

## The ARCHES cross-correlation tool

François-Xavier Pineau<sup>1</sup>

<sup>1</sup>Observatoire Astronomique de Strasbourg, Université de Strasbourg, CNRS

Paris, 1<sup>th</sup> December, 2015



- This talk: “*Cross-correlation tool development & catalogue creation*” (WP4)
- Aims of ARCHES’s WP4:
  - ▶ Create a public  $n$ -catalogues cross-correlation tool:
    - ★ No magic BUT a flexible/multi-purpose/scriptable multi-catalogue xmatch engine
    - ★ Usable as a building block from you own specific code
  - ▶ Use/develop statistical methods to compute probabilities of associations:
    - ★ Astrometry based probabilities only!
    - ★ Can be combined with photometry based probabilities (in a further step)
  - ▶ Use the tool to build ARCHES catalogue(s)
- Beyond the ARCHES project:
  - ▶ tool will be part of the CDS XMatch Service
  - ▶ ⇒ will be maintained, will keep evolving

- This talk: “*Cross-correlation tool development & catalogue creation*” (WP4)
- Aims of ARCHES’s WP4:
  - ▶ Create a public  $n$ -catalogues cross-correlation tool:
    - ★ No magic BUT a flexible/multi-purpose/scriptable multi-catalogue xmatch engine
    - ★ Usable as a building block from you own specific code
  - ▶ Use/develop statistical methods to compute probabilities of associations:
    - ★ Astrometry based probabilities only!
    - ★ Can be combined with photometry based probabilities (in a further step)
  - ▶ Use the tool to build ARCHES catalogue(s)
- Beyond the ARCHES project:
  - ▶ tool will be part of the CDS XMatch Service
  - ▶ ⇒ will be maintained, will keep evolving

- This talk is mainly focused on the probabilistic part
- More details on the tool during the *Hands on* session

- Steps to probabilistic positional xmatch
  - ▶ Make simplifying assumptions
  - ▶ Select candidates: select and group together sources possibly being various detections of a same real source
    - ★ Need for a selection criterion
  - ▶ Make hypothesis: are the sources really from a same real sources or from different real sources?
  - ▶ For each hypothesis:
    - ★ derive the associated *likelihood*
    - ★ derive the associated *prior*
  - ▶ Compute astrometry based probabilities

- Radical simplifying assumptions:
  - ▶ No proper motions
  - ▶ No blending
  - ▶ No clustering (density of sources = Poisson law)
  - ▶ No systematic offsets
  - ▶ You can trust positional uncertainties provided in catalogues

How to select a group of  $n$  sources from  $n$  distinct catalogues as possibly being various observations of a same actual source?

- *Statistical hypothesis testing*

- ▶  $H_0$  (null hypothesis): all  $n$  sources are from the same real source
- ▶  $H_1 = \bar{H}_0$  (alternative hypothesis): at least one source (out of  $n$ ) is spurious

- User input:  $\gamma$ , the probability to accept  $H_0$  while it is true

- ▶  $\gamma$  (I call it completeness) is called *true negative rate*
- ▶ we usually fix  $\gamma = 0.9973$  (99.73%, value of the  $3\sigma$  rule in 1 dimensional pb)
- ▶  $\Leftrightarrow$  fixing the *type I error* = 0.027% = proba to reject null hypothesis while it is true
- ▶ we (theoretically) miss 27/10 000 real association

- The criterion used is based on a  $\chi^2$  test of  $2(n - 1)$  degrees of freedom

- Now, a few slides to explain it since it plays a role in probabilities

How to select a group of  $n$  sources from  $n$  distinct catalogues as possibly being various observations of a same actual source?

- *Statistical hypothesis testing*
  - ▶  $H_0$  (null hypothesis): all  $n$  sources are from the same real source
  - ▶  $H_1 = \bar{H}_0$  (alternative hypothesis): at least one source (out of  $n$ ) is spurious
- User input:  $\gamma$ , the probability to accept  $H_0$  while it is true
  - ▶  $\gamma$  (I call it completeness) is called *true negative rate*
  - ▶ we usually fix  $\gamma = 0.9973$  (99.73%, value of the  $3\sigma$  rule in 1 dimensional pb)
  - ▶  $\Leftrightarrow$  fixing the *type I error* = 0.027% = proba to reject null hypothesis while it is true
  - ▶ we (theoretically) miss 27/10 000 real association
- The criterion used is based on a  $\chi^2$  test of  $2(n - 1)$  degrees of freedom
- Now, a few slides to explain it since it plays a role in probabilities



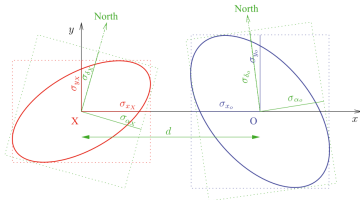
How to select a group of  $n$  sources from  $n$  distinct catalogues as possibly being various observations of a same actual source?

- *Statistical hypothesis testing*
  - ▶  $H_0$  (null hypothesis): all  $n$  sources are from the same real source
  - ▶  $H_1 = \bar{H}_0$  (alternative hypothesis): at least one source (out of  $n$ ) is spurious
- User input:  $\gamma$ , the probability to accept  $H_0$  while it is true
  - ▶  $\gamma$  (I call it completeness) is called *true negative rate*
  - ▶ we usually fix  $\gamma = 0.9973$  (99.73%, value of the  $3\sigma$  rule in 1 dimensional pb)
  - ▶  $\Leftrightarrow$  fixing the *type I error* = 0.027% = proba to reject null hypothesis while it is true
  - ▶ we (theoretically) miss 27/10 000 real association
- The criterion used is based on a  $\chi^2$  test of  $2(n - 1)$  degrees of freedom
- Now, a few slides to explain it since it plays a role in probabilities

- In the classical case (e.g. De Ruiter et al. 1977):
  - ▶ Errors are independent on  $\alpha$  and  $\delta$
  - ▶ Source 1 has errors  $\sigma_{\alpha_1}$  and  $\sigma_{\delta_1}$  on  $\alpha$  and  $\delta$  respectively
  - ▶ Source 2 has errors  $\sigma_{\alpha_2}$  and  $\sigma_{\delta_2}$  on  $\alpha$  and  $\delta$  respectively
  - ▶ The normalized distance (or  $\sigma$ -distance) is defined by:

$$r = \left[ \frac{\Delta\alpha^2}{\sigma_{\alpha_1}^2 + \sigma_{\alpha_2}^2} + \frac{\Delta\delta^2}{\sigma_{\delta_1}^2 + \sigma_{\delta_2}^2} \right]^{1/2}$$

- More generally (see e.g. Pineau et al. 2011)
  - ▶ We assimilate locally the surface of the sphere to the Euclidian plane
  - ▶ The positions of the 2 sources are 2 dimensional vectors:  $\vec{\mu}_1$  and  $\vec{\mu}_2$ .
  - ▶ Errors on  $\vec{\mu}_1$  and  $\vec{\mu}_2$  are oriented ellipses defined by covariance matrices  $\mathbf{V}_1$  and  $\mathbf{V}_2$  respectively:



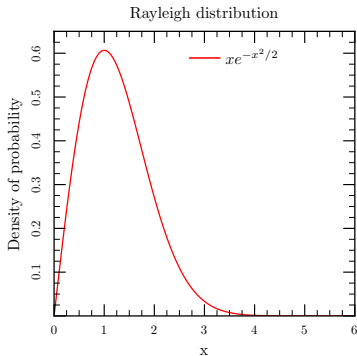
- ▶ The normalized distance becomes (vectorial form):

$$r = \left[ (\vec{\mu}_1 - \vec{\mu}_2)^T (\mathbf{V}_1 + \mathbf{V}_2)^{-1} (\vec{\mu}_1 - \vec{\mu}_2) \right]^{1/2}$$

- ▶  $\Rightarrow$  equation of an ellipse of radius  $r$  and covariance matrix  $\mathbf{V}_1 + \mathbf{V}_2$

- For real associations, i.e. when  $H_0$  is true
  - ▶ The distribution of normalized distances is a Rayleigh distribution of scale  $\sigma = 1$

$$r \stackrel{H_0}{\sim} \text{Rayleigh}$$



- Fixing the completeness  $\gamma \Leftrightarrow$  fixing a normalized distance threshold  $k_\gamma$ :

$$\int_0^{k_\gamma} \text{Rayleigh}(r) dr = \gamma$$

- For  $\gamma = 99.73\%$  (the 1D  $3\sigma$  rule)  $\Rightarrow k_\gamma = 3.4395$  (not 3!)
- So, for 2 sources from 2 distinct catalogues, the selection criterion is

$$[(\vec{\mu}_1 - \vec{\mu}_2)^T (\mathbf{V}_1 + \mathbf{V}_2)^{-1} (\vec{\mu}_1 - \vec{\mu}_2)]^{1/2} \leq k_\gamma$$

- I.e. source 2 kept as candidate if it is inside an error ellipse of covariance matrix  $\mathbf{V} = \mathbf{V}_1 + \mathbf{V}_2$  and of radius  $k_\gamma$ , centered around source 1.
- $\Rightarrow$  the surface area of the acceptance region is  $|\mathbf{V}_1 + \mathbf{V}_2|^{1/2} \pi k_\gamma^2$

- Fixing the completeness  $\gamma \Leftrightarrow$  fixing a normalized distance threshold  $k_\gamma$ :

$$\int_0^{k_\gamma} \text{Rayleigh}(r) dr = \gamma$$

- For  $\gamma = 99.73\%$  (the 1D  $3\sigma$  rule)  $\Rightarrow k_\gamma = 3.4395$  (not 3!)
- So, for 2 sources from 2 distinct catalogues, the selection criterion is

$$[(\vec{\mu}_1 - \vec{\mu}_2)^T (\mathbf{V}_1 + \mathbf{V}_2)^{-1} (\vec{\mu}_1 - \vec{\mu}_2)]^{1/2} \leq k_\gamma$$

- I.e. source 2 kept as candidate if it is inside an error ellipse of covariance matrix  $\mathbf{V} = \mathbf{V}_1 + \mathbf{V}_2$  and of radius  $k_\gamma$ , centered around source 1.
- $\Rightarrow$  the surface area of the acceptance region is  $|\mathbf{V}_1 + \mathbf{V}_2|^{1/2} \pi k_\gamma^2$

- Fixing the completeness  $\gamma \Leftrightarrow$  fixing a normalized distance threshold  $k_\gamma$ :

$$\int_0^{k_\gamma} \text{Rayleigh}(r) dr = \gamma$$

- For  $\gamma = 99.73\%$  (the 1D  $3\sigma$  rule)  $\Rightarrow k_\gamma = 3.4395$  (not 3!)
- So, for 2 sources from 2 distinct catalogues, the selection criterion is

$$[(\vec{\mu}_1 - \vec{\mu}_2)^T (\mathbf{V}_1 + \mathbf{V}_2)^{-1} (\vec{\mu}_1 - \vec{\mu}_2)]^{1/2} \leq k_\gamma$$

- I.e. source 2 kept as candidate if it is inside an error ellipse of covariance matrix  $\mathbf{V} = \mathbf{V}_1 + \mathbf{V}_2$  and of radius  $k_\gamma$ , centered around source 1.
- $\Rightarrow$  the surface area of the acceptance region is  $|\mathbf{V}_1 + \mathbf{V}_2|^{1/2} \pi k_\gamma^2$

Now, a different version of the same story more easily generalisable to  $n$ -catalogues.



- I have 2 sources from 2 distinct catalogues, I suppose  $H_0$  is *true*
- Maximum Likelihood Estimate (MLE) of the position of the real source  $\rightsquigarrow$  the weighted mean position

$$\vec{\mu}_{\Sigma} = \mathbf{V}_{\Sigma}(\mathbf{V}_1^{-1}\vec{\mu}_1 + \mathbf{V}_2^{-1}\vec{\mu}_2)$$

in which

$$\mathbf{V}_{\Sigma} = (\mathbf{V}_1^{-1} + \mathbf{V}_2^{-1})^{-1}$$

- The error on this MLE is ...  $\mathbf{V}_{\Sigma}$
- The result is the same with a (by block) Weighted Least Squares method
- We can now define the Mahalanobis distance:

$$D_M = \left[ \sum_{i=1}^2 (\vec{\mu}_i - \vec{\mu}_{\Sigma})^T \mathbf{V}_i^{-1} (\vec{\mu}_i - \vec{\mu}_{\Sigma}) \right]^{1/2} \underset{H_0}{\sim} \chi_{dof=2}$$

- I have 2 sources from 2 distinct catalogues, I suppose  $H_0$  is *true*
- Maximum Likelihood Estimate (MLE) of the position of the real source  $\rightsquigarrow$  the weighted mean position

$$\vec{\mu}_{\Sigma} = \mathbf{V}_{\Sigma}(\mathbf{V}_1^{-1}\vec{\mu}_1 + \mathbf{V}_2^{-1}\vec{\mu}_2)$$

in which

$$\mathbf{V}_{\Sigma} = (\mathbf{V}_1^{-1} + \mathbf{V}_2^{-1})^{-1}$$

- The error on this MLE is ...  $\mathbf{V}_{\Sigma}$
- The result is the same with a (by block) Weighted Least Squares method
- We can now define the Mahalanobis distance:

$$D_M = \left[ \sum_{i=1}^2 (\vec{\mu}_i - \vec{\mu}_{\Sigma})^T \mathbf{V}_i^{-1} (\vec{\mu}_i - \vec{\mu}_{\Sigma}) \right]^{1/2} \underset{H_0}{\sim} \chi_{dof=2}$$

- I have 2 sources from 2 distinct catalogues, I suppose  $H_0$  is *true*
- Maximum Likelihood Estimate (MLE) of the position of the real source  $\rightsquigarrow$  the weighted mean position

$$\vec{\mu}_{\Sigma} = \mathbf{V}_{\Sigma}(\mathbf{V}_1^{-1}\vec{\mu}_1 + \mathbf{V}_2^{-1}\vec{\mu}_2)$$

in which

$$\mathbf{V}_{\Sigma} = (\mathbf{V}_1^{-1} + \mathbf{V}_2^{-1})^{-1}$$

- The error on this MLE is ...  $\mathbf{V}_{\Sigma}$
- The result is the same with a (by block) Weighted Least Squares method
- We can now define the Mahalanobis distance:

$$D_M = \left[ \sum_{i=1}^2 (\vec{\mu}_i - \vec{\mu}_{\Sigma})^T \mathbf{V}_i^{-1} (\vec{\mu}_i - \vec{\mu}_{\Sigma}) \right]^{1/2} \underset{H_0}{\sim} \chi_{dof=2}$$

Let's merge the two approaches.

Doing the math, we find

- the equality

$$|\mathbf{V}_1 + \mathbf{V}_2| = \frac{|\mathbf{V}_1||\mathbf{V}_2|}{|\mathbf{V}_\Sigma|}$$

- the normalized (or  $\sigma$ ) distance = Mahalanobis distance

$$(\vec{\mu}_1 - \vec{\mu}_2)^T (\mathbf{V}_1 + \mathbf{V}_2)^{-1} (\vec{\mu}_1 - \vec{\mu}_2) = \sum_{i=1}^2 (\vec{\mu}_i - \vec{\mu}_\Sigma)^T \mathbf{V}_i^{-1} (\vec{\mu}_i - \vec{\mu}_\Sigma)$$

- Remark: not surprising when we know that  $Rayleigh = \chi_{dof=2}$

Conclusion on the 2 catalogues case

- Selection criterion

$$D_M \leq k_\gamma$$

- $k_\gamma$  defined such as

$$\int_0^{k_\gamma} \chi_{dof=2}(x) dx = \gamma$$

- $\Rightarrow$  Region of acceptance of surface

$$S_{1,2} = \left[ \frac{|\mathbf{V}_1| |\mathbf{V}_2|}{|\mathbf{V}_\Sigma|} \right]^{1/2} \pi k_\gamma^2$$

We easily generalize to  $n$  catalogues:

- Selection criterion

$$D_M = \left[ \sum_{i=1}^n (\vec{\mu}_i - \vec{\mu}_{\Sigma})^T \mathbf{V}_i^{-1} (\vec{\mu}_i - \vec{\mu}_{\Sigma}) \right]^{1/2} < k_{\gamma}$$

- Now

$$D_M \stackrel{H_0}{\sim} \chi_{dof=2(n-1)}$$

- or, equivalently  $D_M^2 \stackrel{H_0}{\sim} \chi_{dof=2(n-1)}^2$
- So  $k_{\gamma}$  now defined such as

$$\int_0^{k_{\gamma}} \chi_{dof=2(n-1)}(x) dx = \gamma$$

- Region of acceptance of volume

$$S_{1,2} = \left[ \prod_{i=1}^n |\mathbf{v}_i| / |\mathbf{v}_\Sigma| \right]^{1/2} V_{2(n+1)}(k)$$

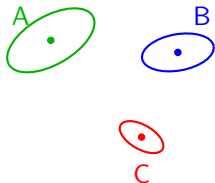
with  $V_{2(n+1)}(k)$ : volume of the  $2(n-1)$ -sphere of radius  $k$



- $\chi_{dof=2(n-1)} \Leftrightarrow (n-1)\chi_{dof=2}$
- $\Rightarrow$  we can perform  $(n-1)$  successives and iteratives xmatches:

$$D_M = \left[ \sum_{i=2}^n (\mu_{\Sigma_{i-1}}^{\vec{}} - \vec{\mu}_i)^T (\mathbf{V}_{\Sigma_{i-1}} + \mathbf{V}_i)^{-1} (\mu_{\Sigma_{i-1}}^{\vec{}} - \vec{\mu}_i) \right]^{1/2}$$

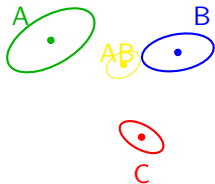
- ▶  $\mu_{\Sigma_{i-1}}^{\vec{}}$ : the weighted mean position of the previous xmatch
- ▶  $\mathbf{V}_{\Sigma_{i-1}}$ : the error on the weighted mean position of the previous xmatch
- Result independant of the xmatch order (for INNER JOIN!)



- $\chi_{dof=2(n-1)} \Leftrightarrow (n-1)\chi_{dof=2}$
- $\Rightarrow$  we can perform  $(n-1)$  successives and iteratives xmatches:

$$D_M = \left[ \sum_{i=2}^n (\mu_{\Sigma_{i-1}}^{\vec{}} - \vec{\mu}_i)^T (\mathbf{V}_{\Sigma_{i-1}} + \mathbf{V}_i)^{-1} (\mu_{\Sigma_{i-1}}^{\vec{}} - \vec{\mu}_i) \right]^{1/2}$$

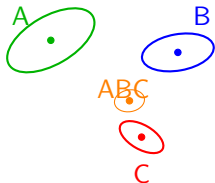
- ▶  $\mu_{\Sigma_{i-1}}^{\vec{}}$ : the weighted mean position of the previous xmatch
- ▶  $\mathbf{V}_{\Sigma_{i-1}}$ : the error on the weighted mean position of the previous xmatch
- Result independant of the xmatch order (for INNER JOIN!)



- $\chi_{dof=2(n-1)} \Leftrightarrow (n-1)\chi_{dof=2}$
- $\Rightarrow$  we can perform  $(n-1)$  successives and iteratives xmatches:

$$D_M = \left[ \sum_{i=2}^n (\mu_{\Sigma_{i-1}}^{\vec{}} - \vec{\mu}_i)^T (\mathbf{V}_{\Sigma_{i-1}} + \mathbf{V}_i)^{-1} (\mu_{\Sigma_{i-1}}^{\vec{}} - \vec{\mu}_i) \right]^{1/2}$$

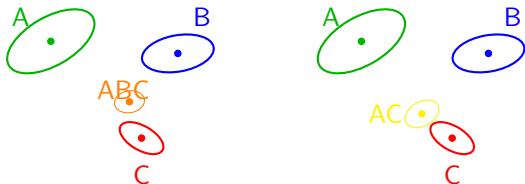
- ▶  $\mu_{\Sigma_{i-1}}^{\vec{}}$ : the weighted mean position of the previous xmatch
- ▶  $\mathbf{V}_{\Sigma_{i-1}}$ : the error on the weighted mean position of the previous xmatch
- Result independant of the xmatch order (for INNER JOIN!)



- $\chi_{dof=2(n-1)} \Leftrightarrow (n-1)\chi_{dof=2}$
- $\Rightarrow$  we can perform  $(n-1)$  successives and iteratives xmatches:

$$D_M = \left[ \sum_{i=2}^n (\mu_{\Sigma_{i-1}}^{\vec{}} - \vec{\mu}_i)^T (\mathbf{V}_{\Sigma_{i-1}} + \mathbf{V}_i)^{-1} (\mu_{\Sigma_{i-1}}^{\vec{}} - \vec{\mu}_i) \right]^{1/2}$$

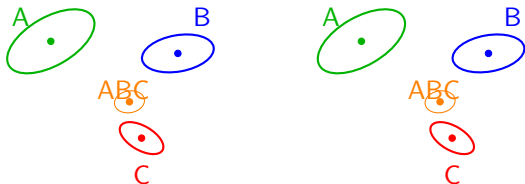
- ▶  $\mu_{\Sigma_{i-1}}^{\vec{}}$ : the weighted mean position of the previous xmatch
- ▶  $\mathbf{V}_{\Sigma_{i-1}}$ : the error on the weighted mean position of the previous xmatch
- Result independant of the xmatch order (for INNER JOIN!)



- $\chi_{dof=2(n-1)} \Leftrightarrow (n-1)\chi_{dof=2}$
- $\Rightarrow$  we can perform  $(n-1)$  successives and iteratives xmatches:

$$D_M = \left[ \sum_{i=2}^n (\mu_{\Sigma_{i-1}}^{\vec{}} - \vec{\mu}_i)^T (\mathbf{V}_{\Sigma_{i-1}} + \mathbf{V}_i)^{-1} (\mu_{\Sigma_{i-1}}^{\vec{}} - \vec{\mu}_i) \right]^{1/2}$$

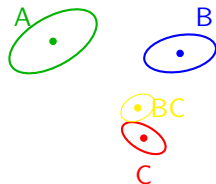
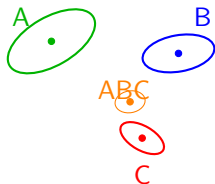
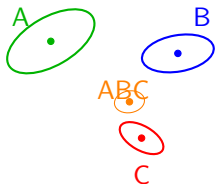
- ▶  $\mu_{\Sigma_{i-1}}^{\vec{}}$ : the weighted mean position of the previous xmatch
- ▶  $\mathbf{V}_{\Sigma_{i-1}}$ : the error on the weighted mean position of the previous xmatch
- Result independant of the xmatch order (for INNER JOIN!)



- $\chi_{dof=2(n-1)} \Leftrightarrow (n-1)\chi_{dof=2}$
- $\Rightarrow$  we can perform  $(n-1)$  successives and iteratives xmatches:

$$D_M = \left[ \sum_{i=2}^n (\mu_{\Sigma_{i-1}}^{\vec{}} - \vec{\mu}_i)^T (\mathbf{V}_{\Sigma_{i-1}} + \mathbf{V}_i)^{-1} (\mu_{\Sigma_{i-1}}^{\vec{}} - \vec{\mu}_i) \right]^{1/2}$$

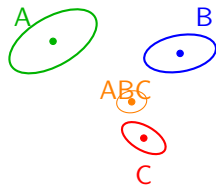
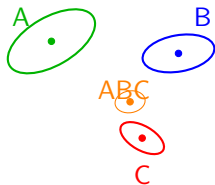
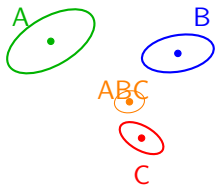
- ▶  $\mu_{\Sigma_{i-1}}^{\vec{}}$ : the weighted mean position of the previous xmatch
- ▶  $\mathbf{V}_{\Sigma_{i-1}}$ : the error on the weighted mean position of the previous xmatch
- Result independant of the xmatch order (for INNER JOIN!)



- $\chi_{dof=2(n-1)} \Leftrightarrow (n-1)\chi_{dof=2}$
- $\Rightarrow$  we can perform  $(n-1)$  successive and iteratives xmatches:

$$D_M = \left[ \sum_{i=2}^n (\mu_{\Sigma_{i-1}}^{\vec{}} - \vec{\mu}_i)^T (\mathbf{V}_{\Sigma_{i-1}} + \mathbf{V}_i)^{-1} (\mu_{\Sigma_{i-1}}^{\vec{}} - \vec{\mu}_i) \right]^{1/2}$$

- ▶  $\mu_{\Sigma_{i-1}}^{\vec{}}$ : the weighted mean position of the previous xmatch
- ▶  $\mathbf{V}_{\Sigma_{i-1}}$ : the error on the weighted mean position of the previous xmatch
- Result independant of the xmatch order (for INNER JOIN!)



To compute Bayes probabilities, we MUST consider all possible hypothesis.

- Law of total probabilities:

$$\sum_{i=1}^n p(H_i) = 1$$

- For 2 catalogues

- ▶ 2 hypothesis

- ★ AB ( $H_0$ )

- ★ A\_B

A  
•

•  
B

- For 3 catalogues
- ▶ 5 hypothesis





To compute Bayes probabilities, we MUST consider all possible hypothesis.

- Law of total probabilities:

$$\sum_{i=1}^n p(H_i) = 1$$

- For 2 catalogues
  - ▶ 2 hypothesis

- ★ AB ( $H_0$ )
- ★ A\_B

A  
•  
B

- For 3 catalogues
  - ▶ 5 hypothesis

A  
•  
B • C

To compute Bayes probabilities, we MUST consider all possible hypothesis.

- Law of total probabilities:

$$\sum_{i=1}^n p(H_i) = 1$$

- For 2 catalogues
  - ▶ 2 hypothesis
    - ★ AB ( $H_0$ )
    - ★ A\_B



- For 3 catalogues
  - ▶ 5 hypothesis



To compute Bayes probabilities, we MUST consider all possible hypothesis.

- Law of total probabilities:

$$\sum_{i=1}^n p(H_i) = 1$$

- For 2 catalogues

- ▶ 2 hypothesis

- ★ AB ( $H_0$ )
- ★ A\_B



- For 3 catalogues

- ▶ 5 hypothesis

- ★ ABC ( $H_0$ )
- ★ AB\_C
- ★ A\_BC
- ★ A\_B\_C
- ★ A\_B\_C



To compute Bayes probabilities, we MUST consider all possible hypothesis.

- Law of total probabilities:

$$\sum_{i=1}^n p(H_i) = 1$$

- For 2 catalogues

- ▶ 2 hypothesis

- ★ AB ( $H_0$ )
- ★ A\_B

A  
•  
  
•  
B

- For 3 catalogues

- ▶ 5 hypothesis

- ★ ABC ( $H_0$ )
- ★ AB\_C
- ★ AC\_B
- ★ A\_BC
- ★ A\_B\_C

A  
•  
• B      • C

To compute Bayes probabilities, we MUST consider all possible hypothesis.

- Law of total probabilities:

$$\sum_{i=1}^n p(H_i) = 1$$

- For 2 catalogues

- ▶ 2 hypothesis

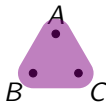
- ★ AB ( $H_0$ )
- ★ A\_B

A  
•  
B  
•

- For 3 catalogues

- ▶ 5 hypothesis

- ★ ABC ( $H_0$ )
- ★ AB\_C
- ★ AC\_B
- ★ A\_BC
- ★ A\_B\_C



To compute Bayes probabilities, we MUST consider all possible hypothesis.

- Law of total probabilities:

$$\sum_{i=1}^n p(H_i) = 1$$

- For 2 catalogues

- ▶ 2 hypothesis

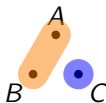
- ★ AB ( $H_0$ )
- ★ A\_B

A  
•  
  
•  
B

- For 3 catalogues

- ▶ 5 hypothesis

- ★ ABC ( $H_0$ )
- ★ AB\_C
- ★ AC\_B
- ★ A\_BC
- ★ A\_B\_C



To compute Bayes probabilities, we MUST consider all possible hypothesis.

- Law of total probabilities:

$$\sum_{i=1}^n p(H_i) = 1$$

- For 2 catalogues

- ▶ 2 hypothesis

- ★ AB ( $H_0$ )
- ★ A\_B

A

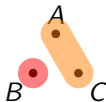
•

B

- For 3 catalogues

- ▶ 5 hypothesis

- ★ ABC ( $H_0$ )
- ★ AB\_C
- ★ AC\_B
- ★ A\_BC
- ★ A\_B\_C



To compute Bayes probabilities, we MUST consider all possible hypothesis.

- Law of total probabilities:

$$\sum_{i=1}^n p(H_i) = 1$$

- For 2 catalogues

- ▶ 2 hypothesis

- ★ AB ( $H_0$ )
- ★ A\_B

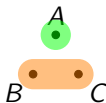
A

B

- For 3 catalogues

- ▶ 5 hypothesis

- ★ ABC ( $H_0$ )
- ★ AB\_C
- ★ AC\_B
- ★ A\_BC
- ★ A\_B\_C





To compute Bayes probabilities, we MUST consider all possible hypothesis.

- Law of total probabilities:

$$\sum_{i=1}^n p(H_i) = 1$$

- For 2 catalogues

- ▶ 2 hypothesis

- ★ AB ( $H_0$ )
- ★ A\_B

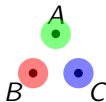
A  
•

•  
B

- For 3 catalogues

- ▶ 5 hypothesis

- ★ ABC ( $H_0$ )
- ★ AB\_C
- ★ AC\_B
- ★ A\_BC
- ★ A\_B\_C



- We generalised for  $n$  catalogues
- The number of hypothesis to be tested is given by the BELL number

Table : Values of the seven first BELL numbers

$n$	2	3	4	5	6	7
$B_n$	2	5	15	52	203	877

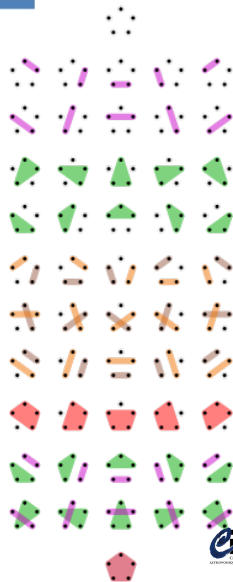
- ▶  $n$  number of catalogues
- ▶  $n=5$  catalogues  $\rightsquigarrow$  52 probabilities to be computed
- $\Rightarrow$  Combinatorial explosion!

- We generalised for  $n$  catalogues
- The number of hypothesis to be tested is given by the BELL number

Table : Values of the seven first BELL numbers

$n$	2	3	4	5	6	7
$B_n$	2	5	15	52	203	877

- ▶  $n$  number of catalogues
- ▶  $n=5$  catalogues  $\rightsquigarrow$  52 probabilities to be computed
- $\Rightarrow$  Combinatorial explosion!

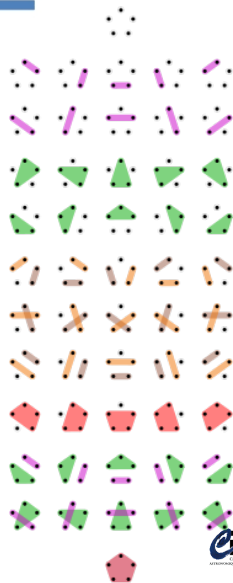


- We generalised for  $n$  catalogues
- The number of hypothesis to be tested is given by the BELL number

Table : Values of the seven first BELL numbers

$n$	2	3	4	5	6	7
$B_n$	2	5	15	52	203	877

- ▶  $n$  number of catalogues
- ▶  $n=5$  catalogues  $\rightsquigarrow$  52 probabilities to be computed
- $\Rightarrow$  Combinatorial explosion!



- Let's call  $x$  the Mahalanobis distance  $D_M$
- Imagine that
  - ▶ we xmatch 2 tables and we obtain  $n_{tot}$  associations
  - ▶ we know the number of spurious associations:  $n_{H_1}$
  - ▶  $\Rightarrow$  we know the number of real associations:  $n_{H_0} = n_{tot} - n_{H_1}$
- Distribution of  $x$  of real associations (likelihood  $p(x|H_0)$ )
  - ▶ almost a  $\chi_{dof=2}(x)$ , because...
  - ▶ ... its integrate over the acceptance domain must be = 1
  - ▶  $\Rightarrow p(x|H_0) = \chi_{dof=2}(x)/\gamma$
- Distribution of  $x$  of spurious associations (likelihood  $p(x|H_1)$ )
  - ▶ Poisson  $\propto x$
  - ▶ again it must integrates to 1 over the domain of acceptance
  - ▶  $\Rightarrow p(x|H_1) = 2x/k_1^2$

- Let's call  $x$  the Mahalanobis distance  $D_M$
- Imagine that
  - ▶ we xmatch 2 tables and we obtain  $n_{tot}$  associations
  - ▶ we know the number of spurious associations:  $n_{H_1}$
  - ▶  $\Rightarrow$  we know the number of real associations:  $n_{H_0} = n_{tot} - n_{H_1}$
- Distribution of  $x$  of real associations (likelihood  $p(x|H_0)$ )
  - ▶ almost a  $\chi_{dof=2}(x)$ , because...
  - ▶ ... its integrate over the acceptance domain must be = 1
  - ▶  $\Rightarrow p(x|H_0) = \chi_{dof=2}(x)/\gamma$
- Distribution of  $x$  of spurious associations (likelihood  $p(x|H_1)$ )
  - ▶ Poisson  $\propto x$
  - ▶ again it must integrates to 1 over the domain of acceptance
  - ▶  $\Rightarrow p(x|H_0) = 2x/k_1^2$

- Let's call  $x$  the Mahalanobis distance  $D_M$
- Imagine that
  - ▶ we xmatch 2 tables and we obtain  $n_{tot}$  associations
  - ▶ we know the number of spurious associations:  $n_{H_1}$
  - ▶  $\Rightarrow$  we know the number of real associations:  $n_{H_0} = n_{tot} - n_{H_1}$
- Distribution of  $x$  of real associations (likelihood  $p(x|H_0)$ )
  - ▶ almost a  $\chi_{dof=2}(x)$ , because...
  - ▶ ... its integrate over the acceptance domain must be = 1
  - ▶  $\Rightarrow p(x|H_0) = \chi_{dof=2}(x)/\gamma$
- Distribution of  $x$  of spurious associations (likelihood  $p(x|H_1)$ )
  - ▶ Poisson  $\propto x$
  - ▶ again it must integrates to 1 over the domain of acceptance
  - ▶  $\Rightarrow p(x|H_0) = 2x/k_1^2$

- Let's call  $x$  the Mahalanobis distance  $D_M$
- Imagine that
  - ▶ we xmatch 2 tables and we obtain  $n_{tot}$  associations
  - ▶ we know the number of spurious associations:  $n_{H_1}$
  - ▶  $\Rightarrow$  we know the number of real associations:  $n_{H_0} = n_{tot} - n_{H_1}$
- Distribution of  $x$  of real associations (likelihood  $p(x|H_0)$ )
  - ▶ almost a  $\chi_{dof=2}(x)$ , because...
  - ▶ ... its integrate over the acceptance domain must be = 1
  - ▶  $\Rightarrow p(x|H_0) = \chi_{dof=2}(x)/\gamma$
- Distribution of  $x$  of spurious associations (likelihood  $p(x|H_1)$ )
  - ▶ Poisson  $\propto x$
  - ▶ again it must integrates to 1 over the domain of acceptance
  - ▶  $\Rightarrow p(x|H_0) = 2x/k_\gamma^2$



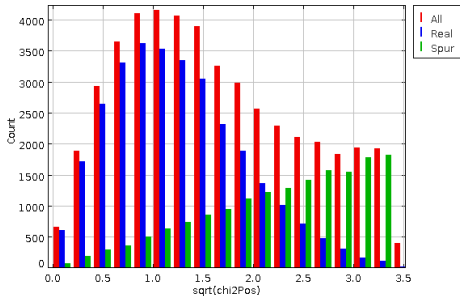
- Let's call  $x$  the Mahalanobis distance  $D_M$
- Imagine that
  - ▶ we xmatch 2 tables and we obtain  $n_{tot}$  associations
  - ▶ we know the number of spurious associations:  $n_{H_1}$
  - ▶  $\Rightarrow$  we know the number of real associations:  $n_{H_0} = n_{tot} - n_{H_1}$
- Distribution of  $x$  of real associations (likelihood  $p(x|H_0)$ )
  - ▶ almost a  $\chi_{dof=2}(x)$ , because...
  - ▶ ... its integrate over the acceptance domain must be = 1
  - ▶  $\Rightarrow p(x|H_0) = \chi_{dof=2}(x)/\gamma$
- Distribution of  $x$  of spurious associations (likelihood  $p(x|H_1)$ )
  - ▶ Poisson  $\propto x$
  - ▶ again it must integrates to 1 over the domain of acceptance
  - ▶  $\Rightarrow p(x|H_0) = 2x/k_\gamma^2$

- Let's call  $x$  the Mahalanobis distance  $D_M$
- Imagine that
  - ▶ we xmatch 2 tables and we obtain  $n_{tot}$  associations
  - ▶ we know the number of spurious associations:  $n_{H_1}$
  - ▶  $\Rightarrow$  we know the number of real associations:  $n_{H_0} = n_{tot} - n_{H_1}$
- Distribution of  $x$  of real associations (likelihood  $p(x|H_0)$ )
  - ▶ almost a  $\chi_{dof=2}(x)$ , because...
  - ▶ ... its integrate over the acceptance domain must be = 1
  - ▶  $\Rightarrow p(x|H_0) = \chi_{dof=2}(x)/\gamma$
- Distribution of  $x$  of spurious associations (likelihood  $p(x|H_1)$ )
  - ▶ Poisson  $\propto x$
  - ▶ again it must integrates to 1 over the domain of acceptance
  - ▶  $\Rightarrow p(x|H_0) = 2x/k_\gamma^2$

- Let's call  $x$  the Mahalanobis distance  $D_M$
- Imagine that
  - ▶ we xmatch 2 tables and we obtain  $n_{tot}$  associations
  - ▶ we know the number of spurious associations:  $n_{H_1}$
  - ▶  $\Rightarrow$  we know the number of real associations:  $n_{H_0} = n_{tot} - n_{H_1}$
- Distribution of  $x$  of real associations (likelihood  $p(x|H_0)$ )
  - ▶ almost a  $\chi_{dof=2}(x)$ , because...
  - ▶ ... its integrate over the acceptance domain must be = 1
  - ▶  $\Rightarrow p(x|H_0) = \chi_{dof=2}(x)/\gamma$
- Distribution of  $x$  of spurious associations (likelihood  $p(x|H_1)$ )
  - ▶ Poisson  $\propto x$
  - ▶ again it must integrates to 1 over the domain of acceptance
  - ▶  $\Rightarrow p(x|H_0) = 2x/k_\gamma^2$

- Blue curve:  $n_{H_0} \times p(x|H_0)$
- Green curve:  $n_{H_1} \times p(x|H_1)$
- Red curve = blue + green
- For an association of given  $x$ :

$$p(H_0|x) = \frac{\text{Blue curve}(x)}{\text{Red curve}(x)}$$



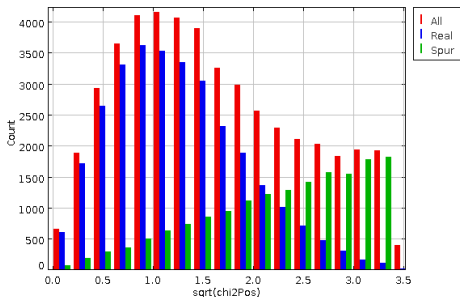
- Bayes formula:

$$p(H_0|x) = \frac{p(H_0)p(x|H_0)}{p(H_0)p(x|H_0) + p(H_1)p(x|H_1)}$$

- Here priors  $p(H_0) = n_{H_0}/n_{tot}$  and  $p(H_1) = n_{H_1}/n_{tot}$

- Blue curve:  $n_{H_0} \times p(x|H_0)$
- Green curve:  $n_{H_1} \times p(x|H_1)$
- Red curve = blue + green
- For an association of given  $x$ :

$$p(H_0|x) = \frac{\text{Blue curve}(x)}{\text{Red curve}(x)}$$



- Bayes formula:

$$p(H_0|x) = \frac{p(H_0)p(x|H_0)}{p(H_0)p(x|H_0) + p(H_1)p(x|H_1)}$$

- Here priors  $p(H_0) = n_{H_0}/n_{tot}$  and  $p(H_1) = n_{H_1}/n_{tot}$

- We know  $n_{tot}$ : number of associations found by the candidate selection criterion
- How to estimate  $n_{H_0}$  or  $n_{H_1}$ ?
  - ▶ one solution is to fit the previous histogram
  - ▶ an analytical solution exists!

- For one random source of catalogue A + one random source of catalogue B distributed on a common surface area S:
  - ▶  $S_{1,2}$ : surface area of the acceptance region
  - ▶ Proba of spurious match:  $S_{1,2}/S$
- For  $n_A$  sources in catalogue A and  $n_B$  sources in catalogue B

$$n_{H_1} = \sum_{i=1}^{n_A} \sum_{j=1}^{n_B} S_{i,j}/S$$

- For circular errors:

$$n_{H_1} = \frac{\pi k \gamma^2}{S} (E\{|V_A|^{1/2}\} + E\{|V_B|^{1/2}\})$$

- ▶  $E\{|V_A|^{1/2}\}$ ,  $E\{|V_B|^{1/2}\}$ : means over cat A and cat B sources respectively
- ▶ Super fast to compute!!

- For one random source of catalogue A + one random source of catalogue B distributed on a common surface area S:
  - ▶  $S_{1,2}$ : surface area of the acceptance region
  - ▶ Proba of spurious match:  $S_{1,2}/S$
- For  $n_A$  sources in catalogue A and  $n_B$  sources in catalogue B

$$n_{H_1} = \sum_{i=1}^{n_A} \sum_{j=1}^{n_B} S_{i,j}/S$$

- For circular errors:

$$n_{H_1} = \frac{\pi k \gamma^2}{S} (E\{|V_A|^{1/2}\} + E\{|V_B|^{1/2}\})$$

- ▶  $E\{|V_A|^{1/2}\}$ ,  $E\{|V_B|^{1/2}\}$ : means over cat A and cat B sources respectively
- ▶ Super fast to compute!!



- For one random source of catalogue A + one random source of catalogue B distributed on a common surface area S:
  - ▶  $S_{1,2}$ : surface area of the acceptance region
  - ▶ Proba of spurious match:  $S_{1,2}/S$
- For  $n_A$  sources in catalogue A and  $n_B$  sources in catalogue B

$$n_{H_1} = \sum_{i=1}^{n_A} \sum_{j=1}^{n_B} S_{i,j}/S$$

- For circular errors:

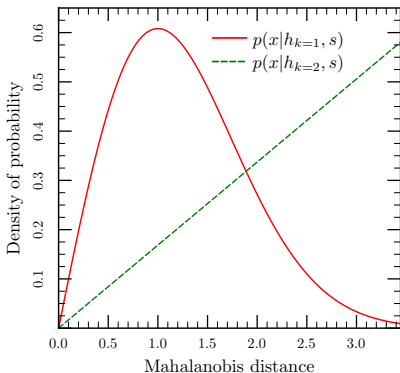
$$n_{H_1} = \frac{\pi k \gamma^2}{S} (E\{|V_A|^{1/2}\} + E\{|V_B|^{1/2}\})$$

- ▶  $E\{|V_A|^{1/2}\}$ ,  $E\{|V_B|^{1/2}\}$ : means over cat A and cat B sources respectively
- ▶ Super fast to compute!!

## Summary for 2 catalogues

- 2 hypotheses
- 2 likelihoods
  - ▶  $H_0 = AB: p(x|H_{AB})$ , Chi of 2 dof
  - ▶  $H_1 = A\_B: p(x|H_{A\_B})$ , Poisson 2D
- 2 priors based on (super fast) geometrical estimates

Likelihoods for  $n = 2$  and  $\gamma = 0.9973$



So now, for 3 catalogues

- 5 hypotheses

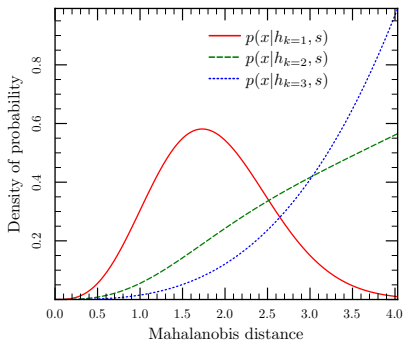
- ▶  $ABC$ , 1 real source
- ▶  $AB\_C$ , 2 real sources
- ▶  $AC\_B$ , 2 real sources
- ▶  $A\_BC$ , 2 real sources
- ▶  $A\_B\_C$ , 3 real sources

- 3 likelihoods

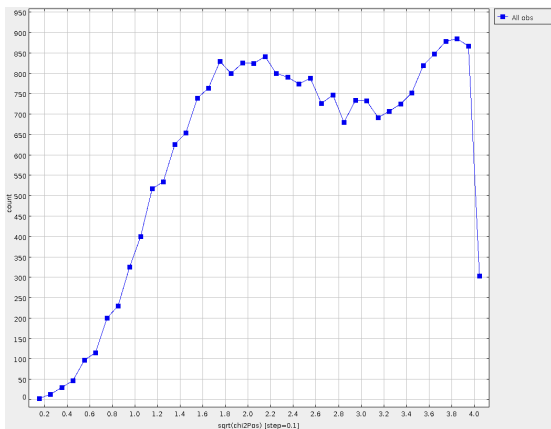
- ▶ 1 likelihood by number of real source
- ▶  $p(x|H_{ABC}) = \chi_{dof=4}(x)/\gamma$ : Chi 4 dof
- ▶  $p(x|H_{AB\_C}) = p(x|H_{AC\_B}) = p(x|H_{A\_BC})$
- ▶  $p(x|H_{A\_B\_C}) = 4x^3/k_\gamma^4$ : Poisson 4D

- And priors?

Likelihoods for  $n = 3$  and  $\gamma = 0.9973$

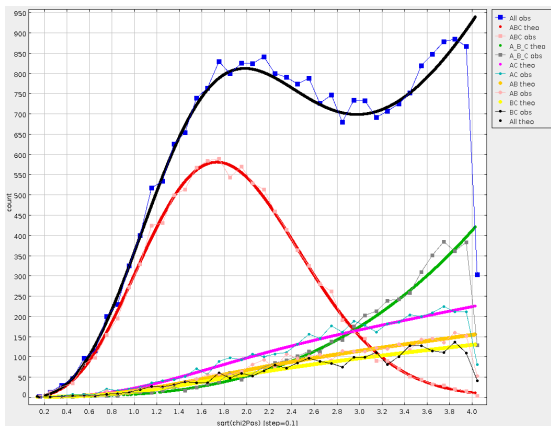


- 5 hypotheses
- $\Rightarrow$  we need 5 priors
  - ▶  $p(H_{ABC}), p(H_{AB..C}), \dots$
- Need 4 xmatches:
  - ▶ A and B  $\rightsquigarrow n_{H_0AB}$
  - ▶ A and C  $\rightsquigarrow n_{H_0AC}$
  - ▶ B and C  $\rightsquigarrow n_{H_0BC}$
  - ▶ A, B and C  $\rightsquigarrow n_{H_0ABC}$
- Having all this, problem solved!



x: Mahalanobis distance  
y: count

- 5 hypotheses
- $\Rightarrow$  we need 5 priors
  - ▶  $p(H_{ABC}), p(H_{AB..C}), \dots$
- Need 4 xmatches:
  - ▶ A and B  $\rightsquigarrow n_{H_0AB}$
  - ▶ A and C  $\rightsquigarrow n_{H_0AC}$
  - ▶ B and C  $\rightsquigarrow n_{H_0BC}$
  - ▶ A, B and C  $\rightsquigarrow n_{H_0ABC}$
- Having all this, problem solved!



x: Mahalanobis distance  
y: count

- This generalises easily for  $n$  catalogues, BUT
  - ▶ Number of hypothesis increases dramatically
  - ▶ Number of priors increases dramatically
  - ▶ Number of xmatches to be performed increases dramatically
- Remark:
  - ▶ instead of computing  $p(H|x)$  one can compute  $p(H|\{\bar{\mu}\}, \{\mathbf{V}\})$
  - ▶ but  $\mathbf{V}$  and magnitudes are NOT independant
  - ▶  $\Rightarrow$  we cannot deals with SEDs separatly (to be investigated)

---

We have been building a general purpose xmatch software implementing all this, but not only.



## Flexible, scalable, efficient

- Flexible and evolutive
  - ▶ dedicated script language
  - ▶ easy to add functionalities
- Efficient
  - tree datastructures
  - multithreading
- Scalable
  - web services (several machines)
  - parallel submission
  - allsky xmatch done cell by cell

## Example of xmatch script

```
# Get XMM sources from a file
get File file=3xmm.fits
where SC_DET_ML<4
set pos ra=SC_RA dec=SC_DEC
set cols *
# Get SDSS DR9 sources from VizierR
get VizierR tabname=V/139/sdss9 mode=cone ...
set pos ra=RAJ2000 dec=DEJ2000
set cols ObjID,/(e_)?[ugriz]mag/,u-g as ug
addmeta ug datatype=float unit=mag ucd=...
# Perform a simple 3" xmatch
xmatch cone dMax=3 join=inner nThreads=48
merge dist mec
# Save intermediary result
save result1.vot votable
# Chain xmatches
get ...
xmatch ...
...
```

## Flexible, scalable, efficient

- Flexible and evolutive
  - ▶ dedicated script language
  - ▶ easy to add functionalities
- Efficient
  - ▶ tree datastructures
  - ▶ multithreading
- Scalable

web services (several machines)  
parallel submission  
allsky xmatch done cell by cell

## Example of xmatch script

```
# Get XMM sources from a file
get File file=3xmm.fits
where SC_DET_ML<4
set pos ra=SC_RA dec=SC_DEC
set cols *
# Get SDSS DR9 sources from VizierR
get VizierR tabname=V/139/sdss9 mode=cone ...
set pos ra=RAJ2000 dec=DEJ2000
set cols ObjID,/(e_)?[ugriz]mag/,u-g as ug
addmeta ug datatype=float unit=mag ucd=...
# Perform a simple 3" xmatch
xmatch cone dMax=3 join=inner nThreads=48
merge dist mec
# Save intermediary result
save result1.vot votable
# Chain xmatches
get ...
xmatch ...
...
```

## Flexible, scalable, efficient

- Flexible and evolutive
  - ▶ dedicated script language
  - ▶ easy to add functionalities
- Efficient
  - ▶ tree datastructures
  - ▶ multithreading
- Scalable
  - ▶ web services (several machines)
  - ▶ parallel submission
  - ▶ allsky xmatch done cell by cell
    - ★ e.g. HEALPix cell
    - ★ one generic script

## Example of xmatch script

```
# Get XMM sources from a file
get File file=3xmm.fits
where SC_DET_ML<4
set pos ra=SC_RA dec=SC_DEC
set cols *
# Get SDSS DR9 sources from VizierR
get VizierR tablename=V/139/sdss9 mode=cone ...
set pos ra=RAJ2000 dec=DEJ2000
set cols ObjID,/(e_)?[ugriz]mag/,u-g as ug
addmeta ug datatype=float unit=mag ucd=...
# Perform a simple 3" xmatch
xmatch cone dMax=3 join=inner nThreads=48
merge dist mec
# Save intermediary result
save result1.vot votable
# Chain xmatches
get ...
xmatch ...
...
```

Algorithm	param	#tbl	prop.mot.	index struct.
chi2 ( $\chi^2$ )	proba	2	$l^1, r^2, b^3$	M/TM-tree

XMatches are chainable: 1  $\chi^2$  xmatch of 4 tables = 3  $\chi^2$  xmatches of 2 tables!

4 to 11 joins (LIRFLIRL/R/F) are supported according to the algorithm.

- <sup>1</sup> l: left table contains extended objects or proper motions;
- <sup>2</sup> r: right table contains extended objects or proper motions;
- <sup>3</sup> b: both left and right tables contain extended objects or proper motion.

Algorithm	param	#tbl	prop.mot.	index struct.
chi2 ( $\chi^2$ )	proba	2	$l^1, r^2, b^3$	M/TM-tree
proba2_vx	proba	2	no (?)	M-tree
proba3_vx	proba	3	no (?)	M-tree
probaN_vx	proba	n	no (?)	M-tree

XMatches are chainable: 1  $\chi^2$  xmatch of 4 tables = 3  $\chi^2$  xmatches of 2 tables!  
 4 to 11 joins (*LIRFLIRL/I'R'F'*) are supported according to the algorithm.

- <sup>1</sup> l: left table contains extended objects or proper motions;
- <sup>2</sup> r: right table contains extended objects or proper motions;
- <sup>3</sup> b: both left and right tables contain extended objects or proper motion.

Algorithm	param	#tbl	prop.mot.	index struct.
chi2 ( $\chi^2$ )	proba	2	$l^1, r^2, b^3$	M/TM-tree
proba2_vx	proba	2	no (?)	M-tree
proba3_vx	proba	3	no (?)	M-tree
probaN_vx	proba	n	no (?)	M-tree
knn	k+dist	2	r, b	kd/M/TM-tree
cone	dist	2	l, r, b	kd/M/TM-tree

XMatches are chainable: 1  $\chi^2$  xmatch of 4 tables = 3  $\chi^2$  xmatches of 2 tables!  
 4 to 11 joins (*LIRFLIRL'I'R'F'*) are supported according to the algorithm.

- <sup>1</sup> l: left table contains extended objects or proper motions;
- <sup>2</sup> r: right table contains extended objects or proper motions;
- <sup>3</sup> b: both left and right tables contain extended objects or proper motion.

Algorithm	param	#tbl	prop.mot.	index struct.
chi2 ( $\chi^2$ )	proba	2	<sup>1</sup> l, <sup>2</sup> r, <sup>3</sup> b	M/TM-tree
proba2_vx	proba	2	no (?)	M-tree
proba3_vx	proba	3	no (?)	M-tree
probaN_vx	proba	n	no (?)	M-tree
knn	k+dist	2	r, b	kd/M/TM-tree
cone	dist	2	l, r, b	kd/M/TM-tree
ext_l <sup>1</sup>	r	2	no	M-tree
ext_r <sup>2</sup>	r	2	no	M-tree
ext_b <sup>3</sup>	r	2	no	M-tree

XMatches are chainable: 1  $\chi^2$  xmatch of 4 tables = 3  $\chi^2$  xmatches of 2 tables!  
 4 to 11 joins (*LIRFL̄RL'I'R'F'*) are supported according to the algorithm.

- <sup>1</sup> l: left table contains extended objects or proper motions;
- <sup>2</sup> r: right table contains extended objects or proper motions;
- <sup>3</sup> b: both left and right tables contain extended objects or proper motion.

Algorithm	param	#tbl	prop.mot.	index struct.
chi2 ( $\chi^2$ )	proba	2	<sup>1</sup> l, <sup>2</sup> r, <sup>3</sup> b	M/TM-tree
proba2_vx	proba	2	no (?)	M-tree
proba3_vx	proba	3	no (?)	M-tree
probaN_vx	proba	n	no (?)	M-tree
knn	k+dist	2	r, b	kd/M/TM-tree
cone	dist	2	l, r, b	kd/M/TM-tree
ext_l <sup>1</sup>	r	2	no	M-tree
ext_r <sup>2</sup>	r	2	no	M-tree
ext_b <sup>3</sup>	r	2	no	M-tree

XMatches are chainable: 1  $\chi^2$  xmatch of 4 tables = 3  $\chi^2$  xmatches of 2 tables!  
 4 to 11 joins ( $LIR\bar{F}\bar{L}\bar{I}\bar{R}'I'R'F'$ ) are supported according to the algorithm.

- <sup>1</sup> l: left table contains extended objects or proper motions;
- <sup>2</sup> r: right table contains extended objects or proper motions;
- <sup>3</sup> b: both left and right tables contain extended objects or proper motion.



- Group XMM FOVs having similar properties (F. Carrera)
  - ▶ tool (STILTS script) dedicated to ARCHES
  - ▶  $\approx 200$  groups in output
    - ★ surface area covered by each group (UNION of FOVs)
- For each group of FOVs
  - ▶ (automatically) write a (quite complex) script
    - ★ another tool dedicated to ARCHES
    - ★ one ARCHES script  $\approx 340$ -800 lines
    - ★ example available during Hands on session
  - ▶ submit the script to the xmatch tool
- Simply concatenates the results

A trap proving (if necessary) that  $n$ -xmatches are complex.

- I want all sources from A having candidates in B or C (or both)
- $\Leftrightarrow$  A inner join (B full join C)
- but imagine, 3 source  $a$ ,  $b$  and  $c$ :
  - ▶  $a$  is  $\chi^2$  compatible with  $b$  but NOT with  $c$
  - ▶  $b$  and  $c$  are  $\chi^2$  compatible
  - ▶ then (B full join C) output contains 1 row:
    - ★ row containing  $b$  and  $c$
  - ▶ then A inner join (B full join C) does not contains any row!
  - ▶ BUT I WANTED the row with  $a$  and  $b$
  - ▶ Solution *ffull join*: result of A ffull join B:
    - ★ row containing  $b$  and  $c$
    - ★ row containing  $b$  alone
    - ★ row containing  $c$  alone
  - ▶ then A inner join (B ffull join C) contains 1 row:
    - ★ row containing  $a$  and  $b$  :)

- Now imagine that  $a$  and  $c$  are  $\chi^2$  compatible and  $a$ ,  $b$  and  $c$  are also  $\chi^2$  compatible
- Then *ffull join* still contains
  - ▶ row containing  $b$  and  $c$
  - ▶ row containing  $b$  alone
  - ▶ row containing  $c$  alone
- But now  $A$  inner join ( $B$  ffull join  $C$ ) contains 3 rows:
  - ▶ row containing  $a$ ,  $b$  and  $c$  :)
  - ▶ row containing  $a$  and  $b$  :(
  - ▶ row containing  $a$  and  $c$  :(
- Need a specific post filter to remove 2 out of 3 rows

- In the framework of ARCHES we have been developing
  - ▶ a flexible, multi-purpose tool
  - ▶ able to xmatch several catalogues in various ways
  - ▶ able to compute probabilities assuming an ideal world
- Provides the basis for an identification process:
  - ▶ needs a layer on top for your particular problem
  - ▶ photometry based proba to be added for proper identification
- Not the end of story:
  - ▶ will be integrated to the CDS XMatch service
  - ▶ opens the road for more complex multi-catalogues proba?
- Script + open tool: good for reproductibility, large parts of the process out of the black box (allows for criticisms), ...