlinearity of species response functions into account from the very beginning. Instead, methods developed for the linear model are used, their misapplication being 'cured' by subsequent 'detrending', in the best case. A need for an improved method is thus obvious.

4 Assumptions

- We will consider only unitary organisms which do not form clusters, spread the progeny using (actively or passively) moving propagules, and having defined, unchanging demographic functions.
- We will adopt the population density as a feasible measure of species' performance.
- We require that a defined proportion of the population gets into the sediment, where it fossilizes. Actually, we will hide several processes under the term "fossilization chance" (see below).
- We will suppose the combined effect of these processes to be 'unpredictable' enough to form a stochastic variate that is characteristic for a given environment.
- We will suppose that sampling is done by taking a constant (up to a measuring error) volume of the sediment, extracting the fossils out of it, and counting individuals belonging to different species.
- Due to the complexity of the problem, we will assume that undimensional analysis is appropriate. Although often unrealistic oversimplification, environments governed by an univariate gradient may occur quite frequently (cf. typical zonal communities of the mountain vegetation, inter-tidal biota, vegetation on moisture gradient, stratified oceanic communities depending on bathymetry, etc.).
- We assume that the typical 'species response function' on an environmental/time gradient is a unimodal function.

5 Processes contributing to abundance

5.1 Population density

Species performance, as determined by the environmental factor(s) and inter-specific competition, maps onto the population density, which is the essential component of the abundance model. There are also other measures of performance, like biomass production (standing crop), and other. However the present study deals with unitary organisms and their counts only; thus we confine to density as a performance measure. The population density is then used as the dependent variable of the SRF.

5.2 Volume sampling

Ideally, we would like to take samples representing constant *area × time*. We assume we are able to take samples, or 'snapshots', of the communities that represent constant *volume* of the sediment, within the sampling error $N(\mu, \sigma)$. Equation 5 shows the relation of the sample volume to time, using sedimentation rate, and sample basal area. Equation 16 incorporates the sample volume into a probabilistic estimate of abundance, y .

Wherever abundance (and hence, density) is estimated by means of counting individuals (a discrete variable), we are to deal with a so-called 'Poisson point process' [?] [?] [?] [?], where the number of individuals in constant amount of sampled *area × time* has a Poisson distribution [?]. Unfortunately, this Poisson sampling process is a source of a substantial error in abundance estimate, unless the samples are very large.

5.3 Sedimentation rate

Sedimentation rate influences the abundances via mapping of time to sediment thickness, i.e., also the sample volume. Thus, it is critical not only to take care of the volume of all samples be (as much as possible) equal, but also the sedimentation rate to be known, together with its variance.

5.4 Fossilization

Several processes influence the chance to fossilize. Predation removes living organisms from the environment, and from the reproduction process, modifying the ecological SRF of the species in a manner similar to competition. If the individual avoids predation, its remains must quickly get buried by the sediment, in order to have a nonzero chance to fossilize. The organism remains not buried quickly are subject to decay and wash-away, that does not let the remain to enter the sediment. Even if buried, erosion, or dissolution of the sediment removes the remains together with the sediment. After surviving all the previous processes, the organismal remains must cope with bioturbation, and dissolution, in order not to be removed from the fossil record.

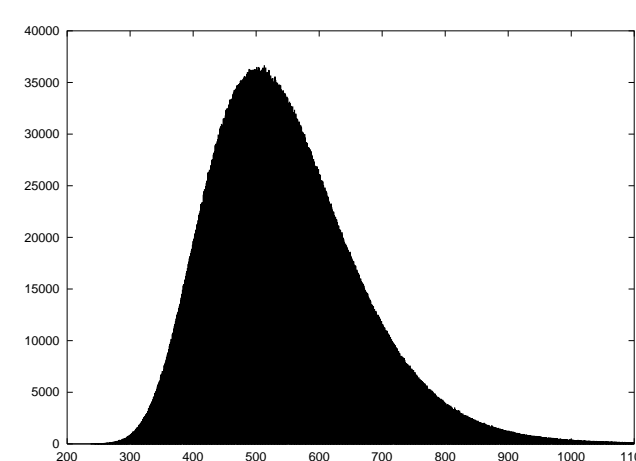
5.5 Contamination and reworking

Although perhaps far too conservative, these effects are modeled here as a constant 'contamination intensity' to allow for nonzero abundances to be sampled in arbitrarily large distance from the species' optimum. This 'contamination intensity' is modeled here as dependent on the species' mean performance; thus, for curious and common species it is high, whereas for stenocratic and rare species it is considered low. Although this model is rough, its ability to allow for substantial nonzero densities far from the optimum make it a suitable one for the purpose of testing the gradient's reconstructability under wide spectrum of conditions.

Reworking of fossils from older strata may produce, if not recognized, a similar effect as contamination does. However, the resistance of fossils to reworking phenomena, not directly the species commonness / scarcity, plays a key role here. Even a prolific species may easily be destroyed during the process(es) of reworking; thus, a reworked scarce species may finally become more abundant than is a reworked common species.

5.6 Generation length

Generation length maps density and time onto abundance. The lower is the generation length of given species, the higher is the abundance compared with another species with the same population density, but longer length of one generation.



Volume sampled constant density process with varying sedimentation rate, fossilization, contamination, and varying length of generation. Parameter setting was $\mu = 0.001$, $\sigma_m = 0.0004$, $r_1 = 0.008$, $r_2 = 0.012$, $c_1 = 0.45$, $c_2 = 0.55$, $\alpha_1 = 0.01$, $\alpha_2 = 0.1$, $\nu_1 = 1.00$, $\sigma_v = 0.15$.

6 Deterministic model

Population densities are measured on the unit area. Density is here vaguely defined as a number of individuals encountered in a "time cut" on a unit area. Total number of individuals (of the species under consideration) that inhabited some area during a time interval, is

$$e = \frac{dta}{g} \quad (1)$$

where e is the total number of individuals encountered in a time interval on an area, d is a 'realized' density, t is time interval within which the sampled volume of sediment was formed, a is the basal area of the sample, and g is duration of one generation for the given species. We expect that g is constant for the given species. The model should account also for a part of density caused by 'alien' individuals. We model the contamination as proportional to both contamination intensity, and mean population density of the proper species

$$d = d_p + \alpha \cdot d_m \quad (2)$$

where d_p is 'pure' population density, d_m is mean population density of given species across all samples, and α is the contamination intensity. Then,

$$e = \frac{dta}{g} = \frac{(d_p + \alpha d_m)ta}{g} \quad (3)$$

Usually, we cannot measure time directly, instead, we can estimate it from the volume of the sediment that has deposited under varying sedimentation rate. Thus,

$$v = ah \quad (4)$$

$$r = \frac{h}{t} \quad (5)$$

$$t = \frac{h}{r} = \frac{v}{ar} \quad (6)$$

$$e = \frac{(d_p + \alpha d_m)ta}{g} = \frac{(d_p + \alpha d_m)va}{arg} = \frac{(d_p + \alpha d_m)v}{rg} \quad (7)$$

where v is the volume of the sediment, h is height (thickness) of the sedimentary column, r is the sedimentation rate, t is time under which the sampled volume of sediment was formed. Unfortunately, we usually cannot arrange sampling so that the sample represents a constant time interval. Fossilization chance c is a proportion of individuals which actually fossilizes. Includes predation, gaps in sedimentation, and removal from the sediment due to any conceivable process. Observed abundance y is then:

$$y = ec = \frac{(d_p + \alpha d_m)vc}{rg} \quad (8)$$

However the above variables are not deterministic, rather stochastic. With unchanged connotation we will now proceed to the probabilistic (stochastic) model.

7 Probabilistic model

Density, d , is a Poisson variate (output of a Poisson process of "colonization", with intensity d_p of the area under consideration). The sample volume, v , although tried to be kept constant, is subject to the sampling error, which has (approximately) a Normal distribution (this is not clear, as Gaussian distribution was derived for location, not scale variables; cf. Sivia 1996, p. 112). Sedimentation rate, r , is unknown, but supposed to lie within the limits r_1, r_2 within one sample. Fossilization chance, c , is unknown as well, but is supposed to lie within the limits c_1, c_2 within one sample. Contamination intensity, α , is unknown as well, but is supposed to lie within the limits α_1, α_2 for all samples and species. We model r and c by a uniform distribution as well. The length of generation for a given species is supposed to be known with certain uncertainty, which has (approximately) a Normal distribution (see above). Thus we put:

$$d = P_o(d_p) + P_o(d_m \cdot U(r_1, \alpha_1)) \quad (9)$$

$$e = P_o \left(\frac{(d_p + \alpha d_m \cdot U(r_1, \alpha_1)) \cdot v}{rg} \right) \quad (10)$$

$$y = B_i(c, e) \quad (11)$$

$$v = N(\mu_v, \sigma_v) \quad (12)$$

$$r = U(r_1, r_2) \quad (13)$$

within sample,

$$c = U(c_1, c_2) \quad (14)$$

within sample,

$$o = U(o_1, o_2) \quad (15)$$

$$g = N(\mu_g, \sigma_g) \quad (16)$$

within as well as among samples. From 10 we have:

$$y = B_i \left(P_o \left(\frac{(d_p + \alpha d_m \cdot U(r_1, \alpha_1)) \cdot N(\mu_v, \sigma_v)}{N(\mu_g, \sigma_g) \cdot U(r_1, r_2)} \right) \cdot U(c_1, c_2) \right) \quad (17)$$

where $P_o(\lambda)$ is a Poisson variate with parameter λ , $B_i(r, m)$ is a Poisson variate with parameters r, m , $N(\mu, \sigma)$ is a Normal (Gaussian) variate with mean μ , and standard deviation σ , and $U(l, k)$ is a uniform variate on the interval (l, k) .

Estimating density

Abundance is a result of a compound process, with density being only one of its parameters. It should be stressed here again that abundance and density is not the same, and the two terms should not be mixed! Only density has an ecological meaning; however, for the paleobiologist, it is accessible only via abundance.

We shall use abundance pdf as the likelihood function for estimation of density posterior pdf from Eq. ???. Thus, we need to derive the abundance pdf first. We will use Eq. 16; hence, using product rule (Eq. ??):

$$prob(y|c, e, I) = \frac{e!}{y!(e-y)!} e^y (1-e)^{e-y} \quad (18)$$

$$prob(y, c|e, I) = prob(y|c, e, I) \times prob(c|I) = \frac{e!}{y!(e-y)!} e^y (1-e)^{e-y} \times \frac{1}{c_2 - c_1} \quad (19)$$

$$prob(y, c, e|o, v, r, I) = prob(y, c|e, I) \times \frac{\left(\frac{d_p + \alpha d_m \cdot \alpha}{rg} \right)^e \exp\left(-\frac{d_p + \alpha d_m \cdot \alpha}{rg}\right)}{e!} \quad (20)$$

$$prob(y, c, e, o, v, r, g, I) = prob(y, c, e|o, v, r, g, I) \times \frac{1}{o_2 - o_1} \quad (21)$$

$$prob(y, c, e, o, v, r, g, I) = prob(c, e, o|v, r, g, I) \times \frac{1}{\sigma_v \sqrt{2\pi}} \exp\left(-\frac{(v - \mu_v)^2}{2\sigma_v^2}\right) \quad (22)$$

$$prob(y, c, e, o, v, r, g, I) = prob(y, c, e, o, v, r, g, I) \times \frac{1}{r_2 - r_1} \quad (23)$$

$$prob(y, c, e, o, v, r, g, I) = prob(y, c, e, o, v, r, g, I) \times \frac{1}{\sigma_g \sqrt{2\pi}} \exp\left(-\frac{(g - \mu_g)^2}{2\sigma_g^2}\right) \quad (24)$$

From marginalization we have:

$$prob(y, c, e, o, v, r | I) = \int_0^{+\infty} prob(y, c, e, o, v, r, g | I) dg \quad (25)$$

$$prob(y, c, e, o, v | I) = \int_0^{+\infty} \int_0^{+\infty} prob(y, c, e, o, v, r | I) dr \quad (26)$$

$$e = \sum_{i=1}^m \sum_{j=1}^m \delta_{ij} \cdot \varphi_{\pi(i) \cdot j} \quad (27)$$

$$\lambda = d_p + \frac{\alpha d_m}{\sigma_g} \cdot \frac{v}{g} \quad (28)$$

so that finally we have

$$prob(y | I) = \int_0^{+\infty} \int_{r_1}^{+\infty} \int_0^{+\infty} \int_0^{+\infty} \int_0^{+\infty} \int_0^{+\infty} prob(y, c, d, o, v, r, g | I) dg dr do dd dd dc \quad (29)$$

From the Bayes' theorem, as given in Eq. ??:

$$prob(hypothesis | data, I) \propto prob(data | hypothesis, I) \times prob(hypothesis | I) \quad (30)$$

Using Eq. 24 as the likelihood, and Jeffreys' prior (Eq. ??) for abundance, we get the posterior pdf for d_p :

$$prob(d_p | y, I) \propto \int_0^{+\infty} \int_0^{+\infty} \int_0^{+\infty} \int_0^{+\infty} \int_0^{+\infty} \int_0^{+\infty} prob(y, c, d, o, v, r, g | I) dg dr do dd dd dc \times \frac{1}{y} \quad (31)$$

Reconstruction

We have designed the measure that quantifies a 'compactness' of a pattern created by a series of SRFs, and we will try to show that it approximates, or quantifies, unimodality quite well. This is because the unimodal pattern is 'compact' in the sense that high values tend to be clustered together, whilst low values surround them.

Formalizing the task we get an objective function:

$$c = \sum_{i=1}^m \sum_{j=1}^m \delta_{ij} \cdot \varphi_{\pi(i) \cdot j} \cdot \alpha_{ij} \quad (32)$$

thus

$$c = \sum_{i=1}^m \sum_{j=1}^m \delta_{ij} \cdot \varphi_{\pi(i) \cdot j} \quad (33)$$

where

$$\varphi_{\pi(i) \cdot j} = \sum_{k=1}^m \delta_{ij} \cdot \varphi_{\pi(i) \cdot k} \cdot \beta_{\pi(j) \cdot k} \quad (34)$$

thus

$$\varphi_{ij} = \sum_{k=1}^m \delta_{ij} \cdot \beta_{jk} \quad (35)$$

However, Eq. 27 is an instance of a notoriously known graph-theoretical problem, the Quadratic Assignment Problem (QAP), where c is the cost function, $|\delta|$ is the distance matrix, $|\varphi|$ is the flow matrix, $\pi(i)$ is the position of the vertex i in the permutation π , and β_{jk} is the performance of species k in sample j . There exists a vast literature concerning the QAP.

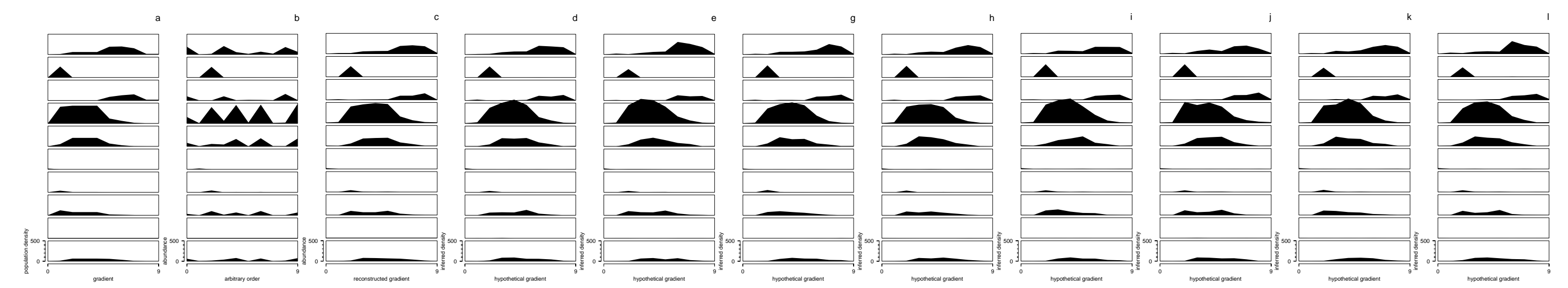
Efficiency

To assess the efficiency of the QAP method of seriation, a series of comparisons with several traditional methods have been run. The QAP method appeared superior to the traditional ones in a well-informed case: densities were a priori known, and the gradient was sampled using regular spacing of samples. The testing array consisted of 10 samples, spaced regularly on the gradient, composed of up to 200 species.

8 Results

8.1 Artificial example

To test the whole procedure, an artificial example has been generated. Gaussian model with randomly spaced samples has been used. Thirty independent estimates were generated for density of every species at each site. From this array, 62 matrices of inferred Bayesian densities (cases 000 to 061) were generated. The first 62 cases produced a handful of hypotheses, which are presented below.



Artificial example a, simulated coenocline, densities along the gradient; b, abundances sampled at sampling sites (a); c, coenocline and the gradient reconstructed from sampled abundances using the QAP algorithm; d - h, several hypotheses that occurred first when running probabilistic QAP seriation on inferred densities; note that the coenocline pattern generally holds while there are minor changes in sample ordering (the seriation of the samples) between hypotheses.

8.2 Case study

Data for this case study are derived from the fauna of land snails of the Holocene locality of Za krizem Cave near Svaty Jan pod Skalou, Bohemia. Age of the fauna is Holocene. The locality and the malacofauna was described by Lozek & Horáček 1993 [?], where also the quantitative data on fossils' abundances come from.

LAYER	SPECIES SAMPLE	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100	101	102	103	104	105	106	107	108	109	110	111	112	113	114	115	116	117	118	119	120	121	122	123	124	125	126	127	128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143	144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175	176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191	192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207	208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223	224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239	240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255	256	257	258	259	260	261	262	263	264	265	266	267	268	269	270	271	272	273	274	275	276	277	278	279	280	281	282	283	284	285	286	287	288	289	290	291	292	293	294	295	296	297	298	299	300	301	302	303	304	305	306	307	308	309	310	311	312	313	314	315	316	317	318	319	320	321	322	323	324	325	326	327	328	329	330	331	332	333	334	335	336	337	338	339	340	341	342	343	344	345	346	347	348	349	350	351	352	353	354	355	356	357	358	359	360	361	362	363	364	365	366	367	368	369	370	371	372	373	374	375	376	377	378	379	380	381	382	383	384	385	386	387	388	389	390	391	392	393	394	395	396	397	398	399	400	401	402	403	404	405	406	407	408	409	410	411	412	413	414</
-------	----------------	---	---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-------