

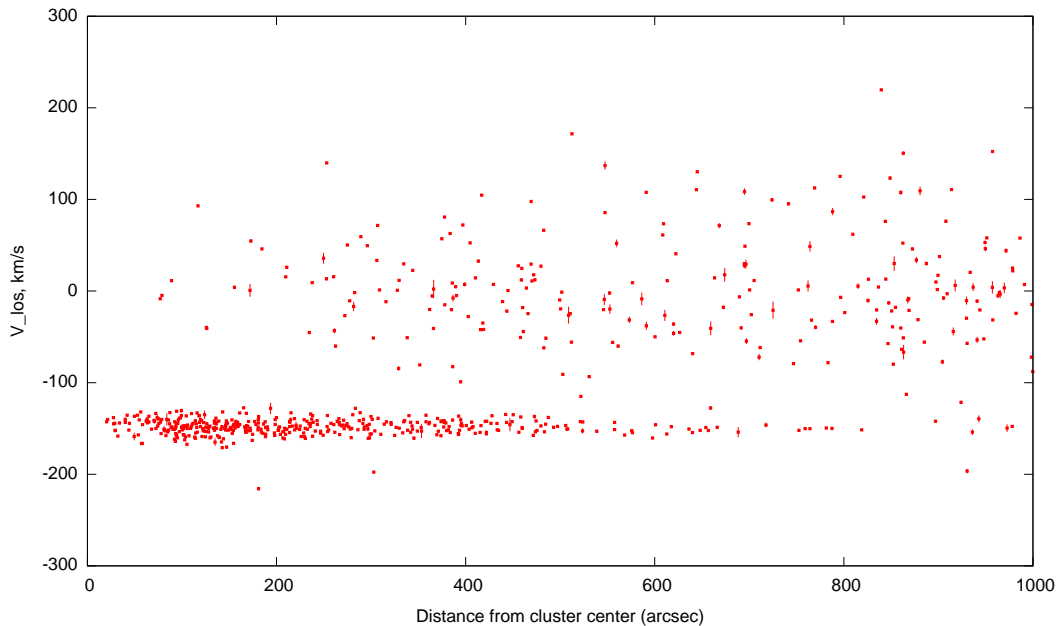
Mixture models and classification

Eugene Vasiliev

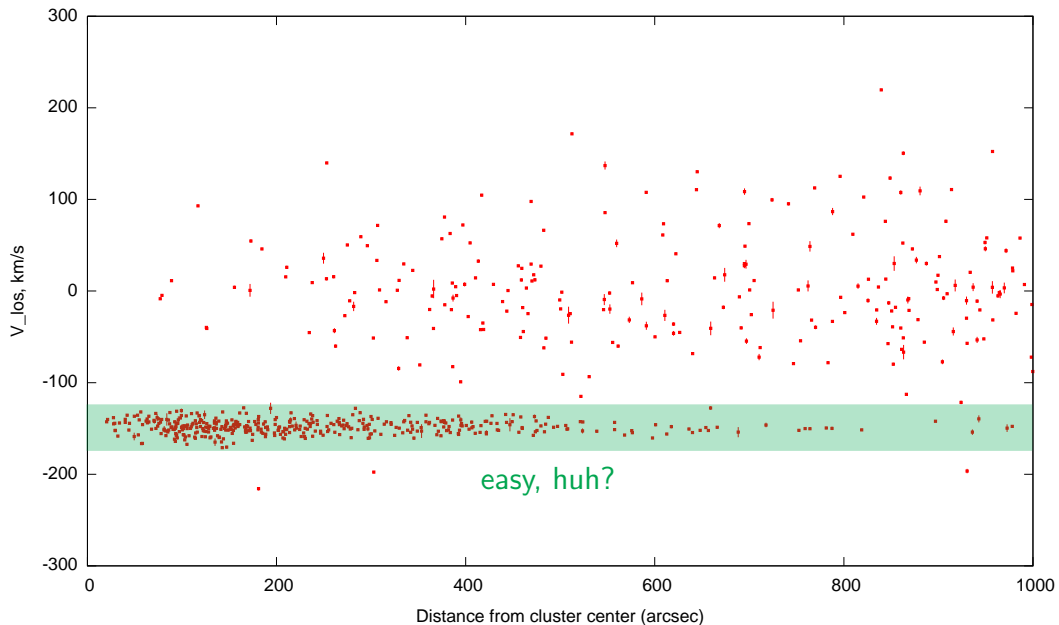
Institute of Astronomy, Cambridge

Heidelberg Physics Graduate Days, October 2020

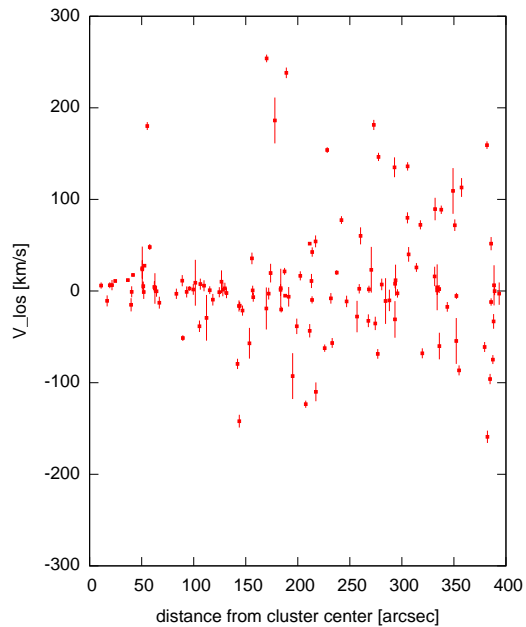
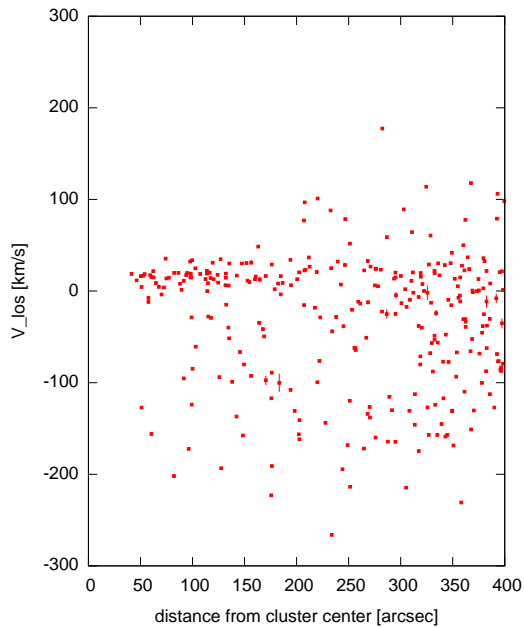
Line-of-sight velocities of globular clusters



Line-of-sight velocities of globular clusters



Line-of-sight velocities of globular clusters



Outline

1. Likelihood-based inference on model parameters
2. Single-component model with measurement errors
3. Mixture model and outlier rejection
4. Mixture model with errors
5. Nonparametric models and further complexities

Food for thought:

Hogg, Bovy & Lang, “Data analysis recipes: fitting a model to data”
([arXiv:1008.4686](https://arxiv.org/abs/1008.4686))

Kuhn & Feigelson, “Mixture models in astronomy” ([arXiv:1711.11101](https://arxiv.org/abs/1711.11101))

Likelihood-based analysis: fit a Gaussian to the data

1. The model: a normal distribution with mean μ and variance σ^2 :

$$x_i \sim \mathcal{N}(\mu, \sigma^2) \equiv \frac{1}{\sqrt{2\pi} \sigma} \exp \left[-\frac{(x_i - \mu)^2}{2 \sigma^2} \right]$$

2. The data: N precisely measured values x_i drawn from this distribution.
3. The likelihood function for the observed dataset given the model :

$$\mathcal{L}(\{x_i\} | \mu, \sigma) = \prod_{i=1}^N \mathcal{N}(x_i | \mu, \sigma^2), \quad \text{or}$$
$$\ln \mathcal{L} = -\frac{1}{2} \sum_{i=1}^N \frac{(x_i - \mu)^2}{\sigma^2} - N \ln \sigma - \frac{N}{2} \ln 2\pi$$

4. Vary the parameters (μ, σ) to maximize $\ln \mathcal{L}$: solve $\left\{ \frac{\partial \mathcal{L}}{\partial \mu} = 0, \frac{\partial \mathcal{L}}{\partial \sigma} = 0 \right\}$

Likelihood-based analysis: fit a Gaussian to the data

4. solve $\left\{ \frac{\partial \mathcal{L}}{\partial \mu} = 0, \frac{\partial \mathcal{L}}{\partial \sigma} = 0 \right\}$:

$$\frac{\partial \mathcal{L}}{\partial \mu} = \sum_{i=1}^N \frac{x_i - \mu}{\sigma^2} = 0 \quad \Longrightarrow \quad \mu = \frac{1}{N} \sum_{i=1}^N x_i,$$

$$\frac{\partial \mathcal{L}}{\partial \sigma} = \sum_{i=1}^N \frac{(x_i - \mu)^2}{\sigma^3} - \frac{N}{\sigma} = 0 \quad \Longrightarrow \quad \sigma = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2.$$

5. The likelihood surface near its maximum is a paraboloid, and the covariance matrix of the uncertainties on model parameters is

$$C = \begin{pmatrix} \delta\mu^2 & \rho \delta\mu \delta\sigma \\ \rho \delta\mu \delta\sigma & \delta\sigma^2 \end{pmatrix} = - \begin{pmatrix} \frac{\partial^2 \ln \mathcal{L}}{\partial \mu^2} & \frac{\partial^2 \ln \mathcal{L}}{\partial \mu \partial \sigma} \\ \frac{\partial^2 \ln \mathcal{L}}{\partial \mu \partial \sigma} & \frac{\partial^2 \ln \mathcal{L}}{\partial \sigma^2} \end{pmatrix}^{-1}$$

Inferring the intrinsic dispersion from imprecise measurements

1. The model for the intrinsic distribution and the measurement process:
the true values are drawn from a Gaussian with mean μ and variance σ^2 ,
the observed values x_i are further perturbed by measurement errors δx_i

$$x_i \sim \underbrace{\mathcal{N}(\mu, \sigma^2) * \mathcal{N}(0, \delta x_i^2)} = \mathcal{N}(\mu, \sigma_i^2), \quad \sigma_i^2 = \sigma^2 + \delta x_i^2$$

convolution of two Gaussians is also a Gaussian

2. Write down the likelihood function $\mathcal{L}(\mu, \sigma \mid \{x_i, \delta x_i\})$:

$$\ln \mathcal{L} = -\frac{1}{2} \sum_{i=1}^N \frac{(x_i - \mu)^2}{\sigma^2 + \delta x_i^2} - \frac{1}{2} \sum_{i=1}^N \ln(\sigma^2 + \delta x_i^2) - \frac{N}{2} \ln 2\pi$$

3. Vary the parameters (μ, σ) to maximize $\ln \mathcal{L}$: solve $\left\{ \frac{\partial \mathcal{L}}{\partial \mu} = 0, \frac{\partial \mathcal{L}}{\partial \sigma} = 0 \right\}$ –
no explicit analytic expression, but straightforward to solve numerically
(note that the best-fit σ may be zero if the actual spread of measured values is smaller
than the typical measurement error)

Bayesian formulation

posterior probability
of model parameters

$$\mathcal{P}(\theta | D, M)$$

↑ model parameters
↑ data (measurements)
↑ model

likelihood of measured data given
the model and its parameters

$$\mathcal{P}(D | \theta, M) \mathcal{P}(\theta | M)$$

prior probability of
model parameters

$$\mathcal{P}(D | M)$$

evidence

Bayesian formulation

posterior probability
of model parameters

$\mathcal{P}(\theta | D, M)$

model parameters

data (measurements)

model

likelihood of measured data given
the model and its parameters

$$\mathcal{P}(D | \theta, M) \mathcal{P}(\theta | M)$$

prior probability of
model parameters

$\mathcal{P}(D | M)$

evidence

Posterior is a normalized probability distribution: $\int \mathcal{P}(\theta | D, M) d\theta = 1$,
hence the evidence is a “normalization factor”:

$$\mathcal{P}(D | M) = \int \mathcal{P}(D | \theta, M) \underbrace{\mathcal{P}(\theta | M)}_{\text{often this is a flat prior, i.e. constant [in some range]}} d\theta$$

often this is a flat prior, i.e. constant [in some range]

Evidence is useful only when choosing between alternative models M_1, M_2 .

Treatment of measurement errors

D are measured (observed) data;

T are “true” (intrinsic) values predicted by the model M with parameters θ :
 $\mathcal{P}(T | \theta, M)$ is the predicted distribution of true values;

$\mathcal{P}(D | T)$ is the measurement model: predicted distribution of D given T ;

predicted distribution of observables is a marginalization over the

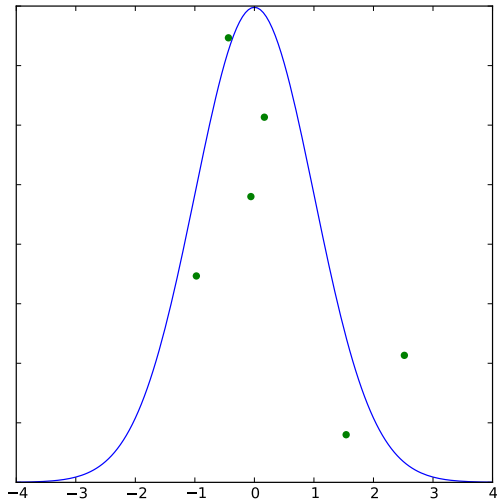
[unknown] true values: $\mathcal{P}(D | \theta, M) = \int \underbrace{\mathcal{P}(D | T) \mathcal{P}(T | \theta, M)}_{\text{convolution with error distribution}} dT$

In the previous example, the model prediction was a normal distribution with two free parameters $\theta = \{\mu, \sigma\}$: $\mathcal{P}(T | \theta, M) = \mathcal{N}(T | \mu, \sigma)$,

and the measurement model was a normal distribution with width δx_i :

$\mathcal{P}(D | T) = \mathcal{N}(D | T, \delta x_i)$.

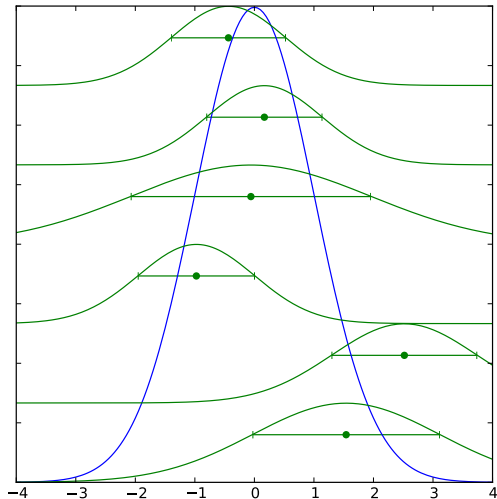
Inferring the intrinsic dispersion



blue: intrinsic distribution

green: true values of sampled points

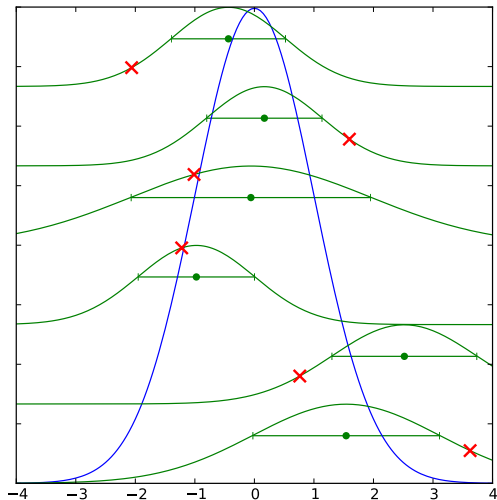
Inferring the intrinsic dispersion



blue: intrinsic distribution

green: true values of sampled points
with measurement uncertainties

Inferring the intrinsic dispersion

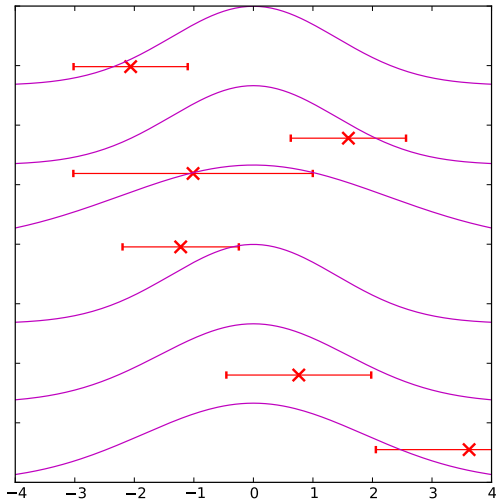


blue: intrinsic distribution

green: true values of sampled points
with measurement uncertainties

red: measured values (perturbed by errors)

Inferring the intrinsic dispersion



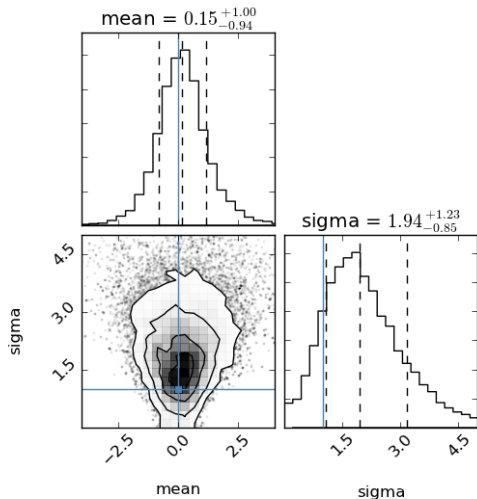
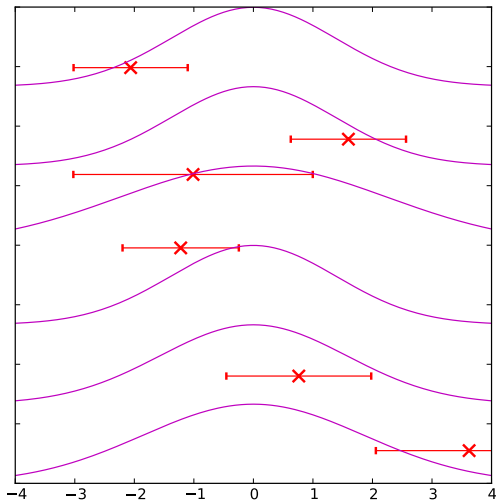
blue: intrinsic distribution

green: true values of sampled points
with measurement uncertainties

red: measured values (perturbed by errors)

magenta: prob.distrib. for each measured point

Inferring the intrinsic dispersion: deconvolution



distribution of inferred model parameters (μ , σ)
obtained from a Markov Chain Monte Carlo run

Membership determination from mixture modelling

N stars with observed properties \mathbf{x}_i , $i = 1..N$
(e.g., position, velocity, colours, etc.);

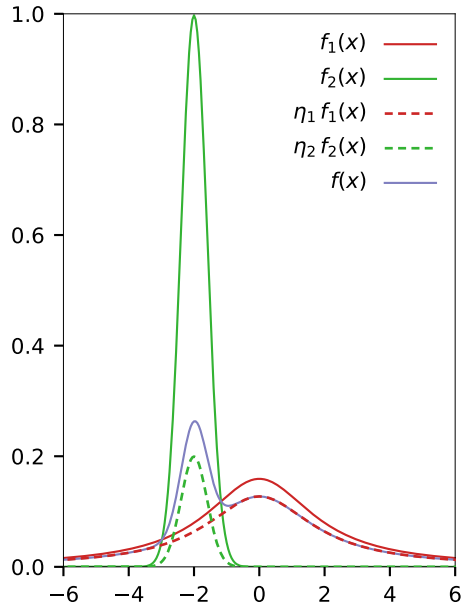
The entire dataset consists of C populations;
 i -th star belongs to the component with
index a_i .

The distribution of values of \mathbf{x} among
stars of c -th component is described
by probability distributions $f_c(\mathbf{x} | \boldsymbol{\theta})$
with some (unknown) parameters $\boldsymbol{\theta}$,
normalized to unity: $\int f_c(\mathbf{x}) d\mathbf{x} = 1$.

The mixture DF is a weighted sum of
component DFs:

$$f(\mathbf{x} | \boldsymbol{\theta}) = \sum_{c=1}^C \eta_c f_c(\mathbf{x} | \boldsymbol{\theta}),$$

where η_c is the fraction of stars in c -th
component, and $\sum_{c=1}^C \eta_c = 1$.



Membership determination from mixture modelling

The log-likelihood of the observed dataset, given the model parameters $(\boldsymbol{\theta}, \boldsymbol{\eta})$, is

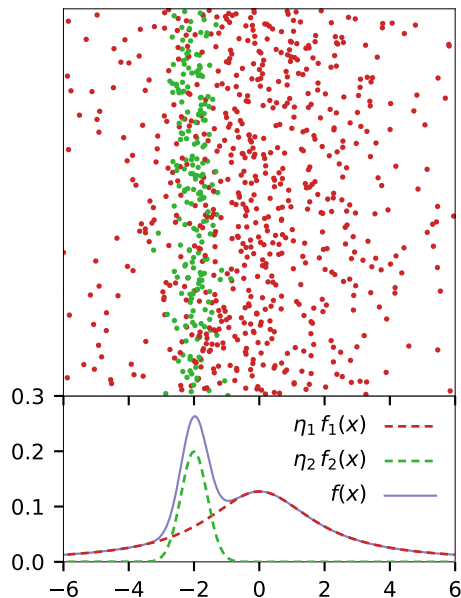
$\ln \mathcal{L} = \sum_{i=1}^N \ln \mathcal{L}_i$, where

$$\mathcal{L}_i \equiv f_{a_i}(\mathbf{x}_i | \boldsymbol{\theta}) = \sum_{c=1}^C \delta_{c a_i} f_c(\mathbf{x}_i | \boldsymbol{\theta}).$$

However, since we do not know the indices a_i , we use the mixture DF:

$$\mathcal{L}_i \equiv f(\mathbf{x}_i | \boldsymbol{\theta}) = \sum_{c=1}^C \eta_c f_c(\mathbf{x}_i | \boldsymbol{\theta}).$$

As usual, the best-fit model parameters may be inferred by maximizing $\ln \mathcal{L}$ (optionally with some priors).

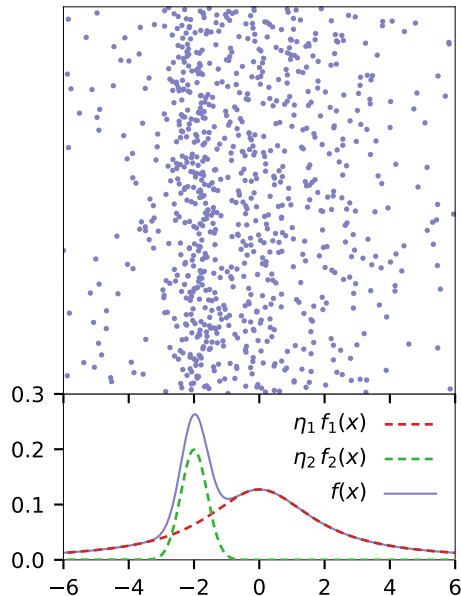


Membership determination from mixture modelling

Assume first that we know the parameters for all DFs θ and their fractions η_c , but do not know which star belongs to which component.

η_c are prior membership probabilities (identical for all stars), while the posterior probabilities for i -th star with measured properties \mathbf{x}_i are

$$p_i^{(c)} = \frac{\eta_c f_c(\mathbf{x}_i | \theta)}{\sum_{k=1}^C \eta_k f_k(\mathbf{x}_i | \theta)} .$$

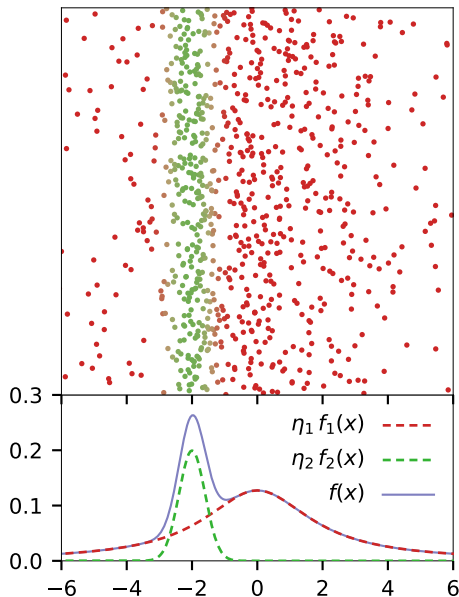


Membership determination from mixture modelling

Assume first that we know the parameters for all DFs θ and their fractions η_c , but do not know which star belongs to which component.

η_c are prior membership probabilities (identical for all stars), while the posterior probabilities for i -th star with measured properties \mathbf{x}_i are

$$p_i^{(c)} = \frac{\eta_c f_c(\mathbf{x}_i | \theta)}{\sum_{k=1}^C \eta_k f_k(\mathbf{x}_i | \theta)} .$$



Membership determination from mixture modelling

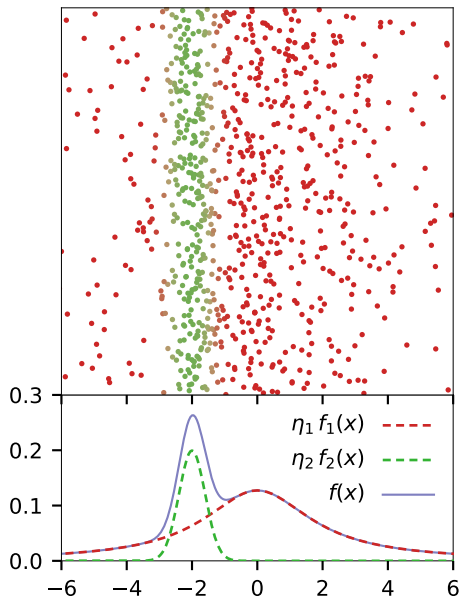
Assume first that we know the parameters for all DFs θ ~~and their fractions η_c~~ , but do not know which star belongs to which component.

η_c are prior membership probabilities (identical for all stars), while the posterior probabilities for i -th star with measured properties \mathbf{x}_i are

$$p_i^{(c)} = \frac{\eta_c f_c(\mathbf{x}_i | \theta)}{\sum_{k=1}^C \eta_k f_k(\mathbf{x}_i | \theta)}.$$

At the same time, $\eta_c = \frac{1}{N} \sum_{i=1}^N p_i^{(c)}$,

so the fractions can be computed alongside membership probabilities.



Membership determination from mixture modelling

Now that we know (probabilistically) the membership of each point $p_i^{(c)}$, we may update the parameters of the DFs θ :

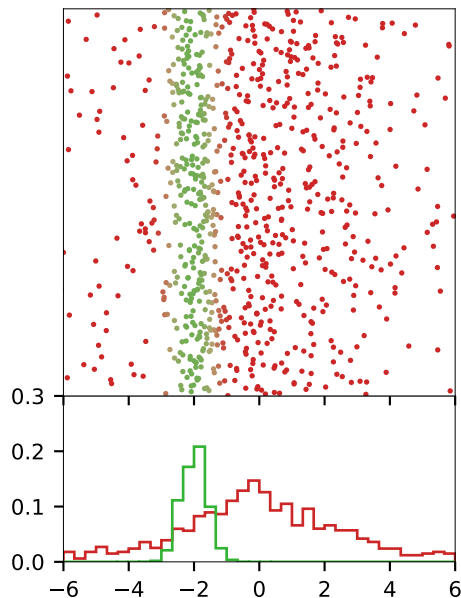
$$\begin{aligned}\theta^{(new)} &= \arg \max_{\theta} (\ln \mathcal{L}) \\ &= \arg \max_{\theta} \left(\sum_{i=1}^N \sum_{c=1}^C p_i^{(c)} \ln f_c(\mathbf{x}_i | \theta) \right).\end{aligned}$$

Fit the parameters θ of each DF f_c to the measured values \mathbf{x}_i , weighted by probabilities $p_i^{(c)}$.

$f_c(\mathbf{x})$ may have any suitable functional form:

- a Gaussian (θ are the mean and dispersion);
- a histogram (θ are the bin heights);
- ...

Repeat these steps until convergence:
this is the expectation/maximization algorithm.



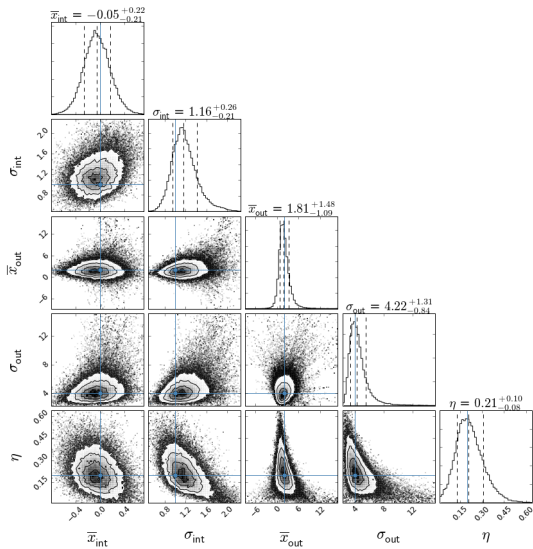
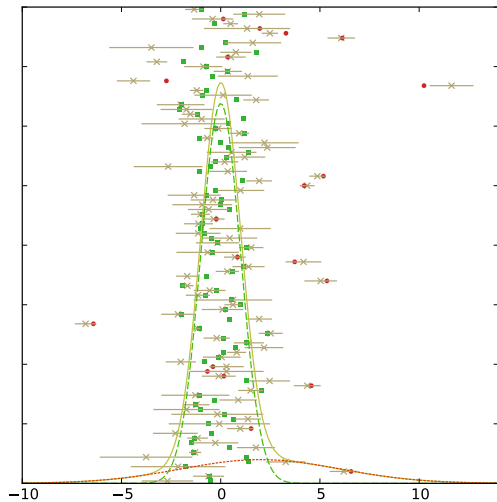
Membership determination from mixture modelling

Expectation/maximization (EM) is one possible way of finding the maximum-likelihood solution for the mixture DF – it gives only the best-fit solution, but no associated uncertainty on η_c and θ_c . One can compute them from the Hessian of the likelihood function $\mathcal{L}(\{\mathbf{x}_i\} \mid \boldsymbol{\eta}, \boldsymbol{\theta})$ at its maximum, or by running a MCMC simulation instead of the EM algorithm. In the latter case, it is also straightforward to marginalize over the nuisance parameters.

Outlier rejection and model fitting with imprecise data

1. Assume a two-component model specified by two intrinsic distributions, convolved with individual measurement errors for each datapoint:
 - ▶ points belonging to the object of interest: $x_i \sim \mathcal{F}_{\text{int}}(\boldsymbol{\theta}) * \mathcal{N}(0, \delta x_i)$
 - ▶ outliers: $x_i \sim \mathcal{F}_{\text{out}}(\boldsymbol{\zeta}) * \mathcal{N}(0, \delta x_i)$ could be a Gaussian, but not necessarily $\boldsymbol{\theta}, \boldsymbol{\zeta}$ are the parameters of the intrinsic distributions (e.g., mean and width of Gaussians)
2. Assume that an (unknown) fraction η of all datapoints are outliers: the probability distribution of the mixture model for datapoint x_i is $x_i \sim \mathcal{F}_{\text{mix},i} \equiv [(1 - \eta) \mathcal{F}_{\text{int}}(\boldsymbol{\theta}) + \eta \mathcal{F}_{\text{out}}(\boldsymbol{\zeta})] * \mathcal{N}(0, \delta x_i)$, and the likelihood of the entire model is $\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\zeta}, \eta \mid \{x_i, \delta x_i\}) = \prod_{i=1}^N \mathcal{F}_{\text{mix},i}$
3. Define suitable priors \mathcal{P} for the nuisance parameters $\boldsymbol{\zeta}, \eta$
4. Obtain the posterior probability distribution for the parameters of interest $\boldsymbol{\theta}$ by marginalizing over the nuisance parameters $\boldsymbol{\zeta}, \eta$:
$$\mathcal{P}(\boldsymbol{\theta} \mid \{x_i\}) = \int \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\zeta}, \eta \mid \{x_i\}) \mathcal{P}(\boldsymbol{\zeta}, \eta) d\boldsymbol{\zeta} d\eta;$$
determine the confidence intervals for $\boldsymbol{\theta}$ from this posterior distribution

Example: measuring the dispersion and pruning outliers



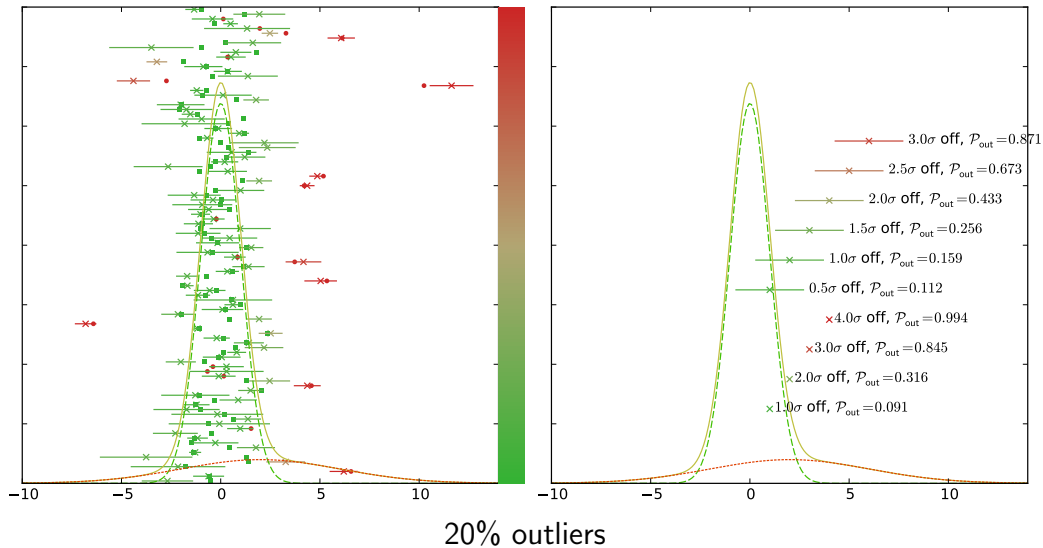
Mixture models and classification

- ▶ Each datapoint has the same prior probability η of being an outlier; however, the posterior probability does depend on the measured value x_i and its uncertainty δx_i , as well as the model parameters θ , ζ :

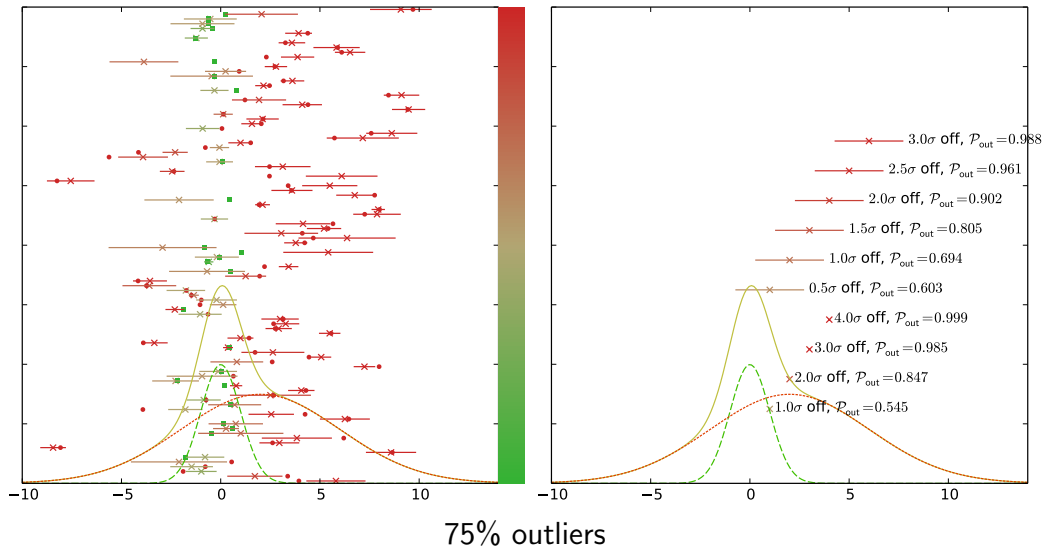
$$\mathcal{P}_{\text{out}}(x_i, \delta x_i \mid \theta, \zeta, \eta) = \frac{[\eta \mathcal{F}_{\text{out}}(\zeta) * \mathcal{N}(0, \delta x_i)](x_i)}{\left[\{(1 - \eta) \mathcal{F}_{\text{int}}(\theta) + \eta \mathcal{F}_{\text{out}}(\zeta)\} * \mathcal{N}(0, \delta x_i) \right](x_i)}$$

- ▶ For the best-fit values of parameters, $\sum_{i=1}^N \mathcal{P}_{\text{out}}(x_i, \delta x_i) = N \eta$
- ▶ There is no single “N- σ ” criterion: if there were no model for outliers, one couldn’t reject a point even when it is 10σ off!
- ▶ The dataset should contain enough contaminants to reliably infer their fraction η and the parameters ζ for \mathcal{F}_{out}
- ▶ Probabilistic membership classification should be carried onward to subsequent modelling procedures, if possible.

Mixture models, classification and rejection of outliers



Mixture models, classification and rejection of outliers

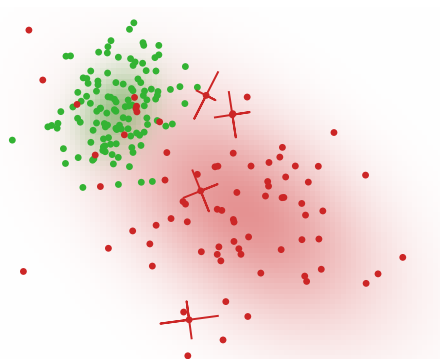


Multidimensional case

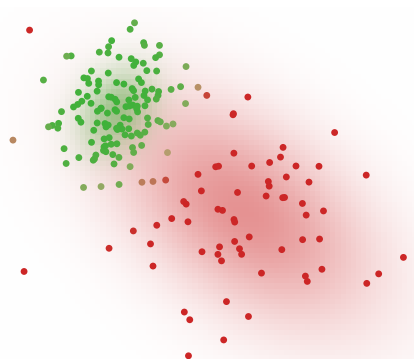
- ▶ D -dimensional Gaussian with mean $\boldsymbol{\mu}$ and covariance matrix Σ :

$$\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \Sigma) = \frac{1}{\sqrt{(2\pi)^D \det \Sigma}} \exp \left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

- ▶ Measurement errors for i -th datapoint described by error covariance matrix δx_i
- ▶ Convolution of two Gaussians is also a Gaussian with covariance $\Sigma + \delta x_i$



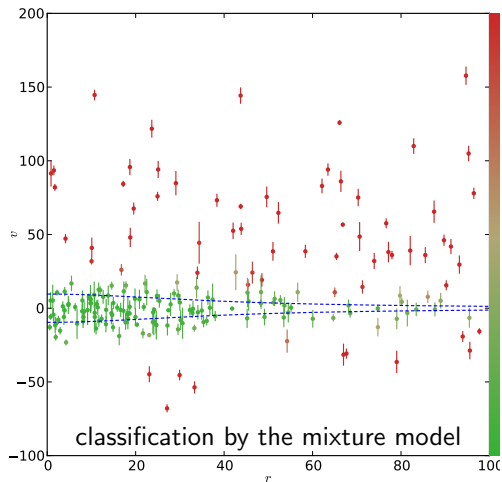
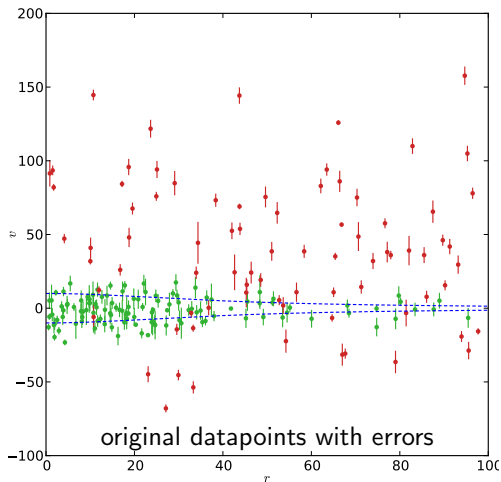
original datapoints with errors



classification by the mixture model

Additional parameters in the model

The probability distribution \mathcal{F}_{int} and the fraction of outliers η may depend on some additional parameters ζ and measured properties $\{\xi_i\}$ (e.g., scale radius a , the distance R_i of a star from the cluster center, etc.)



Fitting dynamical models to discrete-kinematic data

Example: Plummer-like model for the cluster, uniform contamination:

$$\mathcal{F}_{\text{int}}(\{R_i, v_i, \delta v_i\} \mid a, \bar{v}, \sigma_0) = \mathcal{N}\left(v_i \mid \bar{v}, \sqrt{\sigma^2(R_i) + \delta v_i^2}\right) \frac{1 + (R_{\text{max}}/a)^2}{[1 + (R_i/a)^2]^2}$$

scale radius → a cluster mean velocity → \bar{v} central velocity dispersion → σ_0

$$\sigma(R_i \mid a, \sigma_0) \equiv \frac{\sigma_0}{[1 + (R_i/a)^2]^{1/4}}$$

normalized surface density

mean and dispersion of contaminants

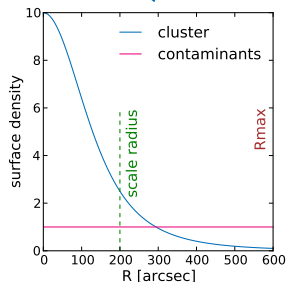
$$\mathcal{F}_{\text{out}}(\{v_i, \delta v_i\} \mid \bar{v}_{\text{out}}, \sigma_{\text{out}}) = \mathcal{N}\left(v_i \mid \bar{v}_{\text{out}}, \sqrt{\sigma_{\text{out}}^2 + \delta v_i^2}\right)$$

Distribution function of the mixture model:

$$\mathcal{F}_{\text{mix}} = (1 - \eta) \mathcal{F}_{\text{int}} + \eta \mathcal{F}_{\text{out}}$$

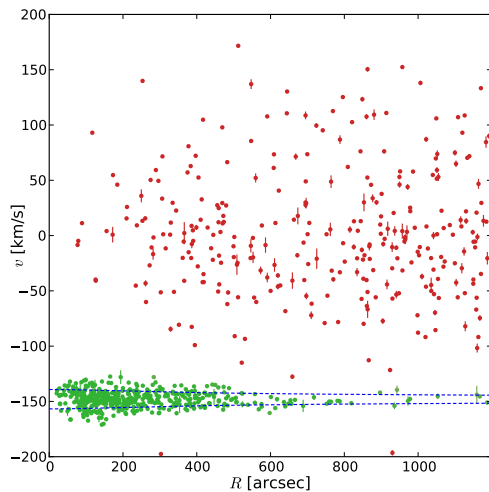
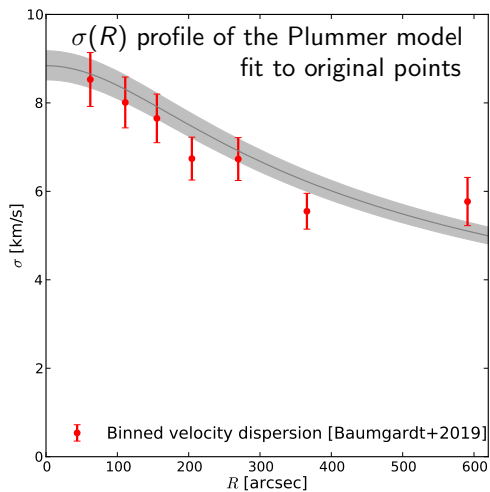
parameters: $a, \bar{v}, \sigma_0, \bar{v}_{\text{out}}, \sigma_{\text{out}}, \eta$

fraction of contaminants ↑ η



Fitting dynamical models to discrete-kinematic data

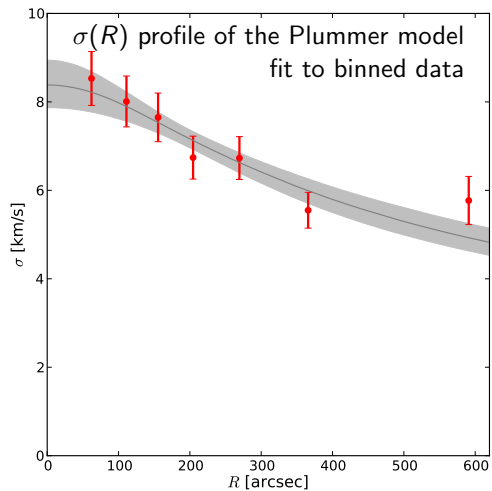
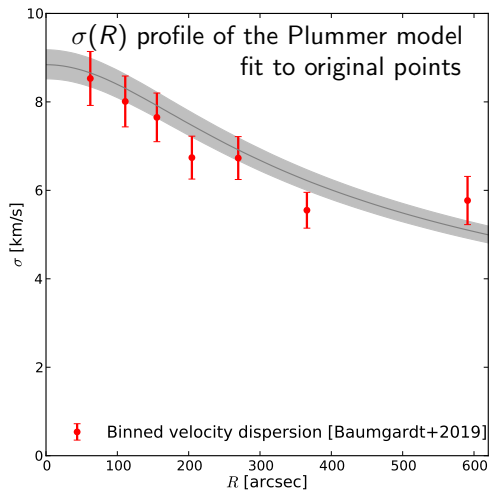
Example: globular cluster NGC 6656



Fitting dynamical models to discrete-kinematic data

Example: globular cluster NGC 6656

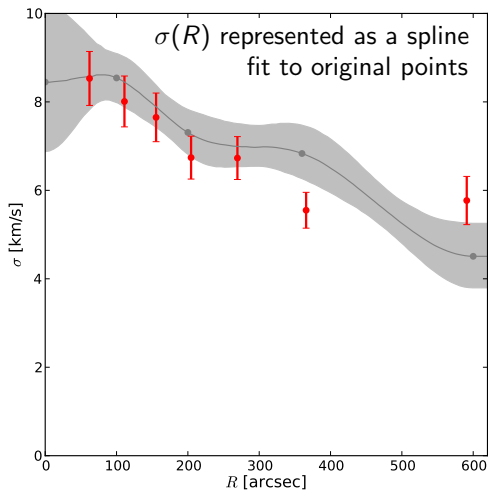
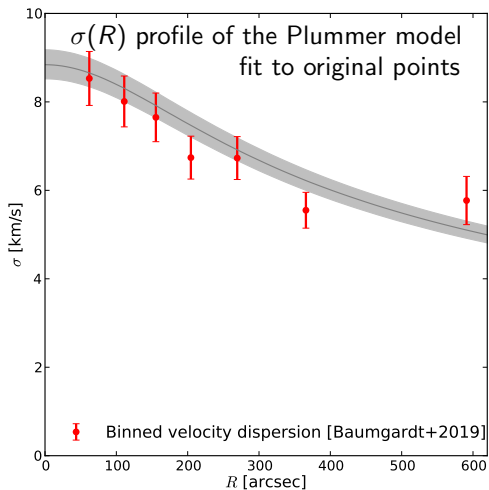
General problem with parametric models:
inferred uncertainty intervals are too small



Fitting “nonparametric” models to discrete-kinematic data

Consider a more flexible (but not dynamically motivated) dispersion profile: $\sigma(R)$ represented as a cubic spline with fixed nodes and free coefficients.

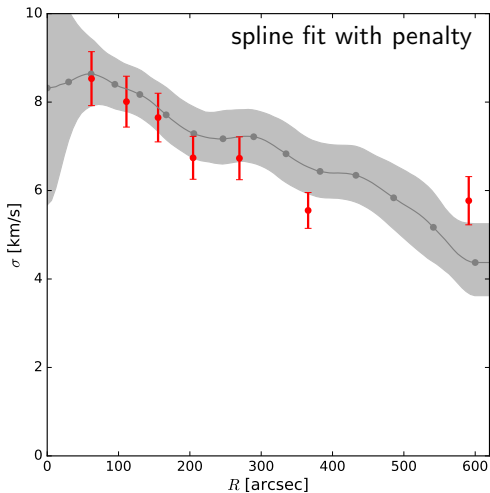
Similar uncertainties as binned data, but in a mathematically consistent model.



Penalization of overly flexible models

If one has too many free parameters, the models would overfit the data.

To prevent this, one needs to impose stronger priors on the parameters, e.g., penalizing large variations between adjacent spline values.



Penalization of overly flexible models

If one has too many free parameters, the models would overfit the data.

To prevent this, one needs to impose stronger priors on the parameters, e.g., penalizing large variations between adjacent spline values.

There are objective methods for determining the optimal value of penalty (smoothing parameter), based on cross-validation:

- ▶ adopt a particular value of the smoothing parameter λ ;
- ▶ split the data sample into two parts – training set and validation set;
- ▶ find the best-fit model for the training set;
- ▶ evaluate the goodness-of-fit of this model for the validation set;
- ▶ average this over different choices of training and validation subsets;
- ▶ adjust the parameter λ to maximize this validation score;

The same approach is used in machine learning to prevent overfitting

Pros and cons of fitting probabilistic mixture models

- + outlier rejection or fitting multiple populations are straightforward
- + easy to account for selection function and measurement uncertainties
- + use all available information (individual datapoints, no binning)
- + can propagate uncertainty in membership into subsequent analysis
- computationally demanding in case of large datasets
- marginalization and convolution often non-analytic \Rightarrow expensive

The most common variant of this technique is the Gaussian mixture modelling, with several implementations available in Python:

- ▶ [scikit-learn](#) (no measurement errors)
- ▶ [Extreme Deconvolution](#) [Bovy, Hogg & Roweis 2011]
- ▶ [XDGMM](#) [Holoien, Marshall & Wechsler 2016] (also in [AstroML](#))
- ▶ [PyGMMis](#) [Melchior & Goulding 2018]

All these variants use the expectation/maximization approach and give only the best-fit parameters, without uncertainties.

Exercise: simple 1d Gaussian mixture models

Scenario: a dataset of N points drawn from a mixture of two populations: broad and narrow.

- ▶ Write a routine for generating the mock dataset (assuming precise measurements).
- ▶ Write another routine for computing the likelihood of the dataset given a model with some free parameters – a sum of two Gaussians with unknown mean, dispersion, and relative weights.
- ▶ Use your favourite method to find the best-fit parameters (e.g., direct maximization of the likelihood function, or the expectation/maximization algorithm, or MCMC).
- ▶ Extend the procedure to the case of measurement errors (preferably, varying between datapoints) – both for the mock data generation and for the fit.
- ▶ Fiducial values: $N \simeq 100$, fraction of broad component (contaminants) between 10 and 90%, measurement uncertainties between 0 and $1 - 2 \times$ the dispersion of the narrow component.

Exercise 2: application to globular clusters

Many globular clusters in the Galaxy have plenty of stars with measured values of line-of-sight velocity; a good starting point is the catalogue of Baumgardt+ 2019: <https://people.smp.uq.edu.au/HolgerBaumgardt/globular/appendix/appendix.html>

The tables for individual clusters contain both members and non-members. For each star, the table contains the value and uncertainty of v_{los} , distance from the cluster centre, and the membership probability – we will not use the latter quantity in the fit, but will compare it with the results of the mixture modelling classifier. Ideally, the distance information can be used to estimate the density profile of the cluster, but the coverage of the spectroscopic dataset is very non-uniform, making this task difficult. We may instead restrict the range of distances to a few arcmin (perhaps depending on the cluster under consideration).

The goal is to measure the mean velocity and its dispersion, and to estimate the membership probability of all stars in the dataset. The results can be compared with the original paper (which used a different approach), shown in the table of [structural parameters](#) (for σ) and [orbits](#) (for $\overline{v_{\text{los}}}$).