# Probabilistic positional cross-identification of catalogs of astrophysical sources: the Aspects code

## Michel Fioc[1,2]

[1]Institut d'astrophysique de Paris

[2]Université Pierre et Marie Curie

EGSG, 2014-12-12

# The cross-identification problem

Consider two catalogs of astrophysical sources

$$K = \{M_1, \ldots, M_n\} \quad \text{and} \quad K' = \{M'_1, \ldots, M'_{n'}\}$$

defined on a common surface of area $S$.

How can one decide, just from the positions and positional uncertainties of $K$-sources (events $c_1, \ldots, c_n$) and $K'$-sources (events $c'_1, \ldots, c'_{n'}$), whether

- $M_i$ is identical to ("is associated with") $M'_j$,                    (event $A_{i,j}$)

or

- $M_i$ has no counterpart in $K'$?                                  (event $A_{i,0}$)

# Probability distribution of the position of a source

Assume that the observed positions of $M_i$ and $M'_j$, $\vec{r}_i$ and $\vec{r}'_j$, are normally distributed around their true positions, $\vec{r}^0_i$ and $\vec{r}'^0_j$:

$$\vec{r}_i \sim \mathcal{N}(\vec{r}^0_i, \Gamma_i) \quad \text{and} \quad \vec{r}'_j \sim \mathcal{N}(\vec{r}'^0_j, \Gamma'_j).$$

If $M_i$ and $M'_j$ are the same point source, their true positions are identical. Then,

$$\vec{r}_{i,j} = \vec{r}'_j - \vec{r}_i \sim \mathcal{N}(0, \Gamma_{i,j}),$$

where $\Gamma_{i,j} = \Gamma_i + \Gamma'_j$, and $\Gamma_i$ and $\Gamma'_j$ are given in the same basis (warning: near the poles, the North direction for $M_i$ differs from that for $M'_j$), i.e.

$$P(c_i \mid c'_j \cap A_{i,j}) = \frac{\exp(-\frac{1}{2}\, \vec{r}^{\mathsf{T}}_{i,j} \cdot \Gamma^{-1}_{i,j} \cdot \vec{r}_{i,j})}{2\,\pi\,\sqrt{\det \Gamma_{i,j}}} = \xi_{i,j}.$$

(For extended sources, one may add to $\Gamma_{i,j}$ a term depending on unknown parameters to account for the possible difference between the true positions.)

If $M_i$ has no counterpart, it is randomly distributed in $S$:

$$P(c_i \mid A_{i,0}) = \frac{1}{S} = \xi_{i,0}.$$

# Naïve "answer"

If $\Gamma_i = \sigma^2$ for all $K$-sources and $\Gamma'_j = \sigma'^2$ for all $K'$-sources,

1. define $R =$ some factor $\times \sqrt{\sigma^2 + \sigma'^2}$;
2. if there is a source $M'_j$ in the disk of radius $R$ centered on $M_i$, then $M'_j$ is the counterpart of $M_i$;
3. if not, $M_i$ has no counterpart.

## Weaknesses

▸ The factor is somewhat arbitrary.

▸ There may be more than one $K'$-source in the disk around $M_i$: which one is the counterpart? The closest one?

▸ $M'_j$ may be the closest source to $M_i$, but another $M_k$ may be the closest source to $M'_j$.

▸ The higher the density of $K'$-sources, the more likely that an unrelated one will be close to $M_i$.

▸ If positional uncertainties are elliptical, there are different ellipses for each $(M_i, M'_j)$: $M'_j$ may be considered as a counterpart for the ellipse defined by $(M_i, M'_j)$, but not for that defined by $(M_i, M'_k)$.

▸ Positional uncertainties are not always known.

# A probabilistic answer

What is really wanted is the *probability* that $M_i$ is associated with $M'_j$ ($j > 0$), or that $M_i$ has no counterpart in $K'$ ($j = 0$), given the positions of all the $K$- and $K'$-sources and the uncertainties on these,

$$C = c_1 \cap \ldots \cap c_n \qquad \text{and} \qquad C' = c'_1 \cap \ldots \cap c'_{n'},$$

i.e.

$$P(A_{i,j} \mid C \cap C').$$

## Unknown parameters

This probability depends on at least one unknown: the *a priori* probability (not knowing $C$ and $C'$) that any $K$-source has a counterpart in $K'$,

$$f = P\Big(\bigcup_{j>0} A_{i,j}\Big) = 1 - P(A_{i,0}).$$

(We will also use the *a priori* probability that any $K'$-source has a counterpart in $K$,

$$f' = P\Big(\bigcup_{i>0} A_{i,j}\Big) = 1 - P(A_{0,j}).)$$

$P(A_{i,j} \mid C \cap C')$ may also depend on other unknowns, such as the positional uncertainty in one catalog, or the combined uncertainty of both.

# Possible assumptions on associations

To compute the probabilities, some model of association must be assumed:

**Several-to-one:** several $K$-sources may be associated with the same $K'$-source, but at most one $K'$-source is associated to a given $K$ source. More precisely,

$$\begin{cases} \text{for all } M_i, \text{ the events } (A_{i,j})_{j\in[\![1,n']\!]} \text{ are exclusive;} \\ \text{for all } M'_j, \text{ the events } (A_{i,j})_{i\in[\![1,n]\!]} \text{ are independent.} \end{cases} \qquad (H_{\text{s:o}})$$

Reasonable if the angular resolution is much poorer in $K'$ than in $K$.

**One-to-several:** symmetric of several-to-one ($K$ and $K'$ swapped). $\qquad (H_{\text{o:s}})$
Appropriate for extended sources looking single at the wavelength of $K$ but breaking up at that of $K'$.

**One-to-one:** any $K$-source has at most one counterpart in $K'$ and reciprocally, i.e.

$$\text{all the events } (A_{i,j})_{i\in[\![1,n]\!],j\in[\![1,n']\!]} \text{ are exclusive.} \qquad (H_{\text{o:o}})$$

One has then $f\,n = f'\,n'$. Natural assumption for point sources if the angular resolution is high in both $K$ and $K'$.

**Several-to-several:** not considered.

In $H_{\text{s:o}}$ and $H_{\text{o:s}}$, catalogs do not play symmetrical roles. Assumption $H_{\text{o:o}}$ is therefore more neutral, but calculations with $H_{\text{s:o}}$ are much simpler and will serve as a guide; they also provide initial values for one-to-one computations.

# Three related problems

▸ For each assumption, **calculate the probability of association** $P(A_{i,j} \mid C \cap C')$ that $M_i$ is the same as $M'_j$ ($j > 0$) or that $M_i$ has no counterpart ($j = 0$), given the coordinates of all sources and the unknown parameters.

▸ **Estimate unknown parameters** from the data, in particular the *a priori* probability $f$ that any $K$-source has some counterpart.

▸ **Select the most likely assumption**, i.e. the most appropriate association model, given the data.

# Computation of $P_{s:o}(A_{i,j} \mid C \cap C')$

$$P(A_{i,j} \mid C \cap C') = \frac{P(A_{i,j} \cap C \mid C')}{P(C \mid C')}.$$

## Computation of the denominator

$M_i$ may be associated to $M'_j$ which may be associated to $M_k$ which may be associated to $M'_\ell$ which may be associated to $M_i \ldots \Rightarrow$ One needs to consider all possible combinations of all the events $A_{k,j_k}$ and order them.

Event

$$\bigcap_{k=1}^{n} \bigcup_{j_k=0}^{n'} A_{k,j_k} = \bigcup_{j_1=0}^{n'} \bigcup_{j_2=0}^{n'} \cdots \bigcup_{j_n=0}^{n'} \bigcap_{k=1}^{n} A_{k,j_k}$$

is certain, so

$$P_{s:o}(C \mid C') = P_{s:o}\left(C \cap \bigcap_{k=1}^{n} \bigcup_{j_k=0}^{n'} A_{k,j_k} \mid C'\right) = \sum_{j_1=0}^{n'} \sum_{j_2=0}^{n'} \cdots \sum_{j_n=0}^{n'} P_{s:o}\left(C \cap \bigcap_{k=1}^{n} A_{k,j_k} \mid C'\right)$$

$$= \sum_{j_1=0}^{n'} \sum_{j_2=0}^{n'} \cdots \sum_{j_n=0}^{n'} P_{s:o}\left(C \mid \bigcap_{k=1}^{n} A_{k,j_k} \cap C'\right) P_{s:o}\left(\bigcap_{k=1}^{n} A_{k,j_k} \mid C'\right).$$

One has

$$P_{\text{s:o}}\left(C \mid \bigcap_{k=1}^{n} A_{k,j_k} \cap C'\right) = \text{cst.} \times \prod_{k=1}^{n} \xi_{k,j_k}$$

and

$$P_{\text{s:o}}\left(\bigcap_{k=1}^{n} A_{k,j_k} \mid C'\right) = P_{\text{s:o}}\left(\bigcap_{k=1}^{n} A_{k,j_k}\right) = \left(\frac{f}{n'}\right)^{q} (1-f)^{n-q},$$

where $q$ is the number of $K$-sources with a counterpart.

Finally,

$$P_{\text{s:o}}(C \mid C') = \text{cst.} \times \sum_{j_1=0}^{n'} \sum_{j_2=0}^{n'} \cdots \sum_{j_n=0}^{n'} \prod_{k=1}^{n} \zeta_{k,j_k} \,,$$

where

$$\zeta_{k,0} := (1-f)\,\xi_{k,0} \quad \text{and} \quad \zeta_{k,j_k} := \frac{f\,\xi_{k,j_k}}{n'} \quad \text{if } j_k > 0.$$

## Computation of the numerator

Similarly,

$$P_{\text{s:o}}(A_{i,j} \cap C \mid C') = \text{cst.} \times \zeta_{i,j} \sum_{j_1=0}^{n'} \cdots \sum_{j_{i-1}=0}^{n'} \sum_{j_{i+1}=0}^{n'} \cdots \sum_{j_n=0}^{n'} \prod_{\substack{k=1 \\ k \neq i}}^{n} \zeta_{k,j_k}.$$

## Ratio

$P_{\text{s:o}}(C \mid C')$ and $P_{\text{s:o}}(A_{i,j} \cap C \mid C')$ may be factorized:

$$\sum_{j_1=0}^{n'} \sum_{j_2=0}^{n'} \cdots \sum_{j_n=0}^{n'} \prod_{k=1}^{n} \zeta_{k,j_k} = \prod_{k=1}^{n} \sum_{j_k=0}^{n'} \zeta_{k,j_k},$$

so

$$
P_{\text{s:o}}(A_{i,j} \mid C \cap C') = \frac{\zeta_{i,j} \prod_{\substack{k=1 \\ k \neq i}}^{n} \sum_{j_k=0}^{n'} \zeta_{k,j_k}}{\prod_{k=1}^{n} \sum_{j_k=0}^{n'} \zeta_{k,j_k}} = \frac{\zeta_{i,j}}{\sum_{k=0}^{n'} \zeta_{i,k}}
$$

$$
= \begin{cases} \dfrac{f\,\xi_{i,j}}{(1-f)\,n'\,\xi_{i,0} + f\sum_{k=1}^{n'}\xi_{i,k}} & \text{if } j > 0, \\[2em] \dfrac{(1-f)\,n'\,\xi_{i,0}}{(1-f)\,n'\,\xi_{i,0} + f\sum_{k=1}^{n'}\xi_{i,k}} & \text{if } j = 0. \end{cases}
$$

In practice, the sums on $k$ may be restricted to sources $M'_k$ close to $M_i$.

# Likelihood and estimation of unknown parameters under $H_{\text{s:o}}$

Maximize the likelihood

$$L = \text{cst.} \times P(C \cap C')$$

to observe all sources at their effective positions. Maximum likelihood estimates $\hat{x}, \hat{y}$, etc., of the unknown parameters $x, y$, etc., are thus obtained by solving

$$\left( \frac{\partial L}{\partial x} \right)_{(x, y, \ldots) = (\hat{x}, \hat{y}, \ldots)} = 0.$$

Under the several-to-one assumption, one has

$$L_{\text{s:o}} = \text{cst.} \times \prod_{i=1}^{n} \sum_{k=0}^{n'} \zeta_{i, k},$$

from which one derives that

$$\frac{\partial \ln L_{\text{s:o}}}{\partial f} = \frac{n\,(1 - f) - \sum_{i=1}^{n} P_{\text{s:o}}(A_{i, 0} \mid C \cap C')}{f\,(1 - f)},$$

so

$$\hat{f}_{\text{s:o}} = 1 - \frac{1}{n} \sum_{i=1}^{n} \hat{P}_{\text{s:o}}(A_{i, 0} \mid C \cap C'), \quad \text{where } \hat{P}_{\text{s:o}} = (P_{\text{s:o}})_{f = \hat{f}_{\text{s:o}}}.$$

As $\hat{f}_{\text{s:o}}$ appears on both sides, we calculate it by a back and forth iteration between the l.h.s. and the r.h.s., starting from some arbitrary $f \in [0, 1]$. (Converges very fast.)

# Theoretical computation of $P_{o:o}(A_{i,j} \mid C \cap C')$

A $K'$-source associated to a $K$-source may not be associated to another one, so

$$P_{o:o}(C \mid C') = \text{cst.} \times \sum_{\substack{j_1=0 \\ j_1 \notin X_0}}^{n'} \sum_{\substack{j_2=0 \\ j_2 \notin X_1}}^{n'} \cdots \sum_{\substack{j_n=0 \\ j_n \notin X_{n-1}}}^{n'} \prod_{k=1}^{n} \eta_{k,j_k} \,,$$

where $X_{k-1}$ is the set of excluded $K'$-sources at depth $k$ in the recursive sum (i.e.

$$X_0 := \varnothing, \quad X_k := (X_{k-1} \cup \{j_k\}) \setminus \{0\},$$

so $X_k$ contains the counterparts already associated with $M_1, \ldots, M_{k-1}$, which may therefore not be associated with $M_k, \ldots, M_n$) and

$$\eta_{k,0} := (1-f)\,\xi_{k,0} \quad \text{and} \quad \eta_{k,j_k} := \frac{f\,\xi_{k,j_k}}{n' - \#X_{k-1}} \quad \text{if } j_k > 0.$$

$P_{o:o}(A_{i,j} \mid C \cap C')$ is computed similarly and

$$P_{o:o}(A_{i,j} \mid C \cap C') = \frac{\zeta_{i,j} \sum_{\substack{j_1=0 \\ j_1 \notin X_0^*}}^{n'} \cdots \sum_{\substack{j_{i-1}=0 \\ j_{i-1} \notin X_{i-2}^*}}^{n'} \sum_{\substack{j_{i+1}=0 \\ j_{i+1} \notin X_i^*}}^{n'} \cdots \sum_{\substack{j_n=0 \\ j_n \notin X_{n-1}^*}}^{n'} \prod_{\substack{k=1 \\ k \neq i}}^{n} \eta_{k,j_k}^*}{\sum_{\substack{j_1=0 \\ j_1 \notin X_0}}^{n'} \sum_{\substack{j_2=0 \\ j_2 \notin X_1}}^{n'} \cdots \sum_{\substack{j_n=0 \\ j_n \notin X_{n-1}}}^{n'} \prod_{k=1}^{n} \eta_{k,j_k}}$$

("$*$" means that $j$ is also excluded if $j > 0$).

# Likelihood and estimation of unknown parameters under $H_{o:o}$

The recursive sums may not be factorized, so the ratio may not be simplified, contrary to the several-to-one case. Because of the combinatorial explosion of the number of terms, $P_{o:o}(A_{i,j} \mid C \cap C')$ and $L_{o:o}$ seem impossible to evaluate.

Assume nonetheless that one can compute $P_{o:o}(A_{i,j} \mid C \cap C')$. Then, one can show that one still has

$$\frac{\partial \ln L_{o:o}}{\partial f} = \frac{n\,(1-f) - \sum_{i=1}^{n} P_{o:o}(A_{i,0} \mid C \cap C')}{f\,(1-f)} \, ,$$

which gives $\hat{f}_{o:o}$ by the back and forth iteration described earlier.

Since $L_{o:o}(f = 0)$ (i.e., when all sources are randomly distributed) is known, $L_{o:o}$ may also be obtained for any $f$ by integrating $\partial \ln L_{o:o}/\partial f$.
One can also compute $L_{o:o}$ like this:

$$L_{o:o} = \text{cst.} \times \prod_{i=1}^{n} \frac{(1-f)\,\xi_{i,0}}{P_{o:o}(A_{i,0} \mid C \cap C' \cap \bigcap_{k=1}^{i-1} A_{k,0})} \, .$$

# Practical computation of $P_{\text{o:o}}(A_{i,j} \mid C \cap C')$

A partially true idea:

**$P_{\text{o:o}}(A_{i,j} \mid C \cap C')$ depends only on the neighbors of $M_i$ and $M'_j$.**

One may indeed expect that, although the numerator $P_{\text{o:o}}(A_{i,j} \cap C \mid C')$ and the denominator $P_{\text{o:o}}(C \mid C')$ depend on distant sources, the effect of these cancels in the *ratio* numerator/denominator.

So, order $K$-sources *by increasing distance* to $M_i$ and rewrite the ratio of the recursive sums in the numerator and denominator up to some depth $\ell$:

$$p_\ell := \frac{\zeta_{i,j} \sum_{\substack{j_2=0 \\ j_2 \notin \widetilde{X}_1^*}}^{n'} \cdots \sum_{\substack{j_\ell=0 \\ j_\ell \notin \widetilde{X}_{\ell-1}^*}}^{n'} \prod_{k=2}^{\ell} \widetilde{\eta}_{k,j_k}^*}{\sum_{\substack{j_1=0 \\ j_1 \notin \widetilde{X}_0}}^{n'} \sum_{\substack{j_2=0 \\ j_2 \notin \widetilde{X}_1}}^{n'} \cdots \sum_{\substack{j_\ell=0 \\ j_\ell \notin \widetilde{X}_{\ell-1}}}^{n'} \prod_{k=1}^{\ell} \widetilde{\eta}_{k,j_k}} \qquad \text{(the tilde is for the reordering)}.$$

One has $p_n = P_{\text{o:o}}(A_{i,j} \mid C \cap C')$. When $\ell \nearrow$, more distant neighbors are progressively included in $p_\ell$. The ratio $p_\ell$ oscillates for small $\ell$, then stabilizes for some value $\ell_{\text{stable}}$. It is therefore tempting to conclude that the sequence $(p_\ell)$ has converged and to set

$$P_{\text{o:o}}^{\text{w}}(A_{i,j} \mid C \cap C') = p_{\ell_{\text{stable}}}$$

(all the more tempting that, when $M_i$ is the only $K$-source considered, $p_1 = P_{\text{s:o}}(A_{i,j} \mid C \cap C')$).

As we will see, this is however *wrong*, as emphasized by the superscript "w".

# One-to-one all-sky simulations: computations with $P^{\mathrm{w}}_{\mathrm{o:o}}$

Known circular positional uncertainties: combined uncertainty $\mathring{\sigma} = \sqrt{\sigma^2 + \sigma'^2}$.
$\mathring{\sigma}$ known $\Rightarrow$ only $f$ must be estimated.
Input values: $f = 1/2$; $\mathring{\sigma} = 10^{-3}$ rad; $n' = 10^5$; $n \in [\![10^3, 10^5]\!]$.

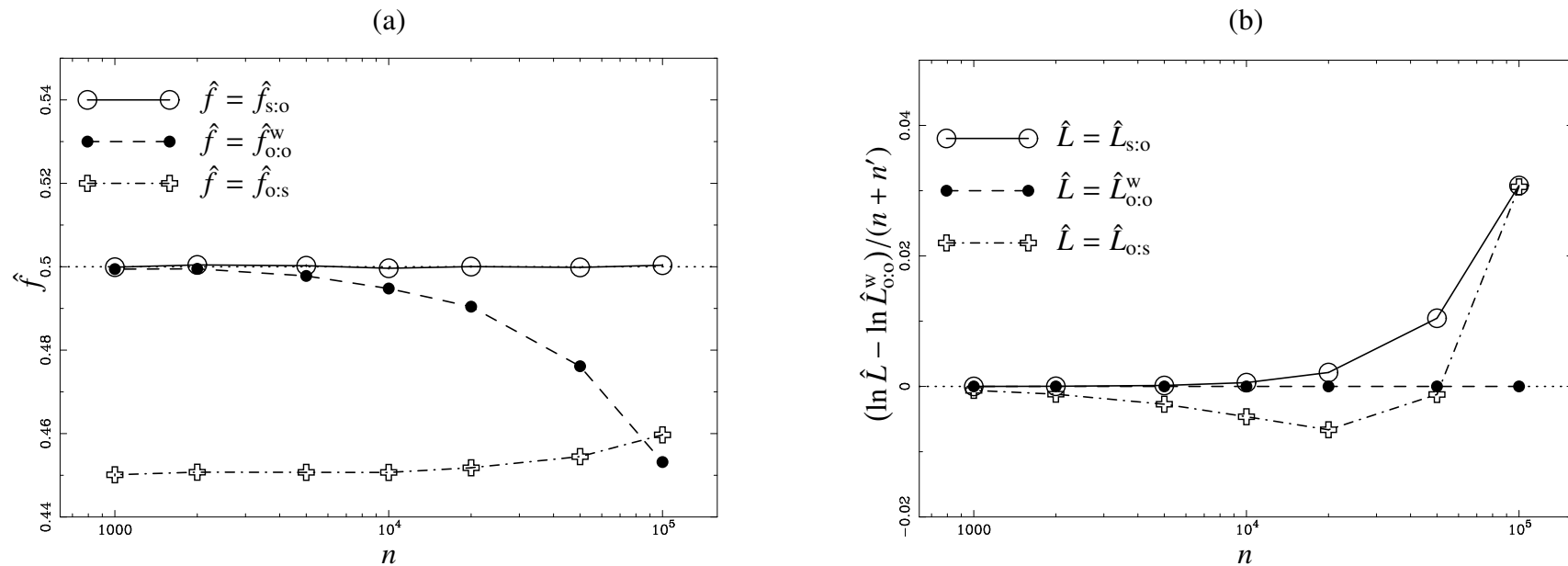(a)                                       (b)



Figure 1.     **(a)** Mean value of different estimators $\hat{f}$ of $f$ as a function of $n$. The dotted line indicates the input value of $f$. (Note: $\hat{f}_{\mathrm{o:s}}$ is derived from the ML estimator $\hat{f}'_{\mathrm{o:s}}$, but is not an ML estimator itself.)
     **(b)** Normalized average maximum value $\hat{L}$ of different likelihoods as a function of $n$, compared to $\hat{L}^{\mathrm{w}}_{\mathrm{o:o}}$.

# Analysis of one-to-one simulations

For one-to-one simulations, one obtains the following:

- Several-to-one estimators always provide closer values to $f$ and $\mathring{\sigma}$ than one-to-one "w" estimators!

- One-to-one "w" estimators are *statistically inconsistent*: the bias does not tend to 0 when $n \nearrow$.
  (Note however that maximum likelihood estimators may be inconsistent in some circumstances. Conditions used to prove their consistency are not applicable here.)

- The maximum of the several-to-one likelihood is larger than that of the "w" one-to-one likelihood.

**Something wrong in the computation of one-to-one probabilities?**

# Reconsideration

After scrutiny of a simple example ($n = n' = 2$), it became clear that *distant sources do matter*, but *only by their number, not their exact positions*: they lock some number of counterparts which may not be associated to $M_i$ and its neighbors.

In the sequence $(p_\ell)$ that should converge to $P_{o:o}(A_{i,j} \mid C \cap C')$, the number $n'$ must be repaced by $n'_{\text{eff}}$, the number of $K'$-sources that may *effectively* be associated with $M_i$ and its $\ell - 1$ nearest neighbors. One has

$$n'_{\text{eff}} = n' - \sum_{\text{distant } M_k} (1 - P_{o:o}[A_{k,0} \mid C \cap C']).$$

Note that for $\ell = n$, one has $n'_{\text{eff}} = n'$ and one recovers the theoretical result for $P_{o:o}(A_{i,j} \mid C \cap C')$.

As $P_{o:o}$ depends on $n'_{\text{eff}}$ which itself depends on $P_{o:o}$, both may be computed with a back and forth iteration, taking $P_{s:o}$ as the initial value of $P_{o:o}$.

(What happened with the partially true idea that only neighbors matter is that, after a transient phase where $p_\ell$ oscillated, a steady state was reached. The ratio $p_\ell$ had however not converged, but was slowly drifting to $p_n = P_{o:o}(A_{i,j} \mid C \cap C')$.)

# Simulations with *known circular* positional uncertainties (revised)

$\mathring{\sigma}$ known. Only $f$ must be estimated.

Input values: $f = 1/2$; $\mathring{\sigma} = 10^{-3}$ rad; $n' = 10^5$; $n \in [\![10^3, 10^5]\!]$.
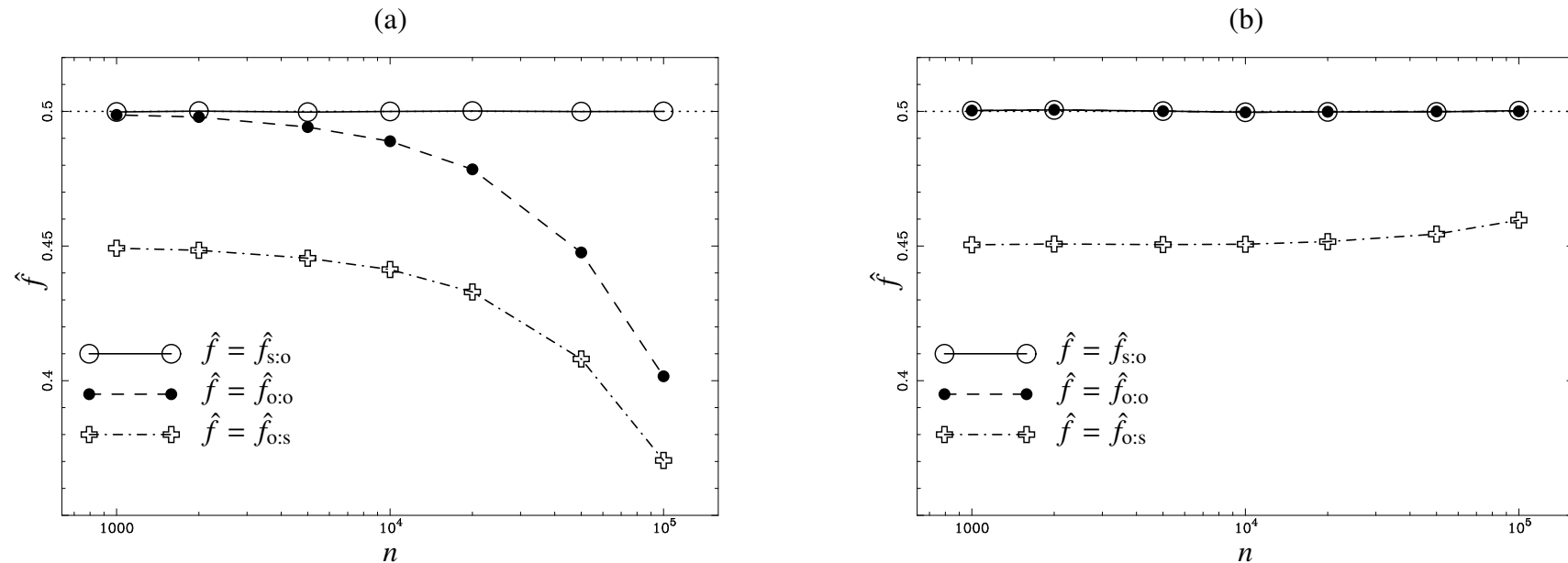


Figure 2.    Mean value of different estimators $\hat{f}$ of $f$ as a function of $n$.

**(a)** Several-to-one simulations.

**(b)** One-to-one simulations.

# Simulations with *unknown circular* positional uncertainties

Both $f$ and $\mathring{\sigma}$ must be estimated.
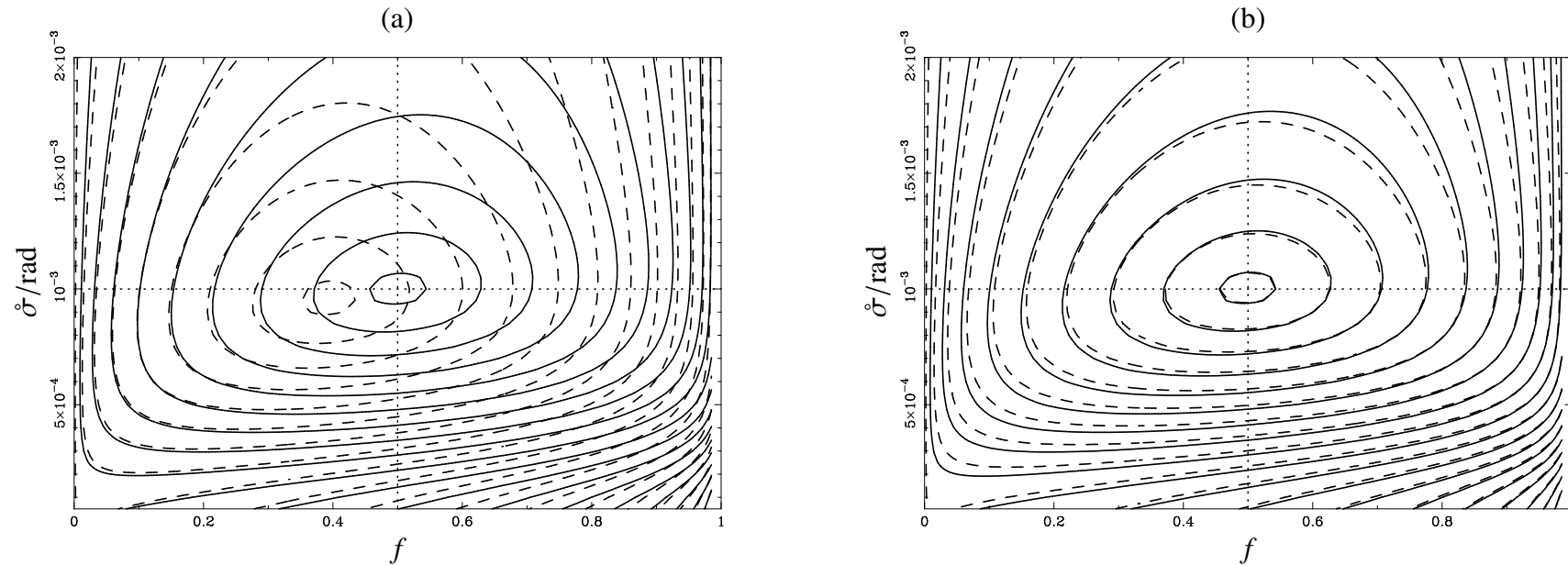Input values: $f = 1/2$; $\mathring{\sigma} = 10^{-3}$ rad; $n = n' = 2 \times 10^4$.



**Figure 3.**    Contour lines of $L_{\text{s:o}}$ (solid) and $L_{\text{o:o}}$ (dashed). The input values of $f$ and $\mathring{\sigma}$ are indicated by dotted lines.

**(a)** Several-to-one simulations.

**(b)** One-to-one simulations.

# Simulations with *unknown elliptical* positional uncertainties

Both $f$ and $\mathring{\sigma}$ must be estimated.

Input values: $f = 1/2$; randomly oriented positional uncertainty ellipses with a semi-major axis of $1.5 \times 10^{-3}$ rad and a semi-minor axis of $0.5 \times 10^{-3}$ rad; $n = n' = 2 \times 10^4$.
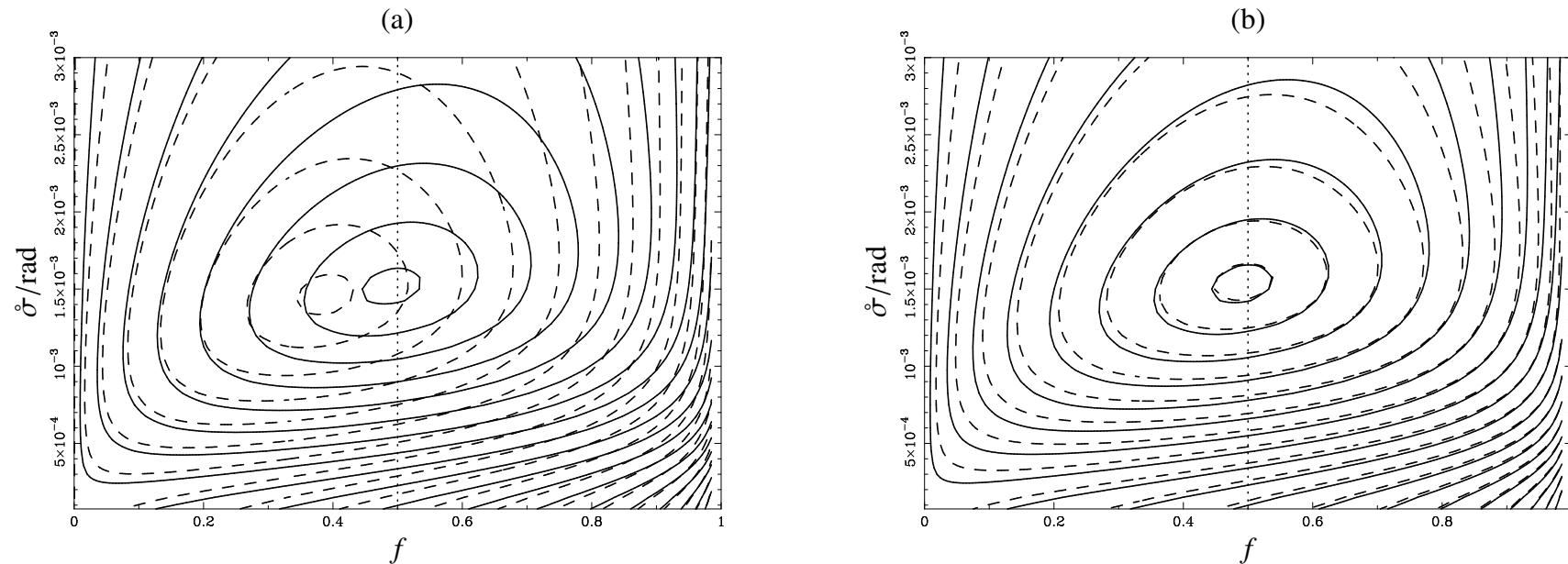
Figure 4.    Contour lines of $L_{s:o}$ (solid) and $L_{o:o}$ (dashed). The input value of $f$ is indicated by a dotted line.

(a) Several-to-one simulations.

(b) One-to-one simulations.

# Comparison of maximum likelihoods (revised)

Circular positional uncertainties. Combined positional uncertainty known.
Input values: $f = 1/2$; $\mathring{o} = 10^{-3}$ rad; $n' = 10^5$; $n \in [\![10^3, 10^5]\!]$.
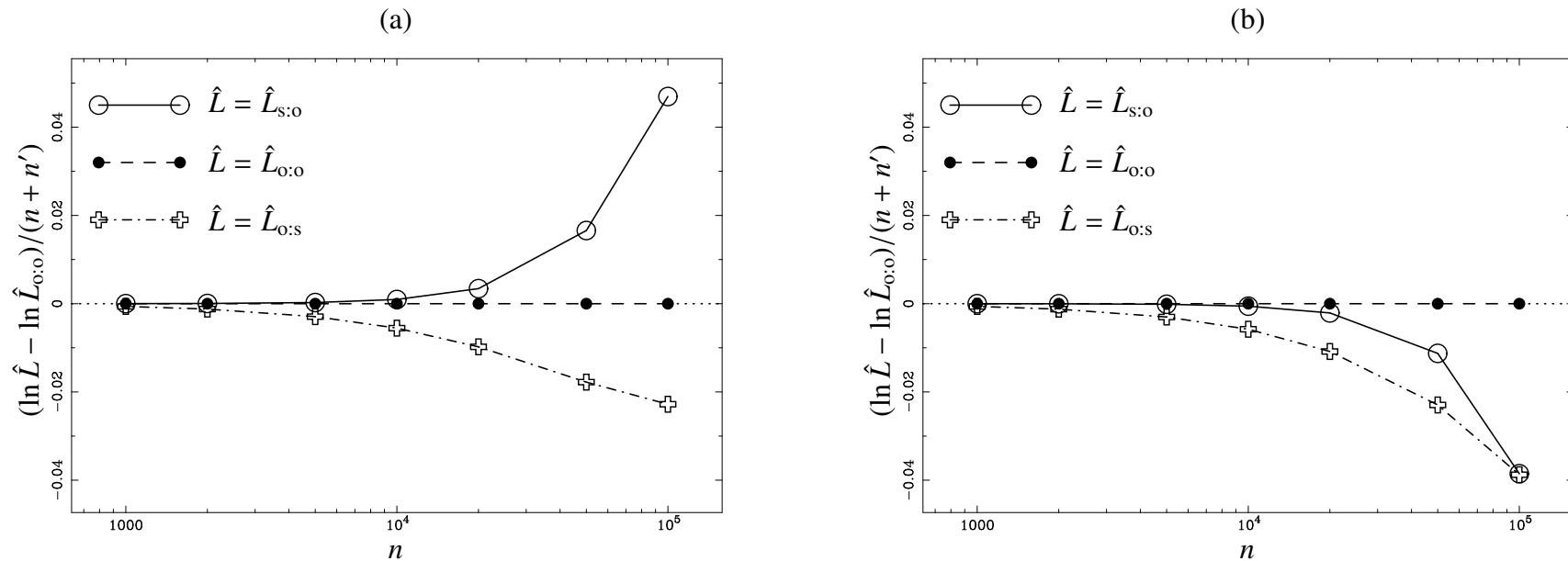
(a)  (b)



**Figure 5.**  Normalized mean value of different likelihoods at their maximum as a function of $n$.

**(a)** Several-to-one simulations.

**(b)** One-to-one simulations.

# Conclusions

- Several-to-one estimators provide unbiased values of $f$ for several-to-one simulations, but also for one-to-one simulations.

- Revised one-to-one estimators provide unbiased values of $f$ for one-to-one simulations, but not for several-to-one simulations (not a problem).

- The same holds for $\mathring{\sigma}$ if it is unknown.

- These estimators are robust: if $\mathring{\sigma}$ is unknown or if positional uncertainties are elliptical, the right value of $f$ is still recovered.

- For several-to-one simulations, $\hat{L}_{s:o} > \hat{L}_{o:o} > \hat{L}_{o:s}$, as expected.

- For one-to-one simulations, $\hat{L}_{o:o} > \hat{L}_{s:o}$ and $\hat{L}_{o:o} > \hat{L}_{o:s}$, as expected.

# The Aspects code

All these simulations were created and analyzed with the Fortran 95 code Aspects ([aspɛ]). Source available at

**www2.iap.fr/users/fioc/Aspects/** .

- ▸ Probabilities of cross-identification.
- ▸ Fraction of sources without counterpart.
- ▸ Likelihood of the association model.
- ▸ Estimation of the positional uncertainty. For extended sources,
    - ▸ possibly different true positions,
    - ▸ size-dependent unknown positional uncertainties.

**References:**

- ▸ Paper: $A\&A$, 566, A8 (`arXiv:1209.5361`);
- ▸ Code documentation and complements: `arXiv:1404.4224`.