# Searching for what no one is looking for

Sebastian Ratzenböck

Data Science @ Uni Vienna

universität
wien

# Searching for what no one is looking for
## Blind searches in Gaia DR2

Sebastian Ratzenböck

Data Science @ Uni Vienna

universität
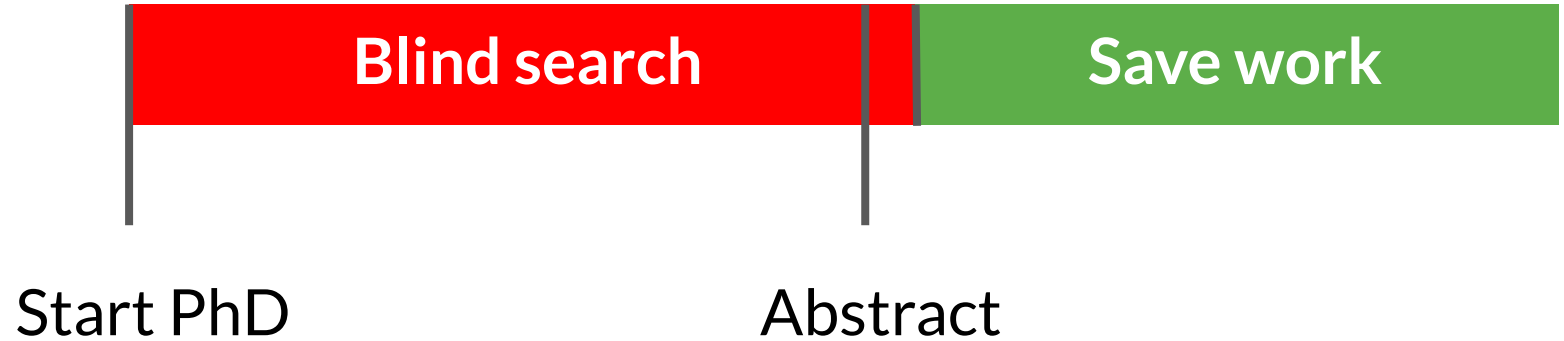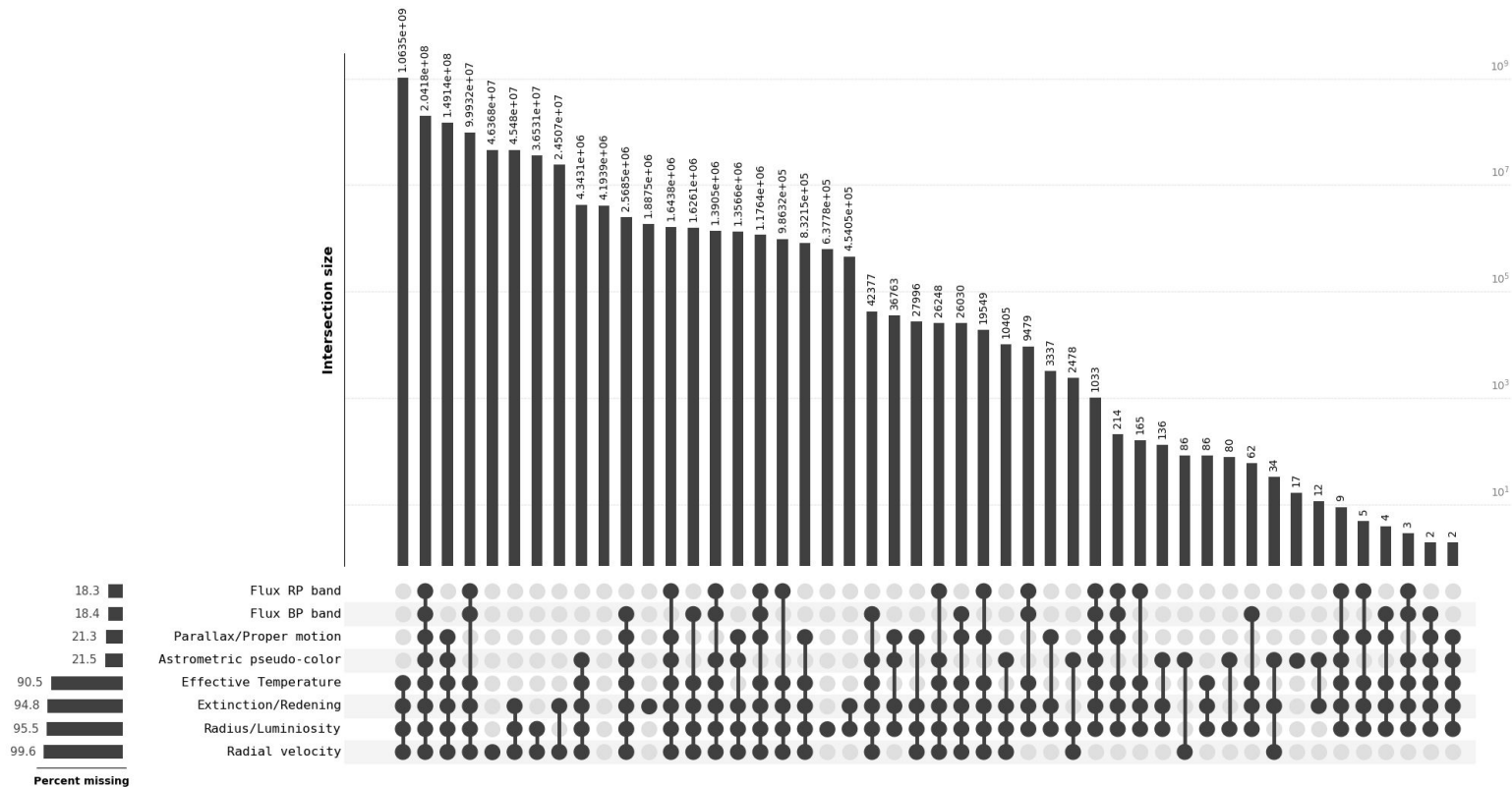wien

AIA, Garching
July 25th, 2019

Start PhD
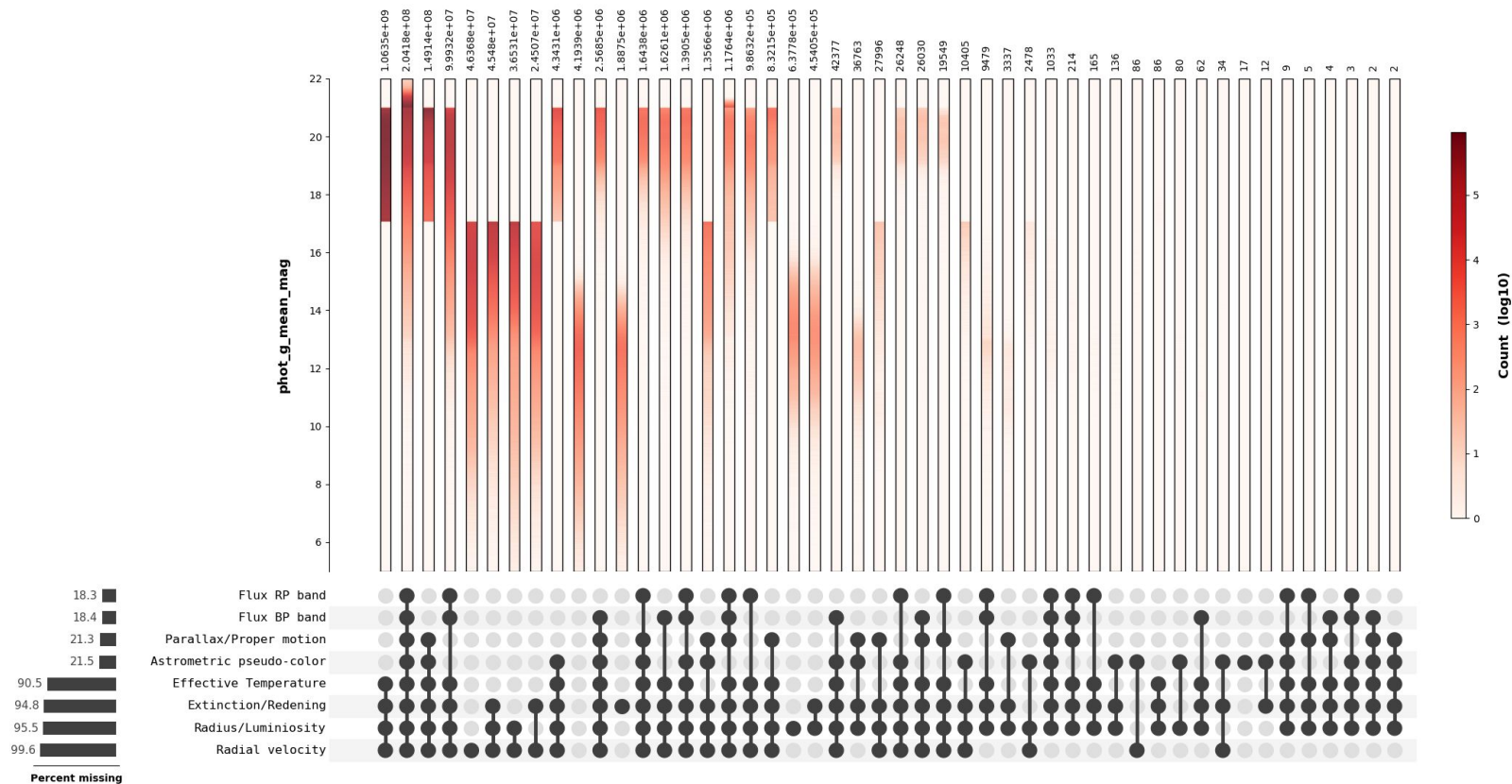
# Part I

# The blind search

# Blind search in Gaia DR2

"Gaia does not exclusively observe stars: *all* objects brighter than G ≈ 20 mag are observed, [...]."

- Gaia: Science Performance

# Data structure: Missing value sets in Gaia

# Subset distributions: G-magnitude

# Clustering pipeline

1.  Give each data point a label

$$f(\vec{x}_i, \vec{\theta}) = s_i \qquad \vec{x}_i \in \mathbb{R}^n, \ \ s_i \in \mathbb{Z}, \ \ \vec{\theta} = (\theta_1, \ldots, \theta_m)$$

2.  Validate clustering
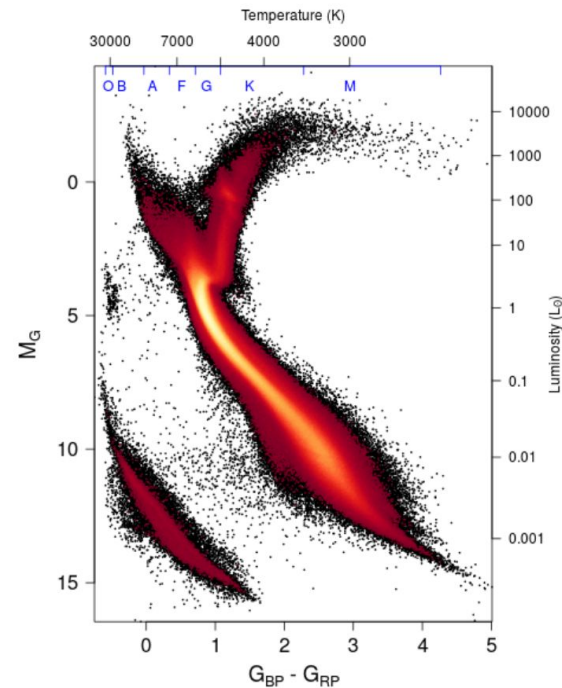
# Clustering *cycle*

1. Give each data point a label

$$f(\vec{x}_i, \vec{\theta}) = s_i \qquad \vec{x}_i \in \mathbb{R}^n, \ \ s_i \in \mathbb{Z}, \ \ \vec{\theta} = (\theta_1, \ldots, \theta_m)$$

2. Validate clustering (then go back to 1.)

# Choosing the algorithm

Requirements

- Deal with non-linearities between variables
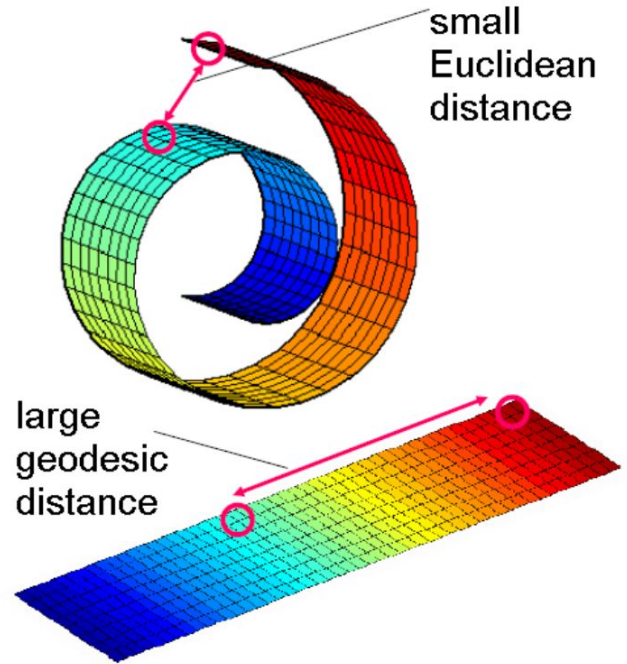  - Flexible model
- Applicable to millions of data points



Source: Gaia Collaboration (2018)

# Reducing dimensions - Manifold learning

If variables **depend** on each other their joint distribution does not span the whole space
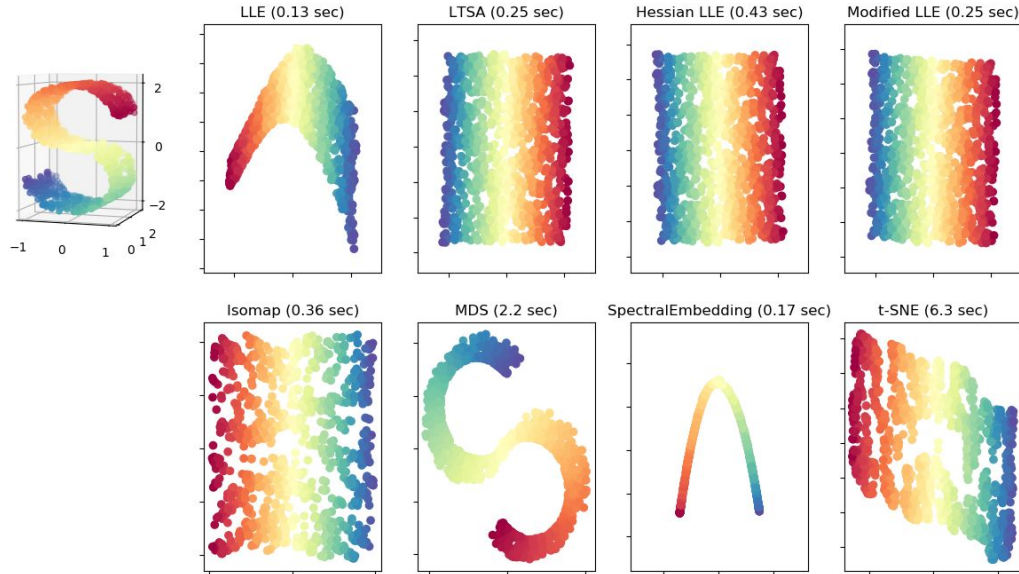
→ data lies on (around) the support of the joint distribution

**Manifold**: underlying support of the data distribution known only through finite sampling



small Euclidean distance

large geodesic distance

# Reducing dimensions - Manifold learning



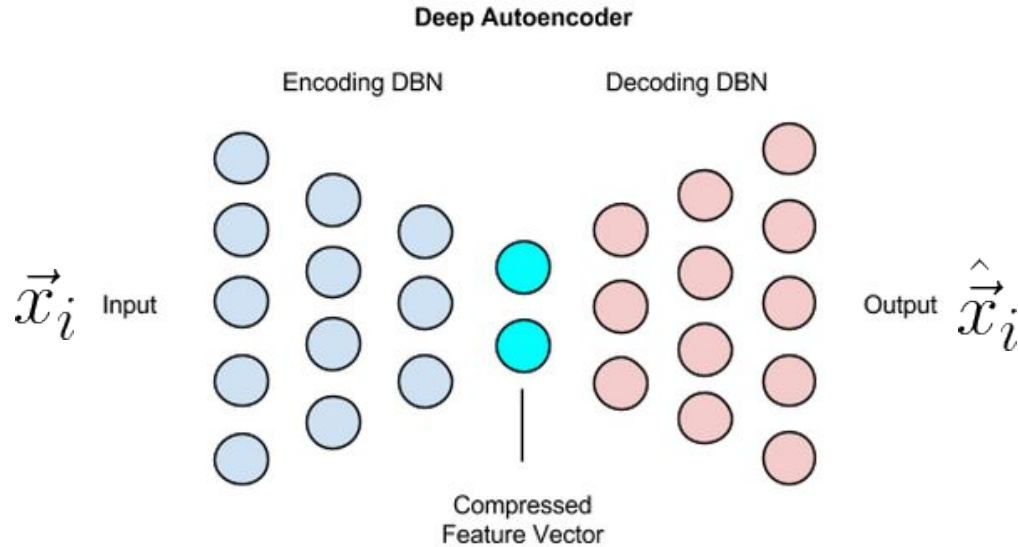Source: Scikit Learn (scikit-learn.org/stable/modules/manifold.html)

# Reducing dimensions - Manifold learning

**BUT:** We generally do **not** know the dimensionality of the intrinsic manifold
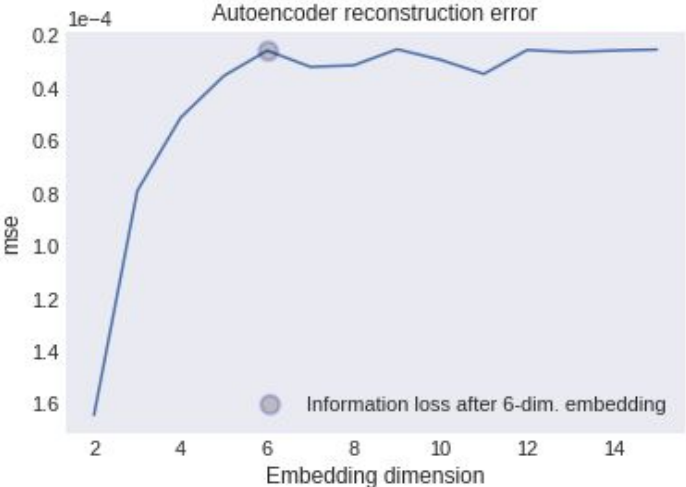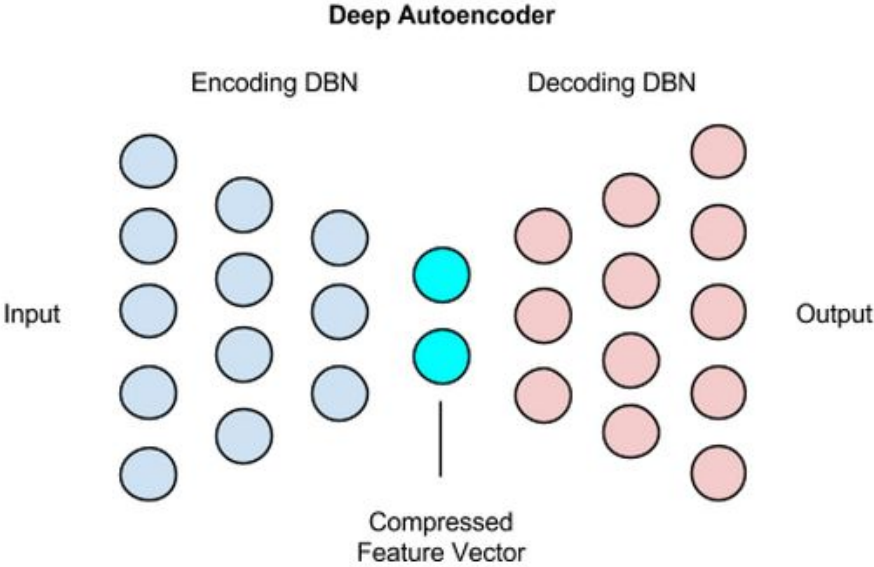
→ Can lose valuable information

# Reducing dimensions - Manifold learning

**BUT:** We generally do **not** know the dimensionality of the intrinsic manifold

→ Can lose valuable information



small Euclidean distance

large geodesic distance

# Reducing dimensions - Autoencoder

**Deep Autoencoder**

Encoding DBN       Decoding DBN

$\vec{x}_i$   Input

Output   $\hat{\vec{x}}_i$

Compressed
Feature Vector

$$Loss = \sum_i ||\vec{x}_i - \hat{\vec{x}}_i||$$

# Reducing dimensions - Autoencoder

# Deep embedded clustering (Xie et al. 2016)

- Use neural net as powerful feature extractor

- Introduce a second training phase where the representation in the mapping to the latent space is optimized for k-Means clustering
  - Set centroids (hyperparameter) in latent space and force points around these centroids to be t-distributed by minimizing KL divergence loss term

# Deep embedded clustering (Xie et al. 2016)



(a) Epoch 0

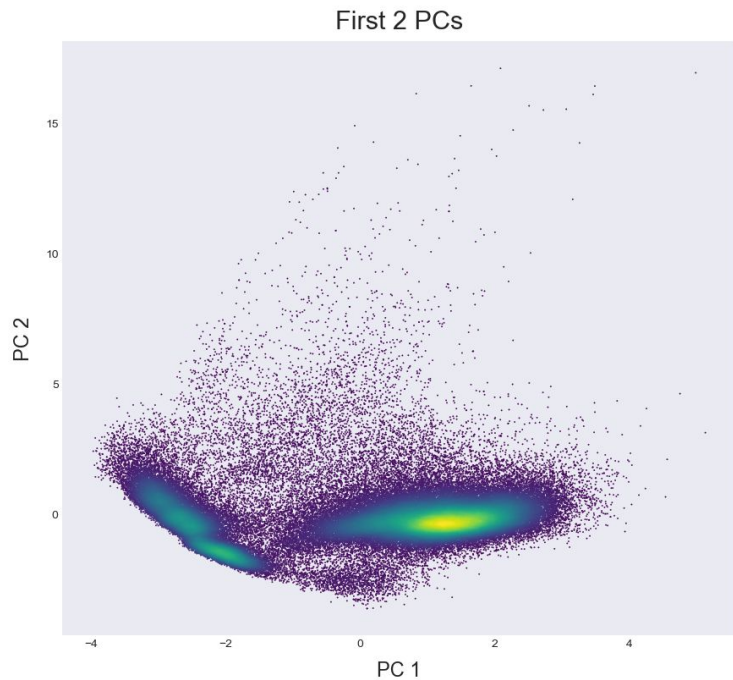(b) Epoch 3

(c) Epoch 6

(d) Epoch 9

(e) Epoch 12

(f) Accuracy vs. epochs

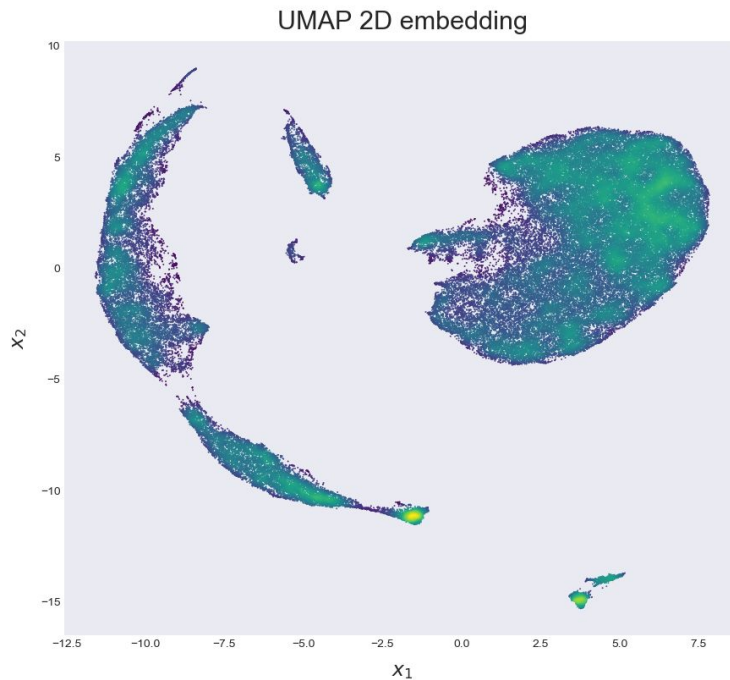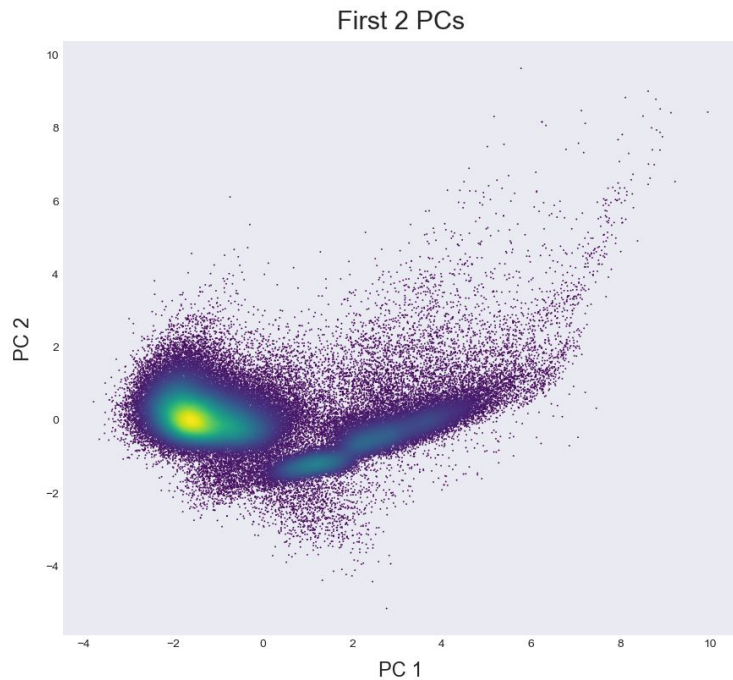# Deep embedded clustering on Gaia data



Epoch 0

# Deep embedded clustering on Gaia data



Epoch 50

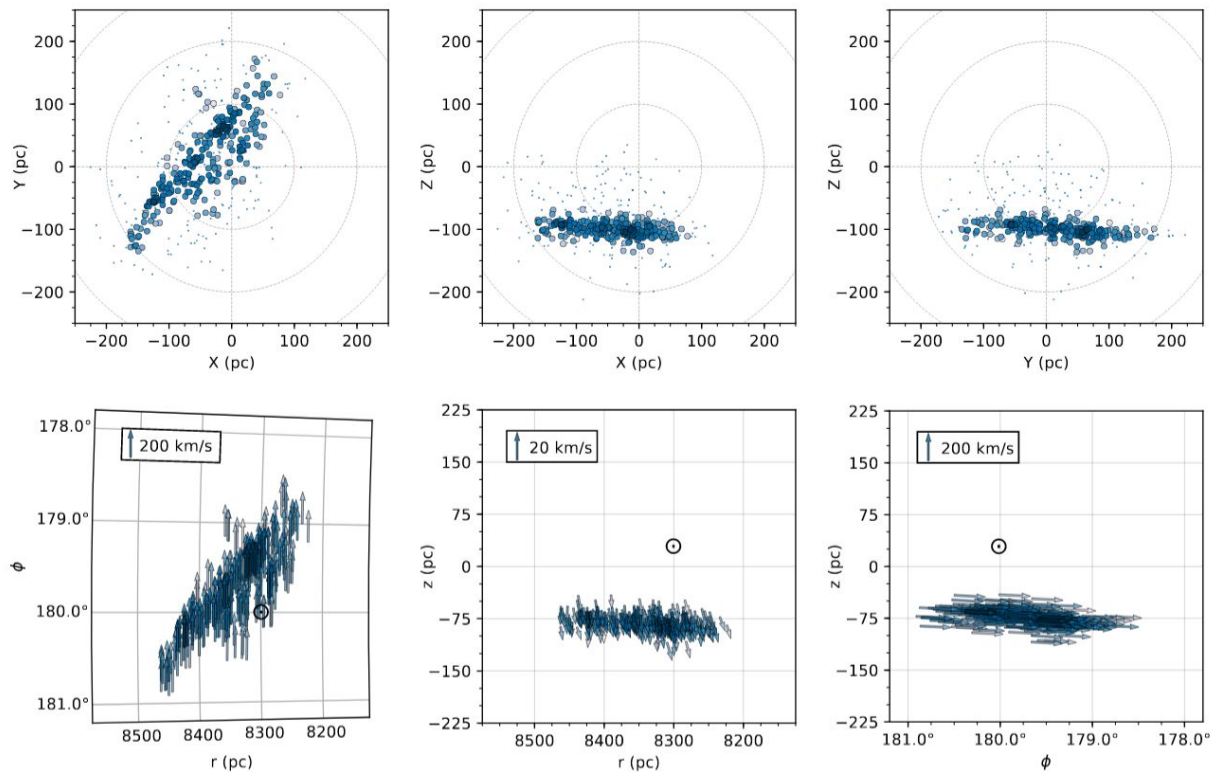# Deep embedded clustering on Gaia data

Epoch 1870

# Part II
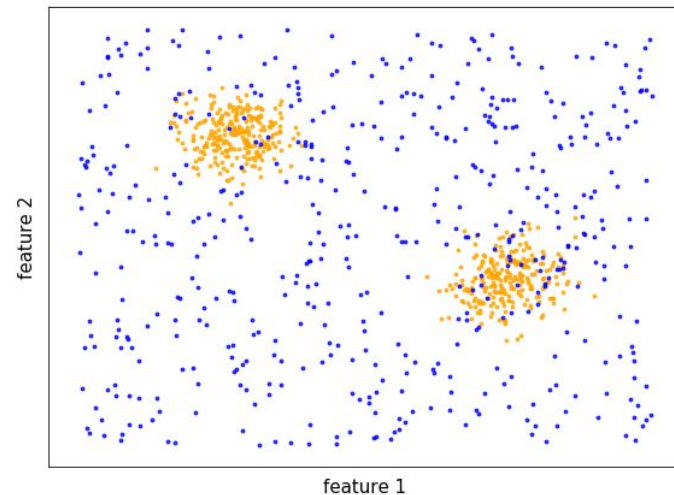
# Feature engineering for robust OC extraction
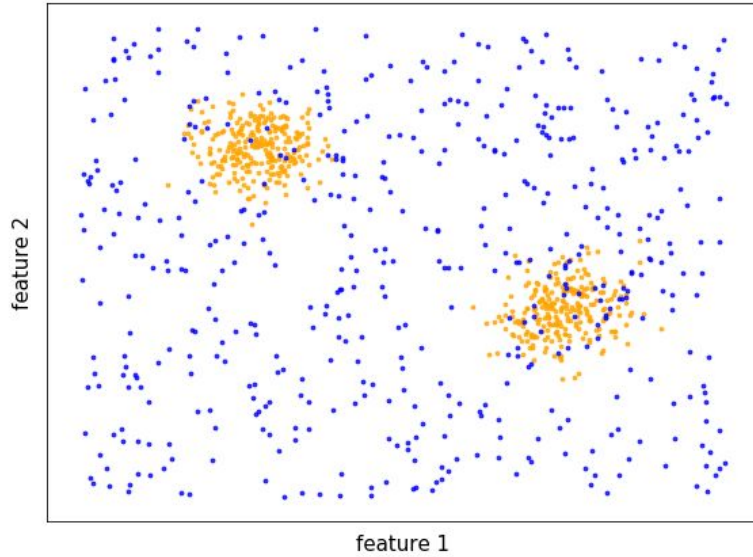
# Stellar clusters



Source: Meingast et al. (2019)
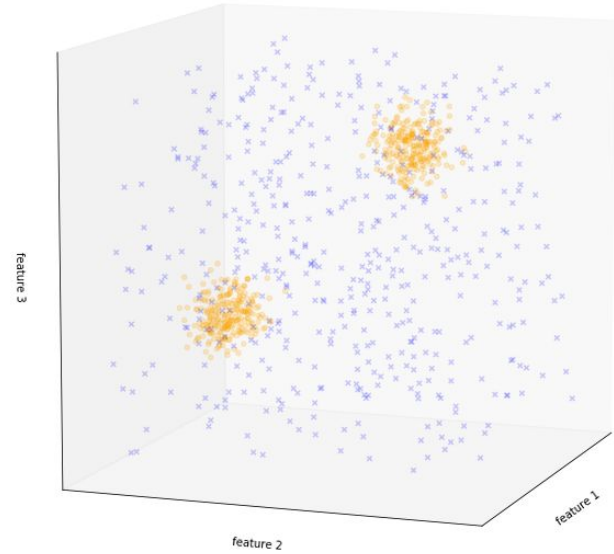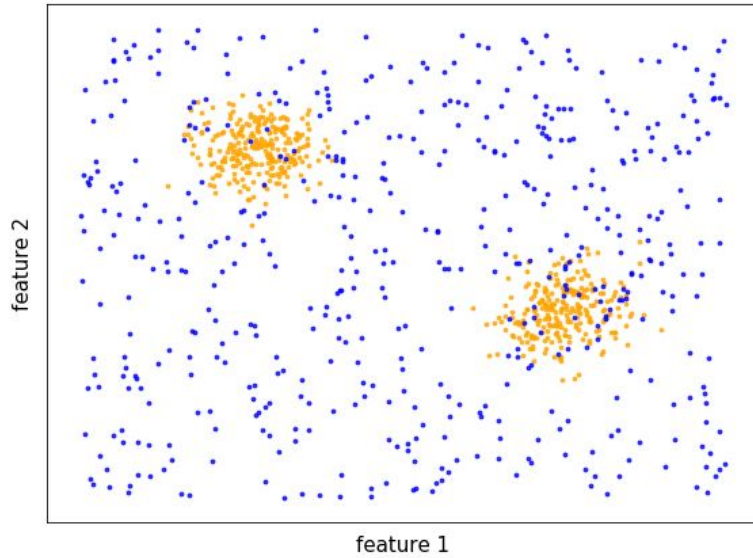
# Stellar clusters - feature space

- 5D feature space: XYZ + PMs

- The feature space is dominated by noise
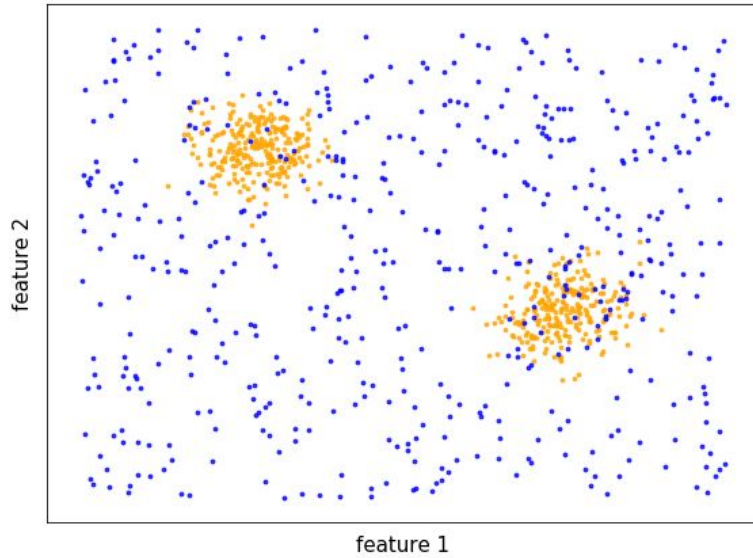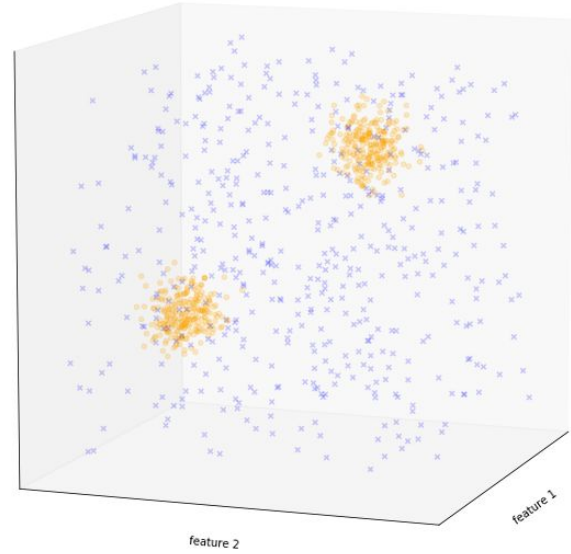- OC have different densities

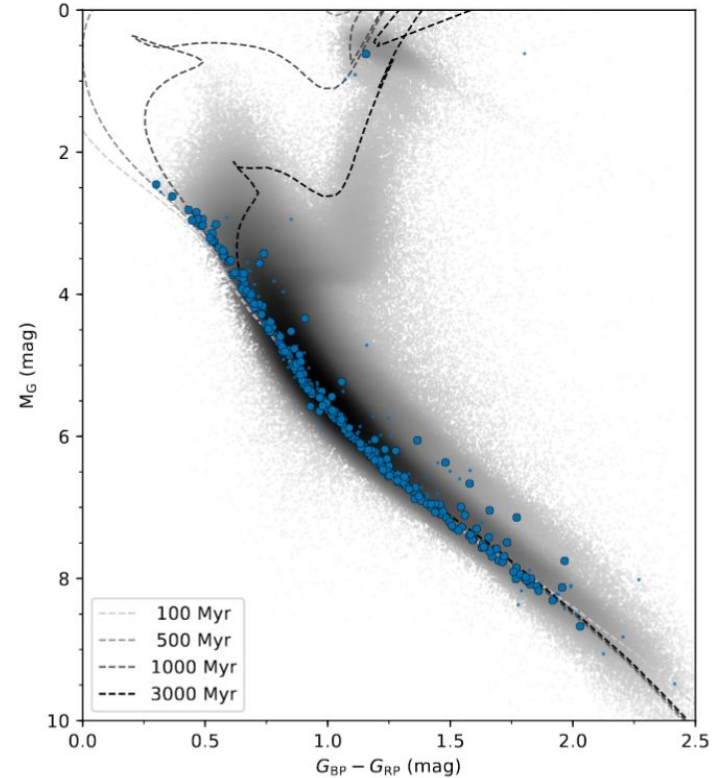# Feature engineering

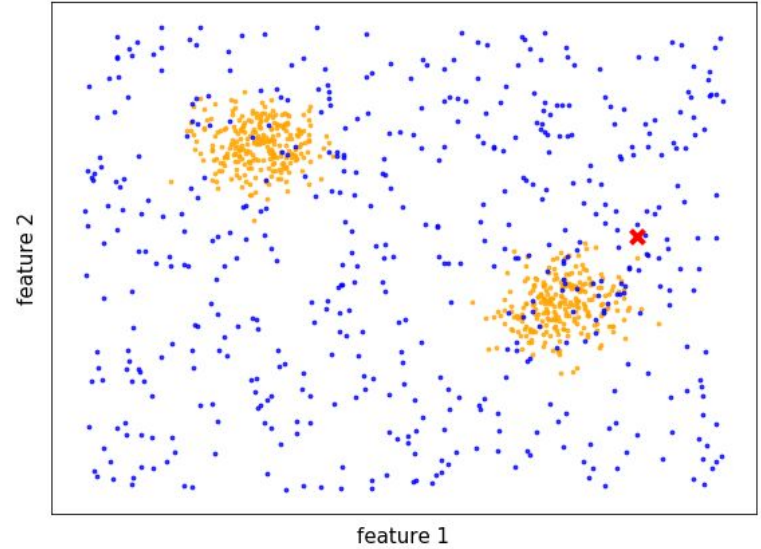# Feature *adding*

# Feature *adding*

# Age engineering

- Fitting a curve to the data which you
  believe are members of your cluster
- Usually quite messy, isochrone models
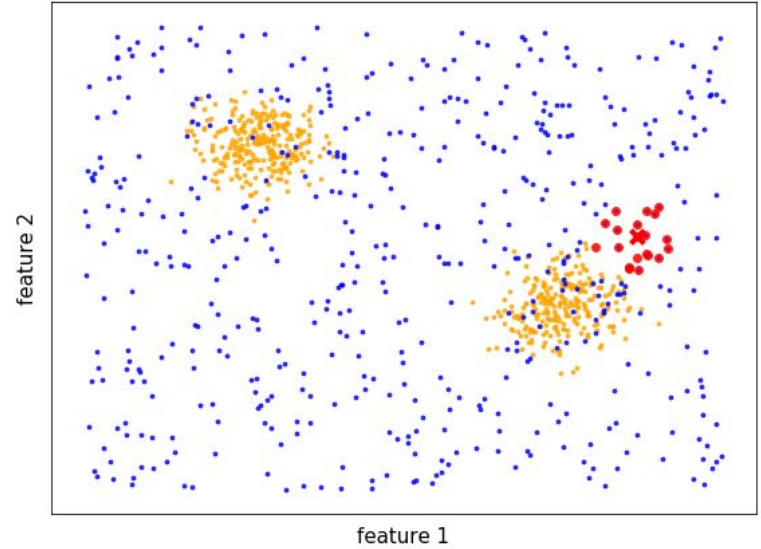  are not perfect



Source: Meingast et al. (2019)

# Age engineering
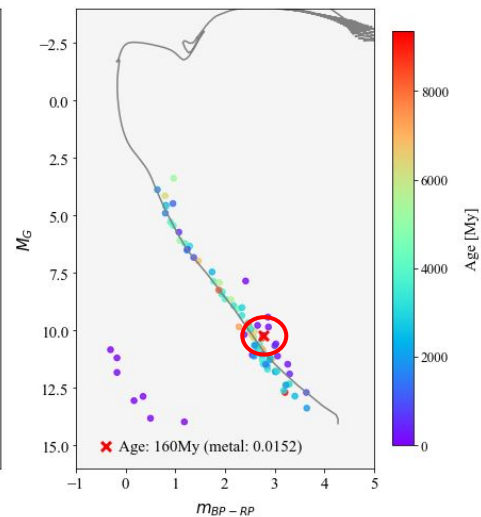
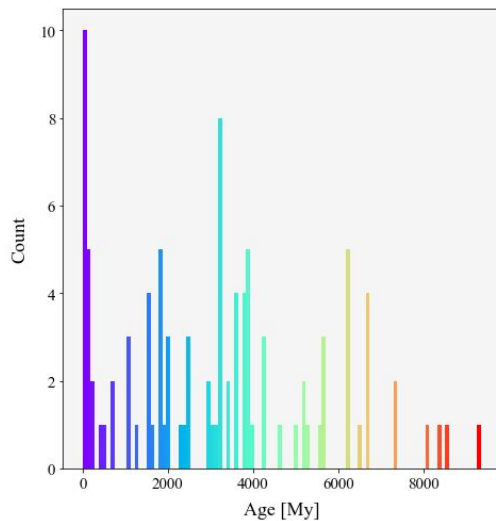1. Take a point from the sample

# Age engineering

1. Take a point from the sample
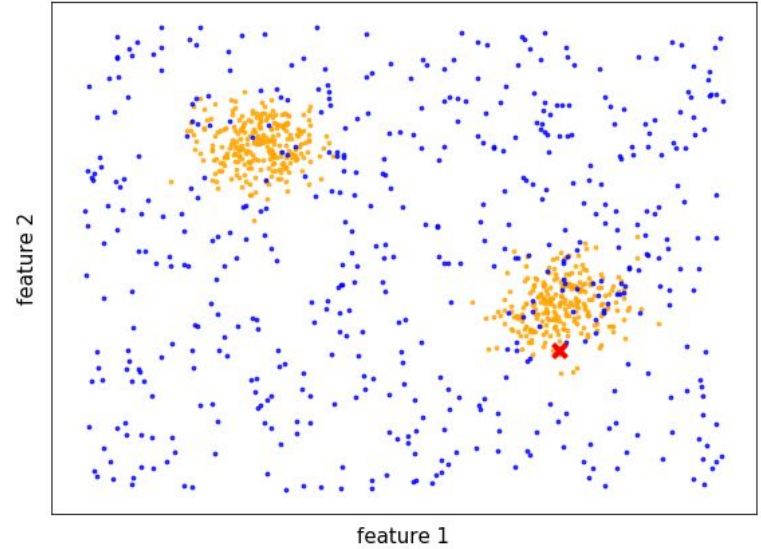
2. Get its neighborhood

# Age engineering

1. Take a point from the sample

2. Get its neighborhood

3. Plot neighborhood points in CMD & fit age to points

# Age engineering

1. Take a point from the sample

# Age engineering

1. Take a point from the sample

2. Get its neighborhood

# Age engineering

1. Take a point from the sample

2. Get its neighborhood

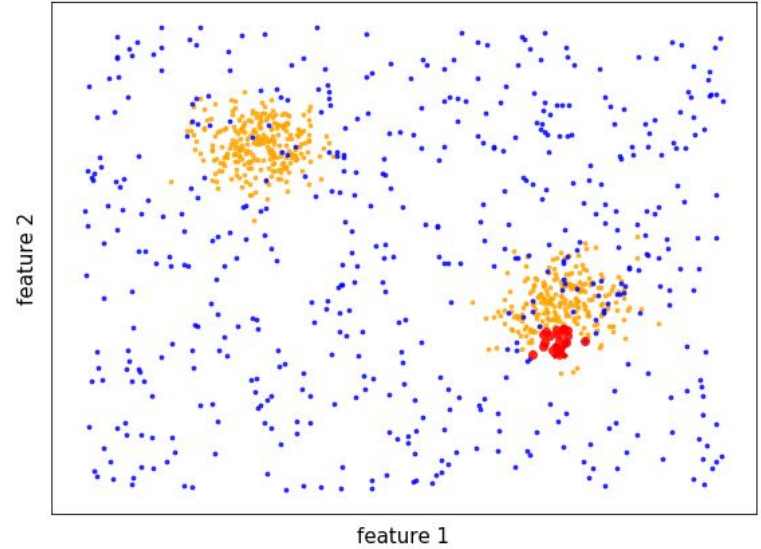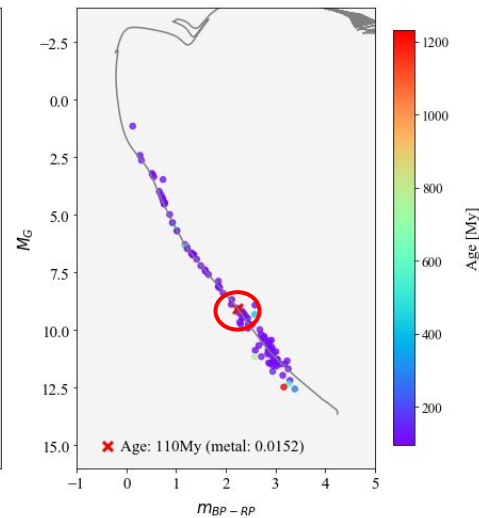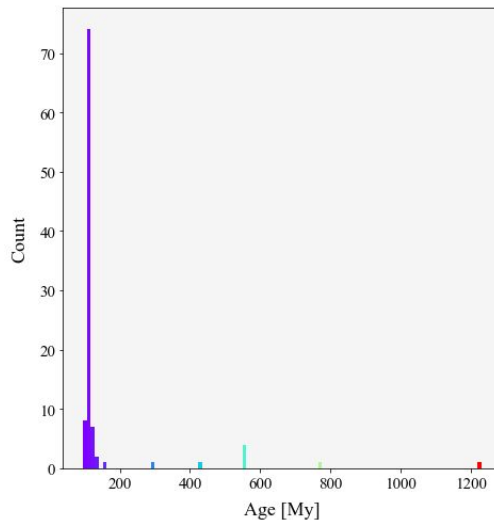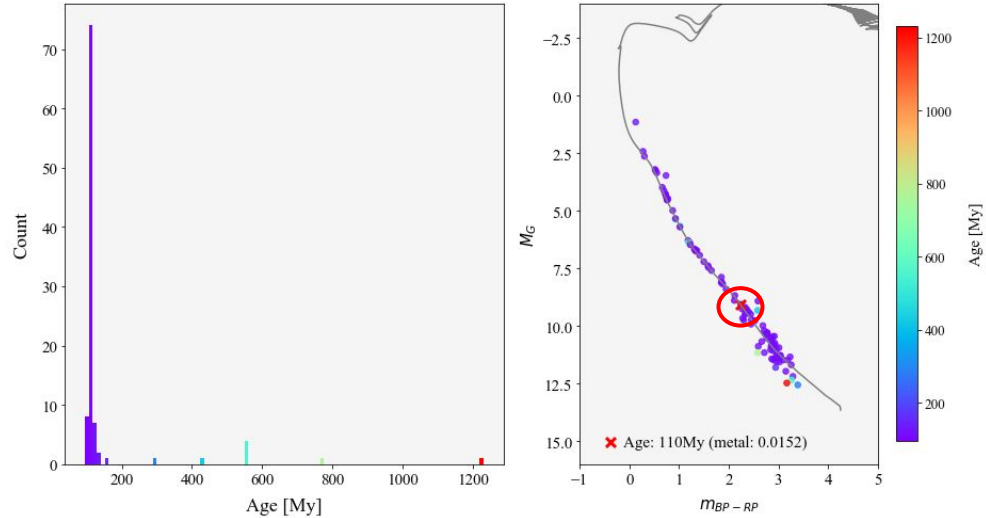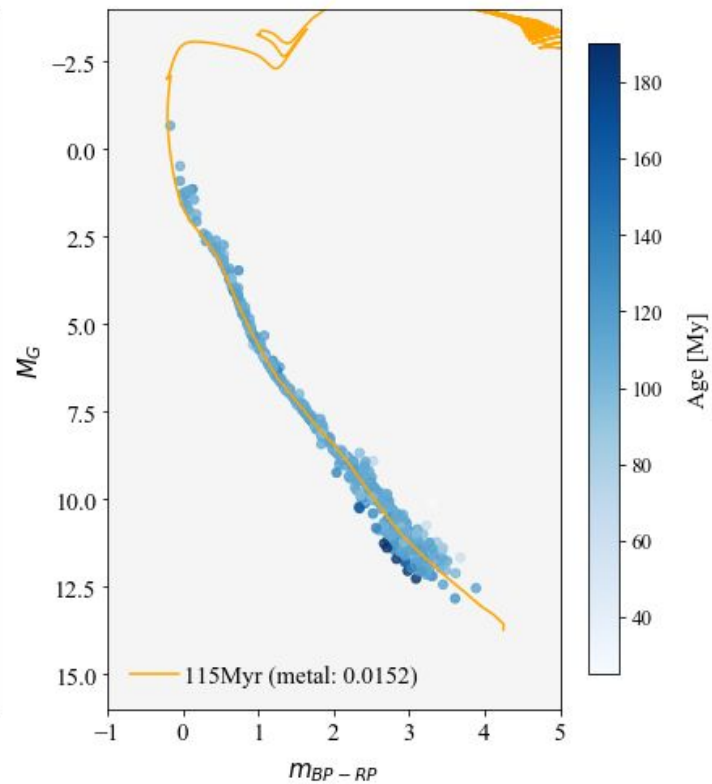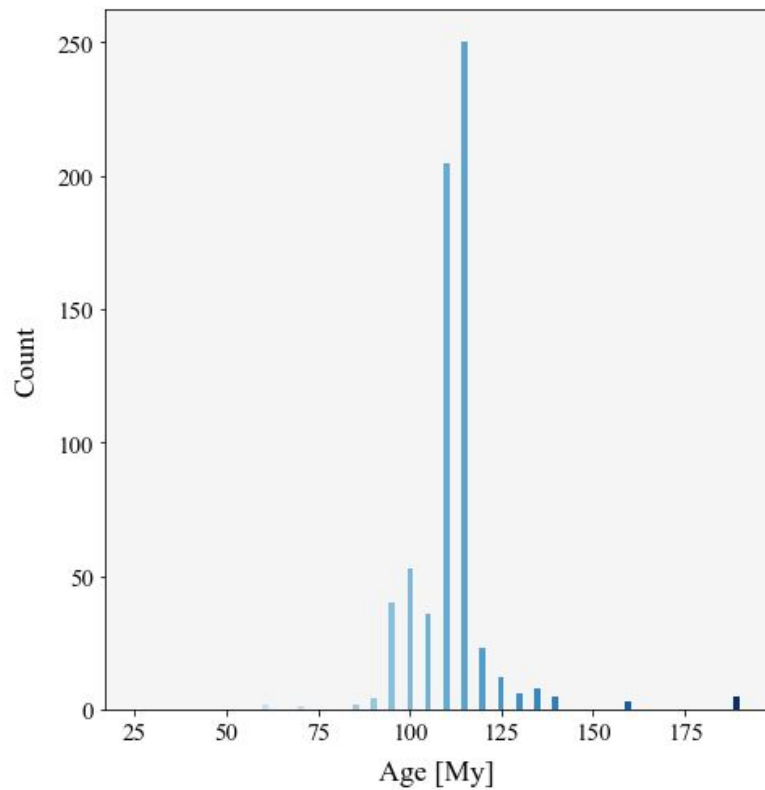3. Plot neighborhood points in CMD & fit age to points

# Age engineering

1. Take a point from the sample

2. Get its neighborhood

3. Plot neighborhood points in CMD & fit age to points



$$d_{ij} = c_x \times \sqrt{(\vec{x}_i - \vec{x}_j)^2} + c_v \times \sqrt{(\vec{v}_i - \vec{v}_j)^2} + c_{age} \times |age_i - age_j| + c_m \times |m_i - m_j|$$

# Results

# Results