# Variable Star Classification

Pierre Dubath

Observatory of the University of Geneva
Gaia CU7 team
Pierre.Dubath@unige.ch
June 1st, 2011

## Random forest automated supervised classification of *Hipparcos* periodic variable stars

P. Dubath,[1,2]* L. Rimoldini,[1,2] M. Süveges,[1,2] J. Blomme,[3] M. López,[4] L. M. Sarro,[5] J. De Ridder,[3] J. Cuypers,[6] L. Guy,[1,2] I. Lecoeur,[1,2] K. Nienartowicz,[1,2] A. Jan,[1,2] M. Beck,[1,2] N. Mowlavi,[1,2] P. De Cat,[6] T. Lebzelter[7] and L. Eyer[1,2]

# Context

- Work in progress of the CU7 team to get ready for the Gaia data analysis

- Using Hipparcos data as a control sample

    ➔ Establish the best classification strategy

- Current choice: classification in three steps

# Variable Star Classification

- Work in progress of the CU7 team to get ready for the Gaia data analysis (major contributions from Isabelle Lecoeur et Lorenzo Rimoldini)

- Random forest classification of Hipparcos periodic variables, Dubath et al. 2011, MNRAS

- Hipparcos unsolved variable classification, Rimoldini et al. in preparation

- Overall performance of a complete Hipparcos variable star classification

Sample of Surveyed Stars

Constant Stars

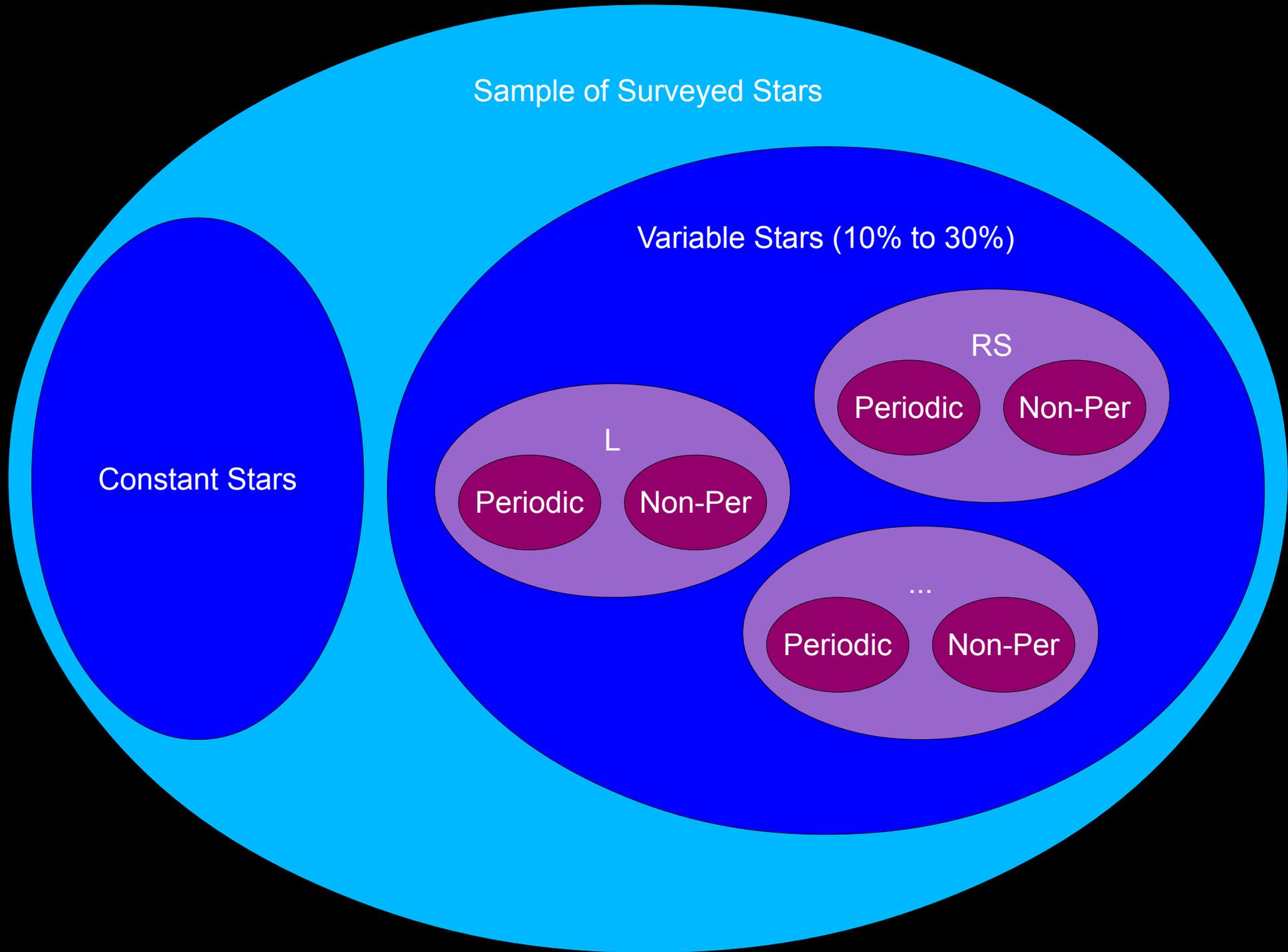Variable Stars (10% to 30%)

L
Periodic
Non-Per

RS
Periodic
Non-Per
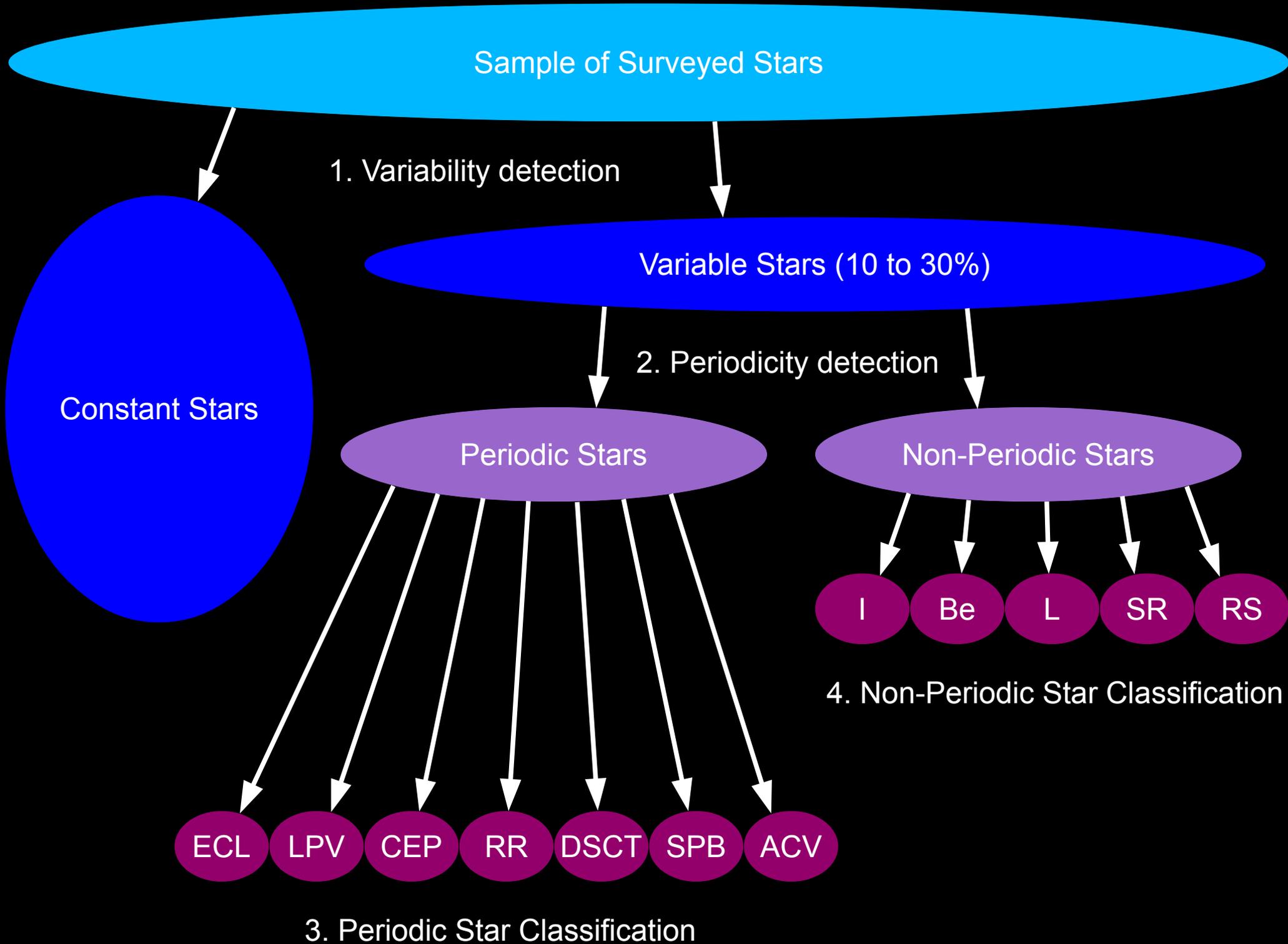
...
Periodic
Non-Per

Sample of Surveyed Stars

1. Variability detection

Constant Stars

Variable Stars (10 to 30%)

2. Periodicity detection

Periodic Stars

Non-Periodic Stars

LPV

ECL

SPB  ACV

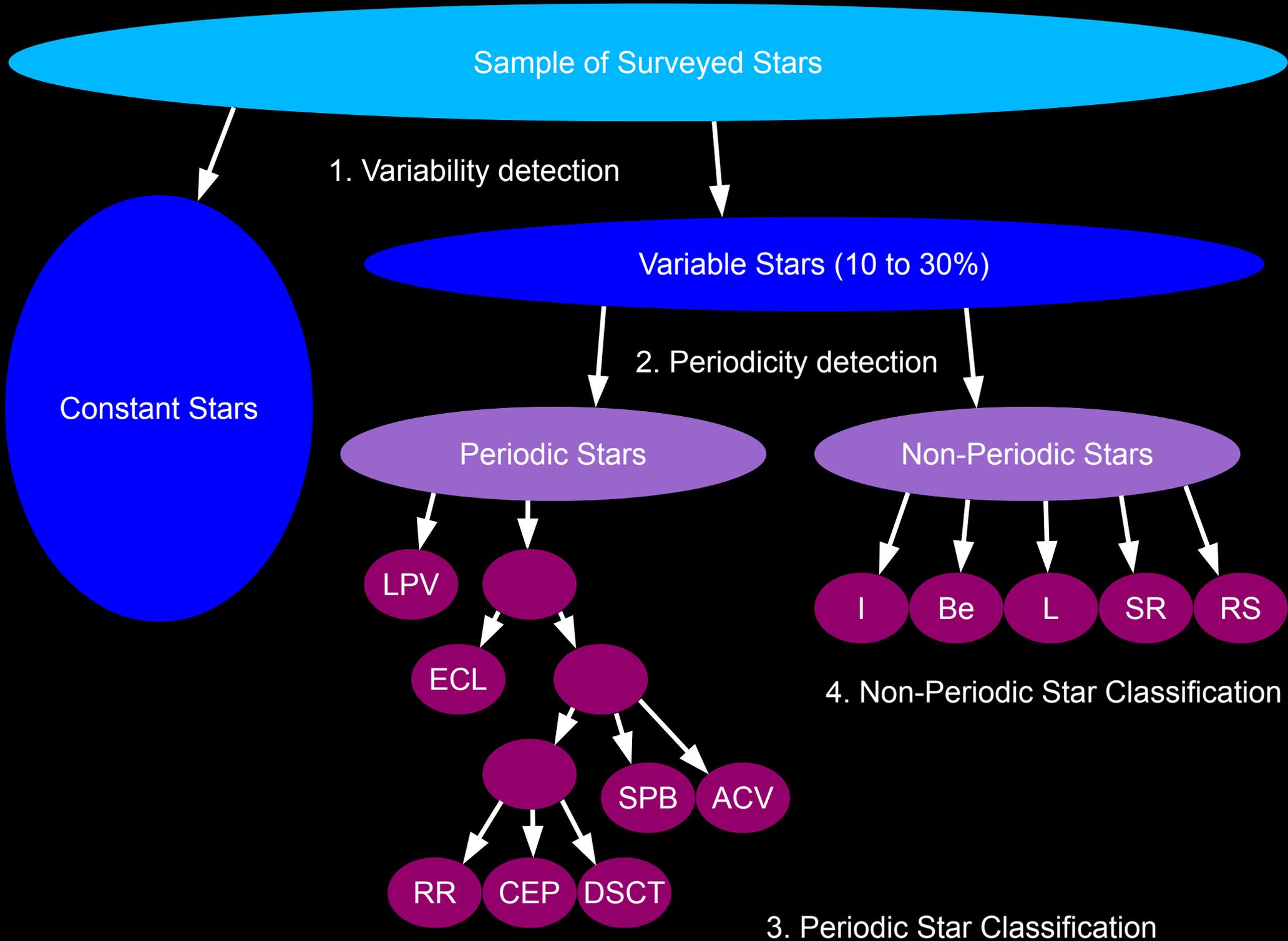RR  CEP  DSCT
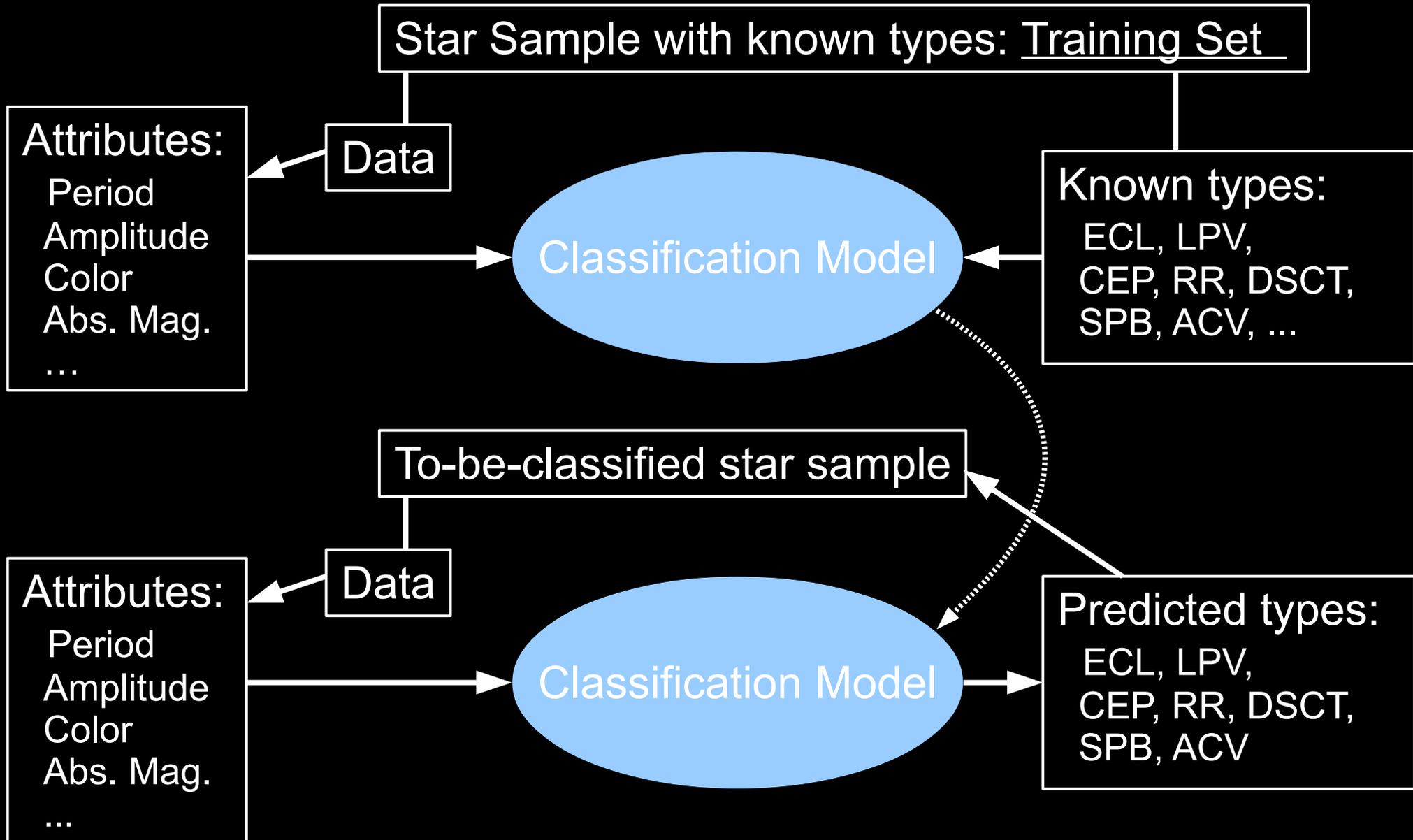
I  Be  L  SR  RS

4. Non-Periodic Star Classification

3. Periodic Star Classification

# Supervised Classification

# Period-Amplitude diagram

# Supervised Classification

Star Sample with known types: <u>Training Set</u>

Attributes:
 Period
 Amplitude
 Color
 Abs. Mag.
 …

Data

Classification Model

Known types:
 ECL, LPV,
 CEP, RR, DSCT,
 SPB, ACV, ...

To-be-classified star sample

Attributes:
 Period
 Amplitude
 Color
 Abs. Mag.
 …

Data

Classification Model

Predicted types:
 ECL, LPV,
 CEP, RR, DSCT,
 SPB, ACV

**Table 1.** Training set composition

| Type | | Num | Main reference |
|---|---|---|---|
| Eclipsing Binary | *EA* | 228 | Hipparcos |
| | *EB* | 255 | Hipparcos |
| | *EW* | 107 | Hipparcos |
| Ellipsoidal | *ELL* | 27 | Hipparcos |
| Long Period Variable | *LPV* | 285 | Lebzelter (p. c.) |
| RV Tauri | *RV* | 5 | AAVSO |
| W Virginis | *CWA* | 9 | AAVSO |
| | *CWB* | 6 | AAVSO |
| Delta Cepheid | *DCEP* | 189 | AAVSO |
|   (first overtone) | *DCEPS* | 31 | AAVSO |
|   (multi mode) | *CEP(B)* | 11 | AAVSO |
| RR Lyrae | *RRAB* | 72 | AAVSO |
| | *RRC* | 20 | AAVSO |
| Gamma Doradus | *GDOR* | 27 | De Cat (p. c.) |
| Delta Scuti | *DSCT* | 43 | AAVSO |
|   (low amplitude) | *DSCTC* | 81 | AAVSO |
| SX Phoenicis | *SXPHE* | 4 | AAVSO |
| Beta Cephei | *BCEP* | 30 | De Cat (p. c.) |
| Slowly Pulsating B star | *SPB* | 81 | De Cat (p. c.) |
| B emmission line star | *BE* | 9 | AAVSO |
| Gamma Cassiopeiae | *GCAS* | 4 | AAVSO |
| Alpha Cygni | *ACYG* | 18 | AAVSO |
| Alpha-2 Canum Venaticorum | *ACV* | 77 | Romanyuk (p. c.) |
| SX Arietis | *SXARI* | 7 | Romanyuk (p. c.) |
| BY Draconis | *BY* | 5 | Eker et al. (2008) |
| RS Canum Venaticorum | *RS* | 30 | Eker et al. (2008) |
| | **Total:** | 1661 | |

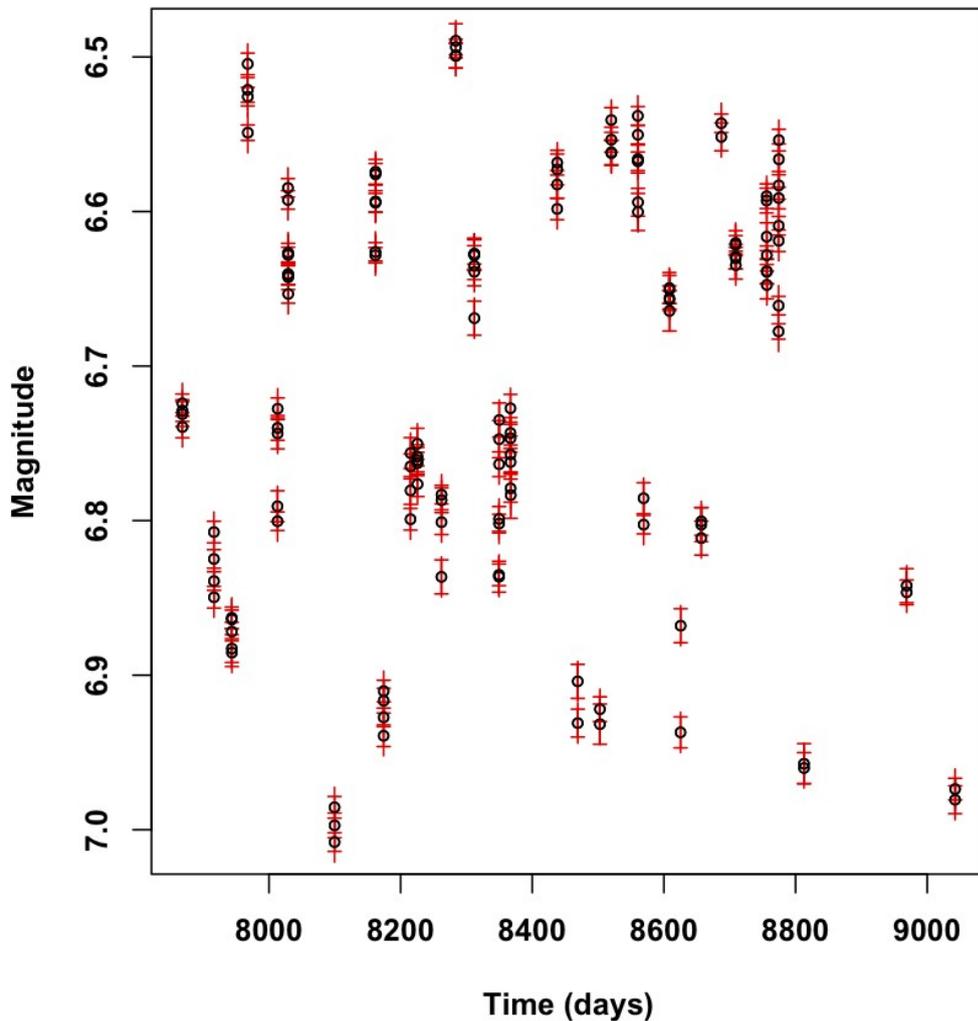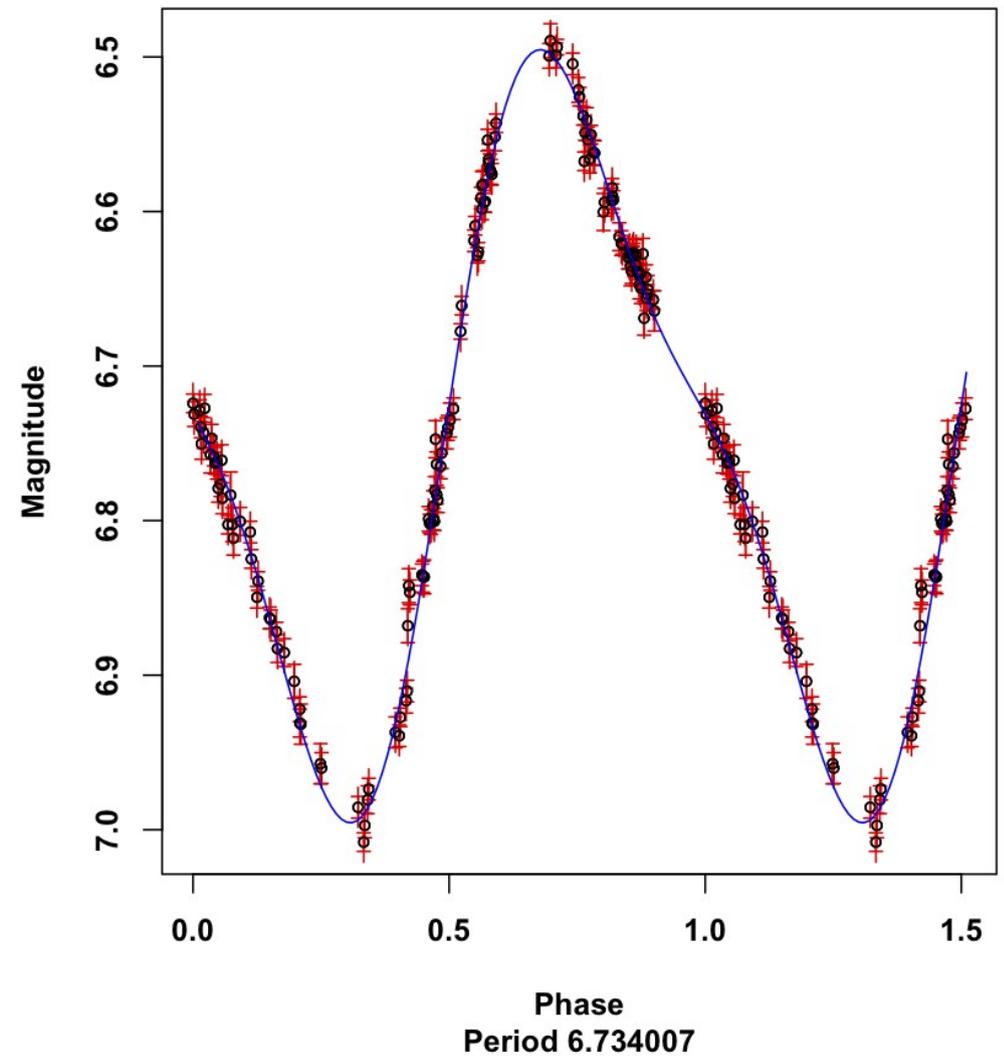AAVSO : Watson, Henden, & Price (2010)

# Lomb-Scargle period search



- (Zechmeister & Kűrster 2009)

# Fourier series modeling



**Original Time Series**
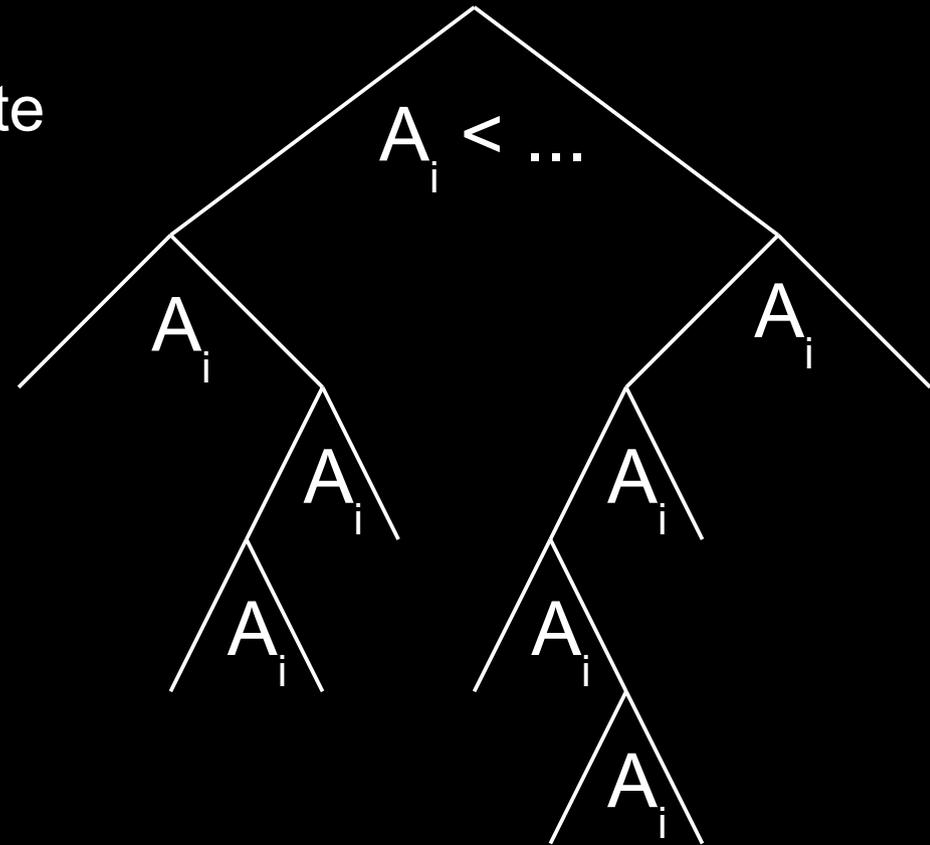
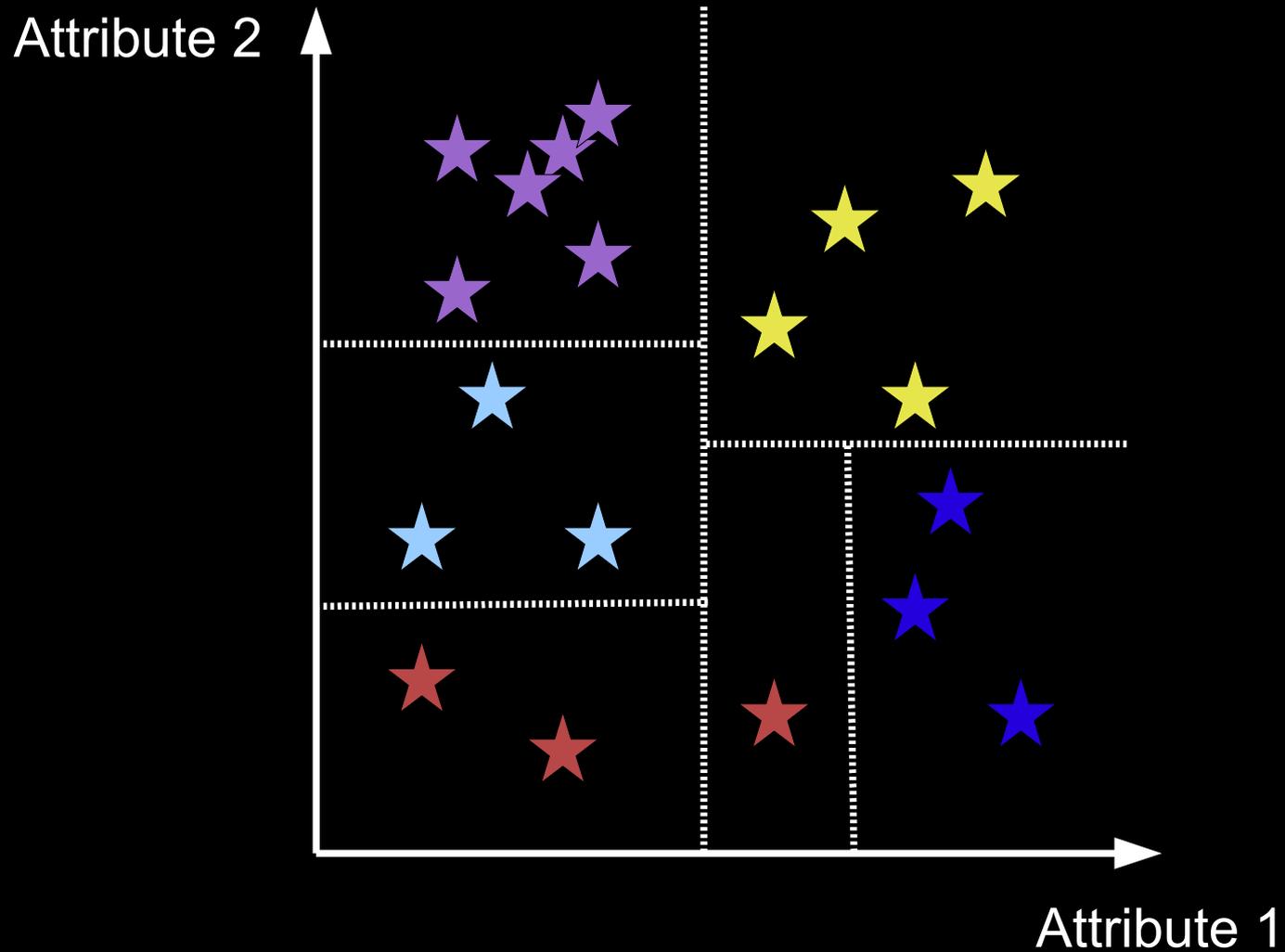**Srv: hipparcos, Src: 60259**

Phase
Period 6.734007

# Supervised Classification

# Random Forest (1/3)

- Classification trees

- Binary partitions using one attribute

- Each split minimize impurity

$A_i < \ldots$

$A_i$

$A_i$

$A_i$

$A_i$

$A_i$

$A_i$

$A_i$

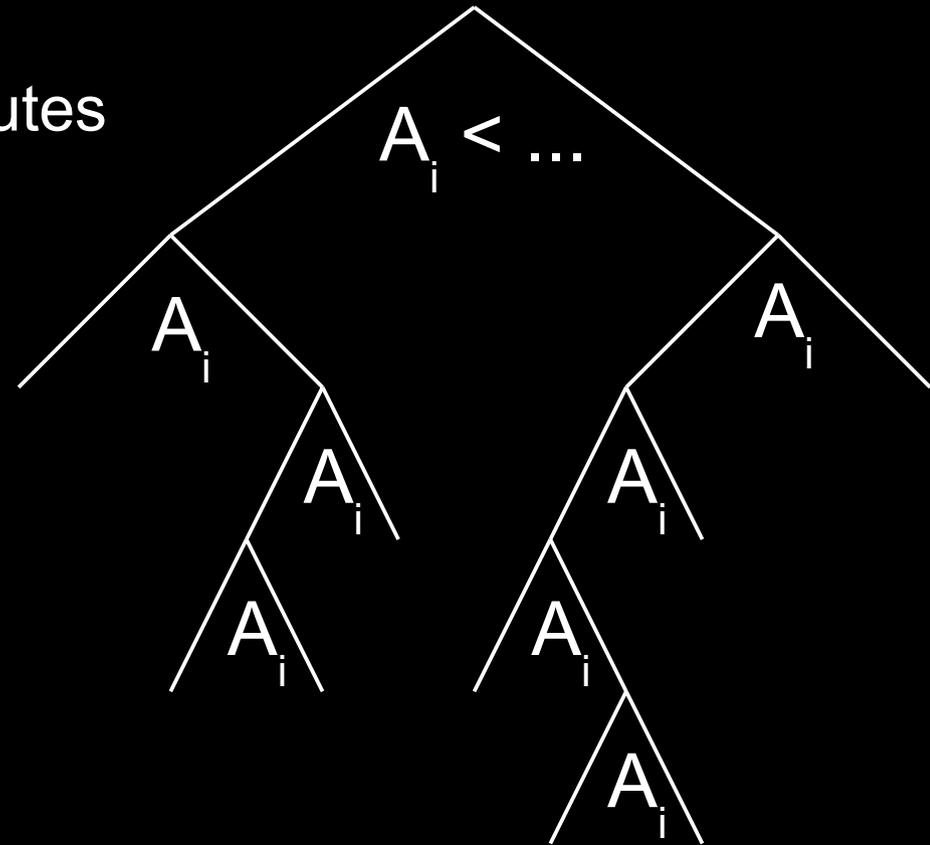# Random Forest (2/3)

# Random Forest (3/3)

- Bootstrapping

- $A_i$ from a random sub-set of attributes

- Average many trees

$A_i < \dots$

$A_i$

$A_i$

$A_i$

$A_i$

$A_i$

$A_i$

$A_i$

# Multistage classification

# Random Forest Attribute Importance



- Attributes with Spearman correlation larger than 80% are trimmed

# Cross-Validation

Training Set sub-set

Attributes:
 Period
 Amplitude
 Color
 Abs. Mag.
 …

Data

Classification Model

Predicted types

Known types

# Random Forest CV error rates

Confusion matrix (rows labelled at right, columns labelled at top):

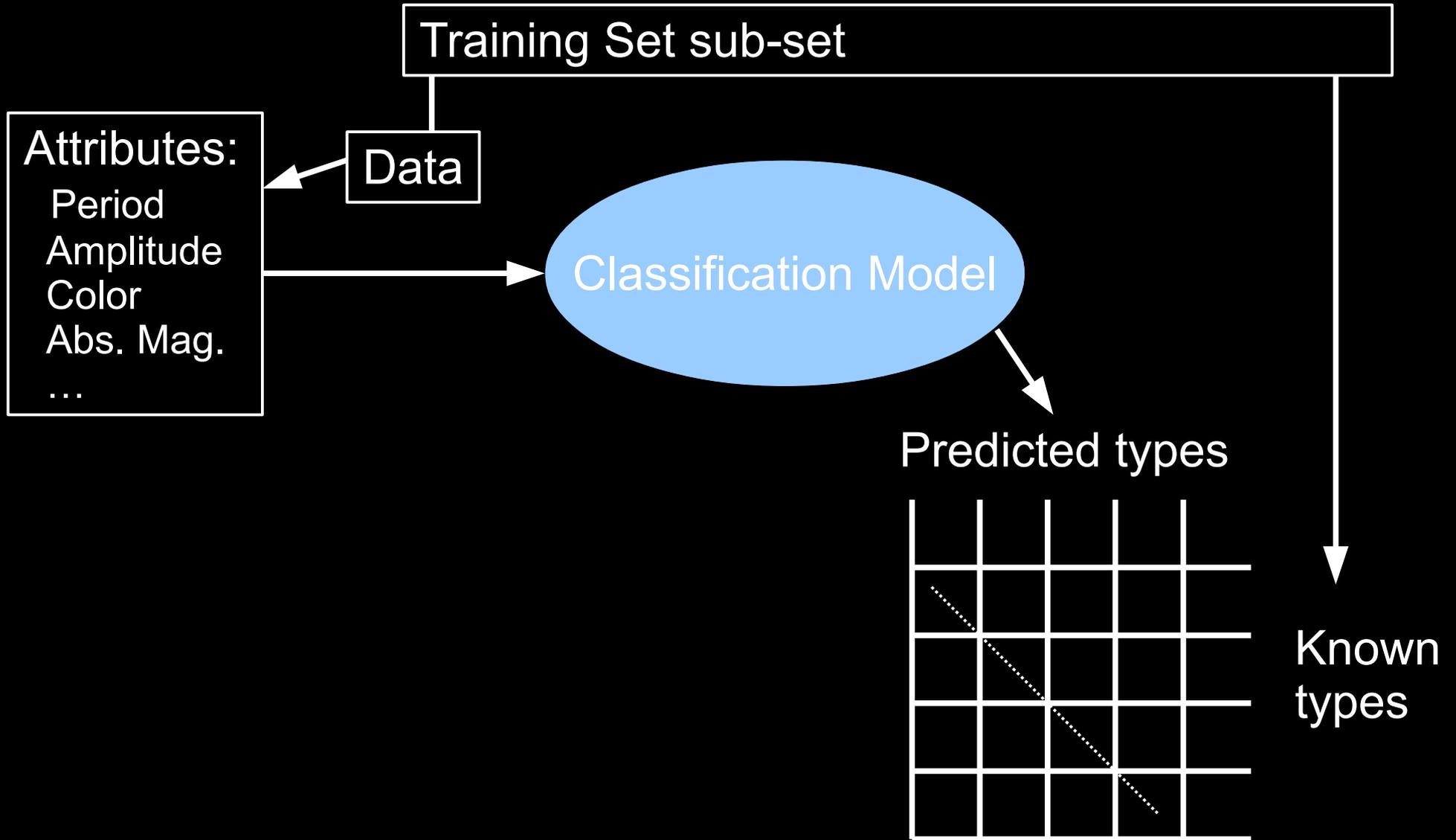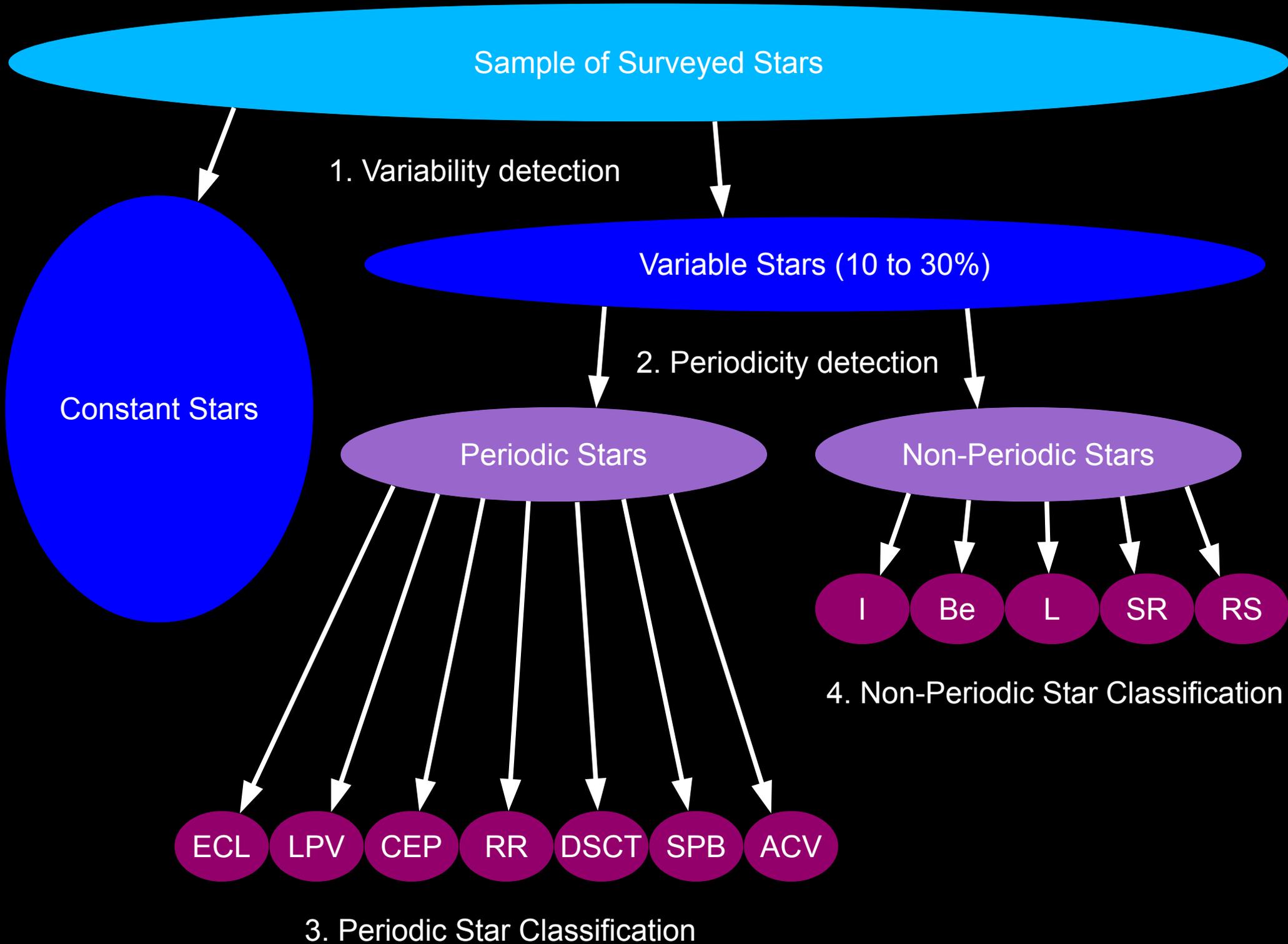| | EA | EB | EW | ELL | LPV | RV | CWA | CWB | DCEP | DCEPS | CEP(B) | RRAB | RRC | GDOR | DSCT | DSCTC | BCEP | SPB | BE+GCAS | ACYG | ACV | SXARI | BY+RS |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| EA | 214 | 13 | | | | | | | | | 1 | | | | | | | | | | | | |
| EB | 19 | 191 | 28 | 2 | 1 | | | | 2 | | | | | 1 | | 4 | | 3 | | 2 | 2 | | |
| EW | | 30 | 76 | | | | | | | 1 | | | | | | | | | | | | | |
| ELL | | 14 | | | 1 | | | | | | | | | 1 | 1 | 3 | | | | | 5 | | 2 |
| LPV | | | | | 285 | | | | | | | | | | | | | | | | | | |
| RV | | 1 | | | 1 | | | | 2 | 1 | | | | | | | | | | | | | |
| CWA | | 2 | | | | 1 | | | 5 | | | | | | | | | | | | | | 1 |
| CWB | | 1 | | | | | | 2 | 2 | 1 | | | | | | | | | | | | | |
| DCEP | | | | | | | | | 183 | 5 | 1 | | | | | | | | | | | | |
| DCEPS | | 1 | | | | | | | 11 | 17 | | | | | | | | | | | | | 2 |
| CEP(B) | | 1 | | | | | | | 4 | | 6 | | | | | | | | | | | | |
| RRAB | | 1 | | | | | | | | | | 69 | 1 | | | | | 1 | | | | | |
| RRC | | 2 | 4 | | | | | | | | | 1 | 12 | | 1 | | | | | | | | |
| GDOR | | | | | | | | | | | | | | 27 | | | | | | | | | |
| DSCT | | 1 | 1 | | | | | | | | | | | 1 | 32 | 12 | | | | | | | |
| DSCTC | | 1 | | | | | | | | | | | | | 1 | 77 | | | | | 2 | | |
| BCEP | | 1 | 1 | | | | | | | | | | | | | 1 | 26 | 1 | | | | | |
| SPB | | | 1 | | | | | | | | | | | | | | 1 | 74 | | 1 | 4 | | |
| BE+GCAS | 1 | | | | | | | | | 1 | | | | | | | | 5 | | 2 | 4 | | |
| ACYG | | 1 | | | | | | | | | | | | | | | | | 1 | 13 | 2 | | 1 |
| ACV | | 3 | | | | | | | | 1 | | | | 1 | | | | 6 | | | 66 | | |
| SXARI | | 2 | | | | | | | | | | | | | | | | 2 | | | 3 | | |
| BY+RS | | 1 | | | | | | | 1 | | | | | | | | | | | | | | 33 |

# Variability detection

- ## Variability criteria

$$\chi^2 = \sum_{i=1}^{N} \left( \frac{x_i - \bar{x}}{\sigma^2} \right)^2$$

| Using errors | Not using errors |
|---|---|
| Chi square | Abbe |
| Skewness | Skewness |
| Kurtosis | Kurtosis |
| Stetson | Inter-quartile range |
| Outlier median | B/R Correlation |

$$Abbe = \frac{1}{2} \frac{\sum_{i=1}^{n-1} \left( x_{i+1} - x_i \right)^2}{\sum_{i=1}^{n} \left( x_i - \bar{x} \right)^2}$$

- ## Compute pValues: probability of the null hypothesis $H_0$ = constant star
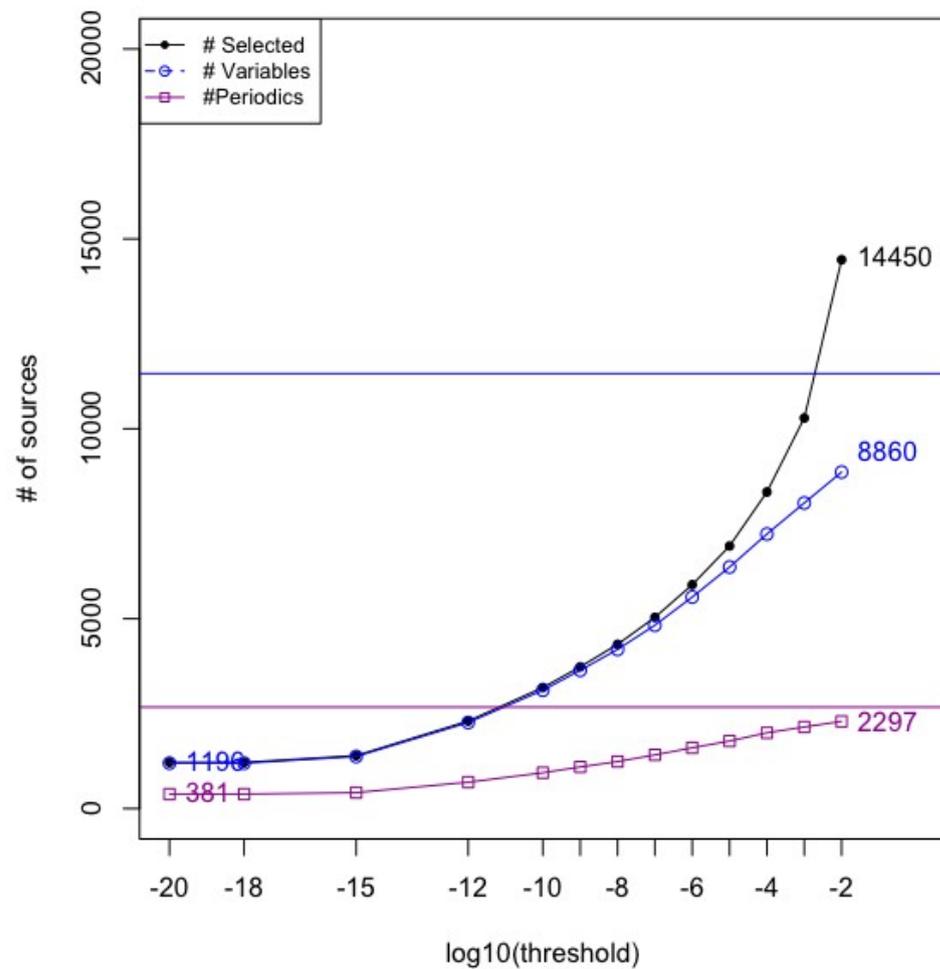
  → If pValue < 1e-4 → variable star

# Inter-quartile range Hypparcos pValues



**Histogram pValues InterquartileRange**

white = all, red = constants, blue = variables

## Chi2 Test

- # Selected
- # Variables
- #Periodics

11387
8587
2564
2667

## Abbe Test

- # Selected
- # Variables
- #Periodics

14450
8860
2297
1196
381

log10(threshold)

# of sources

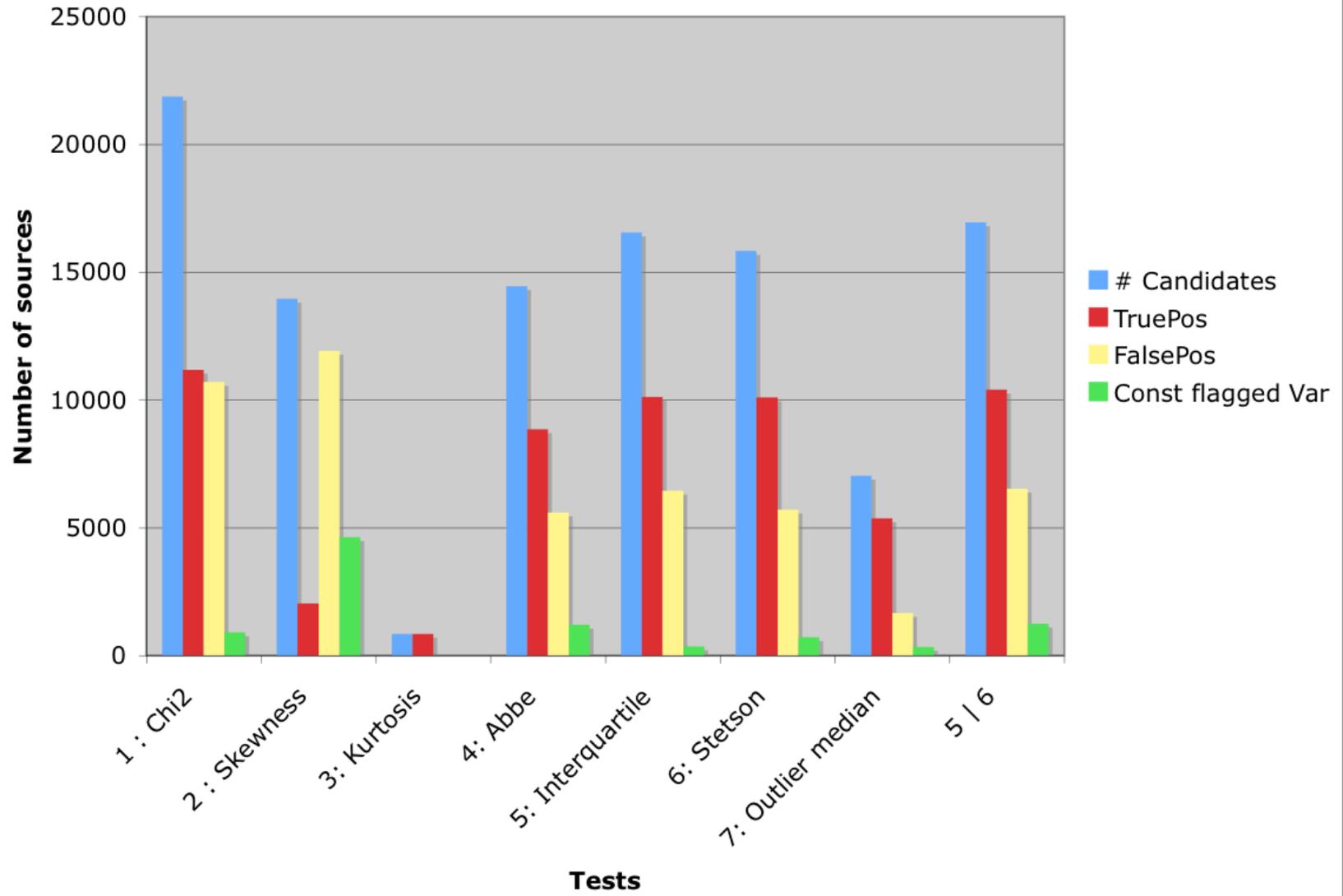**Kurtosis Test**

**Outliermedian Test**

**Comparison of tests to select variables**

Legend: # Candidates, TruePos, FalsePos, Const flagged Var

X-axis (Tests): 1 : Chi2, 2 : Skewness, 3: Kurtosis, 4: Abbe, 5: Interquartile, 6: Stetson, 7: Outlier median, 5 | 6

Y-axis: Number of sources

**Union (Abbe/Stet/IQR/Chi2/Kur/Out)Test**

Legend:
- # Selected
- # Variables
- #Periodics

x-axis: log10(threshold)

y-axis: # of sources
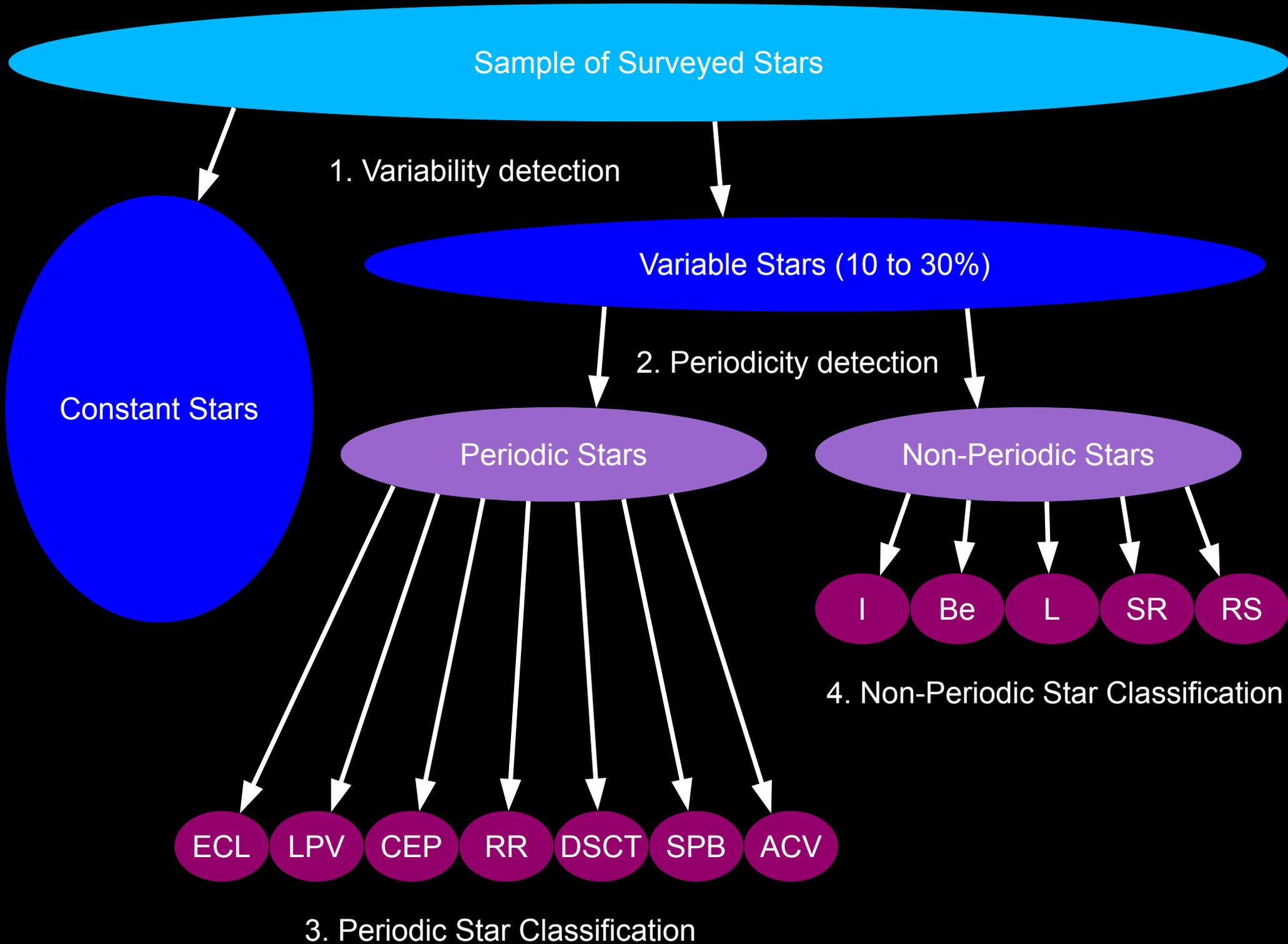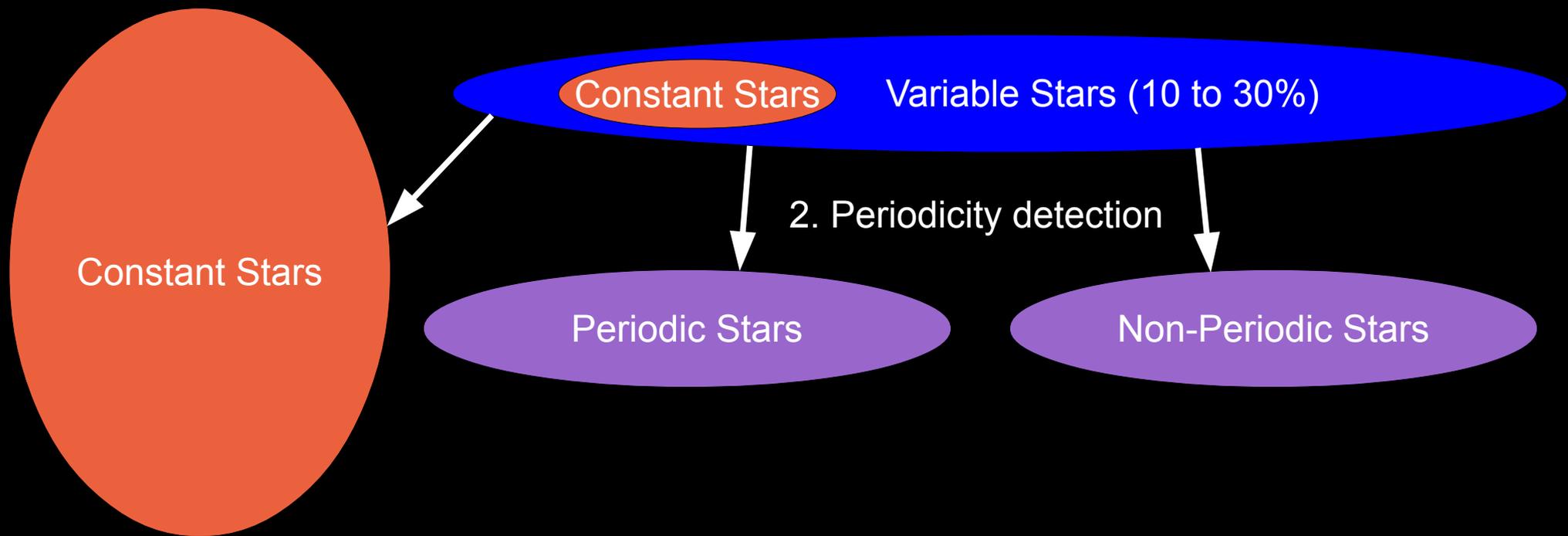
# Choice for variability detection

- Union of Stetson with pValues < 1e-2 and Inter-quartile range with pValues < 1e-3
  - → 17'006 candidates (14.8 % of total)

Sample of Surveyed Stars

1. Variability detection

Constant Stars

Variable Stars (10 to 30%)

2. Periodicity detection

Periodic Stars

Non-Periodic Stars

I   Be   L   SR   RS

4. Non-Periodic Star Classification

ECL   LPV   CEP   RR   DSCT   SPB   ACV

3. Periodic Star Classification

Constant Stars

Variable Stars (10 to 30%)

Constant Stars

2. Periodicity detection
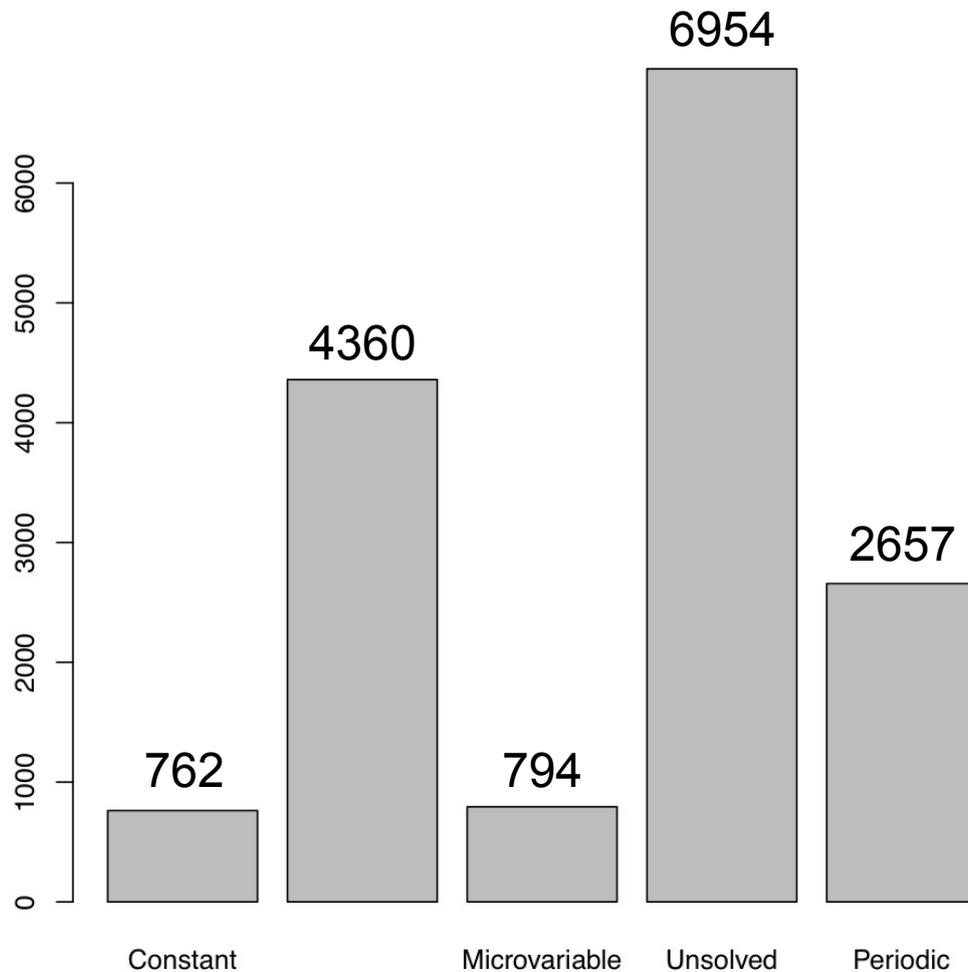
Periodic Stars

Non-Periodic Stars

# Periodicity detection through a supervised classification

- Compute a number of "attributes" characterizing objects and their light-curve

- Use attributes (features) as predicators (variables) in a supervised classifiers

- Train the classifier with a set of stars of known types

- Use a 10-fold Cross-Validation (CV) to evaluate the performance of this approach

# Period search

- Generalized Lomb–Scargle method (Zechmeister & Kűrster 2009)

- Our Sample of 15'527 stars includes 3022 stars with known periods (3022 = 2657 P + 365 U)

- Recovery rate of 77 % (2323 out of 3022)

  ➔ 1644 with correct period

  ➔ 679 with twice the period values

- Recovery rate for the 2657 periodic = 86 % (i.e., 2300)

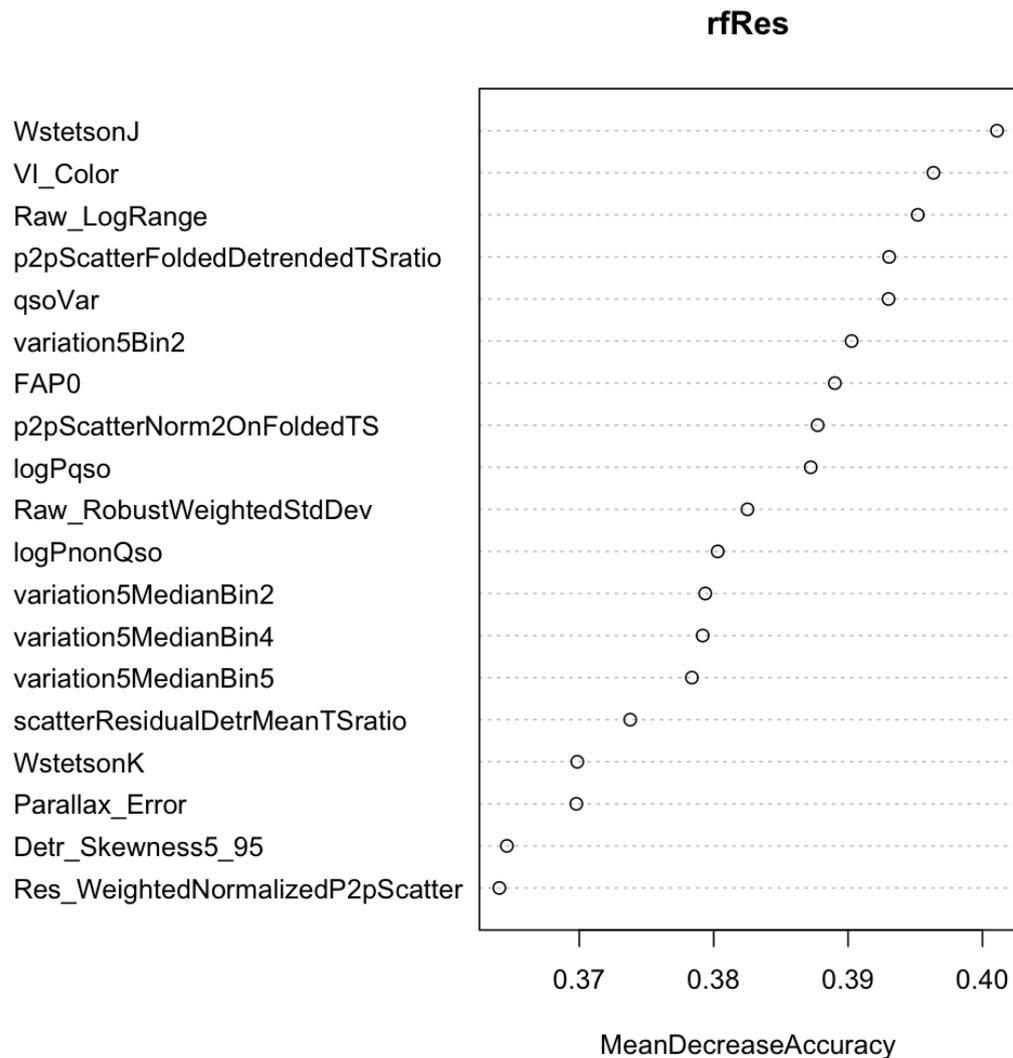# Hipparcos variability types



Total : 15527

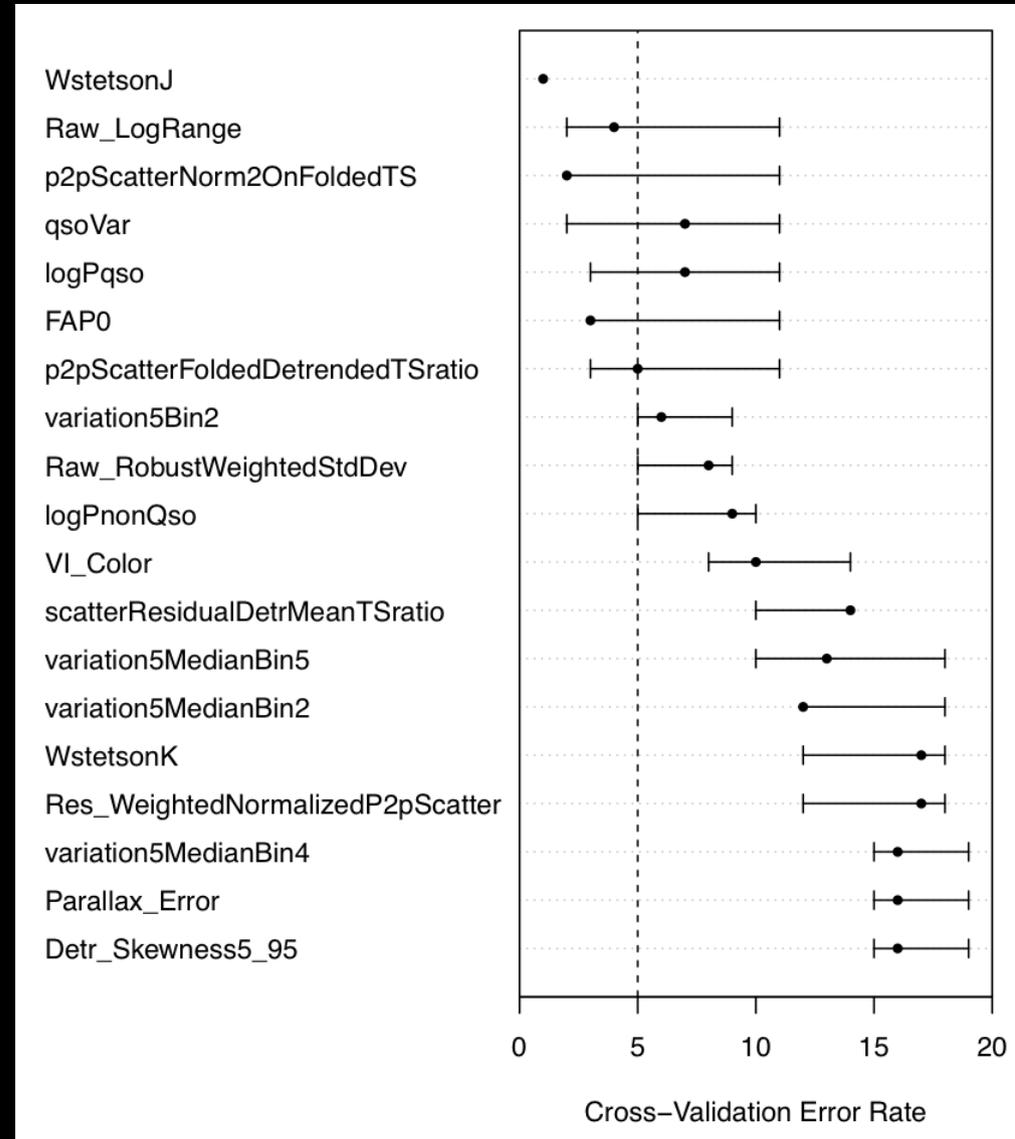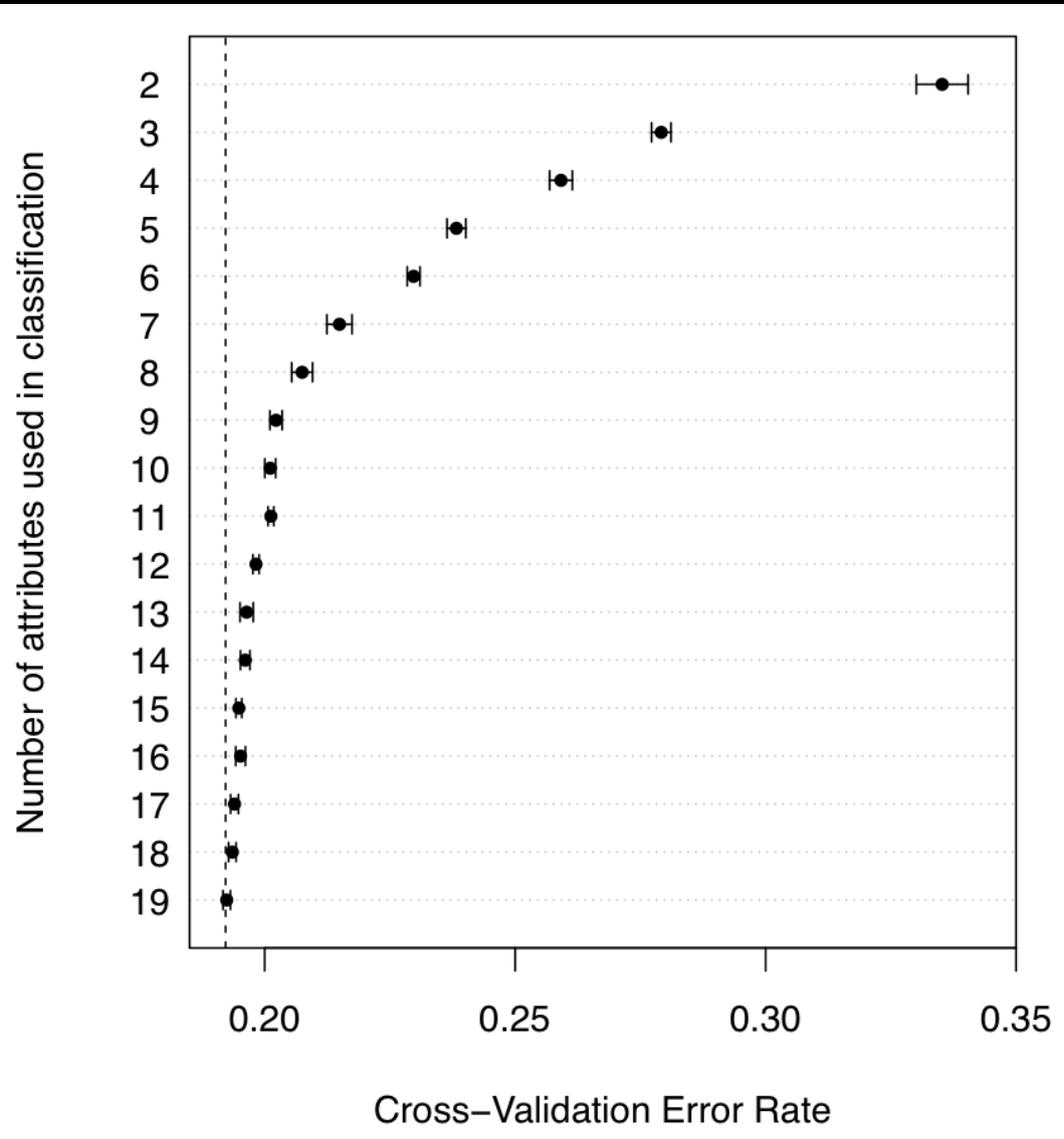P : 2657 != 2712

Hip Types R (662)
         D (816)
   ...removed from TS!

# Random forest attribute importance



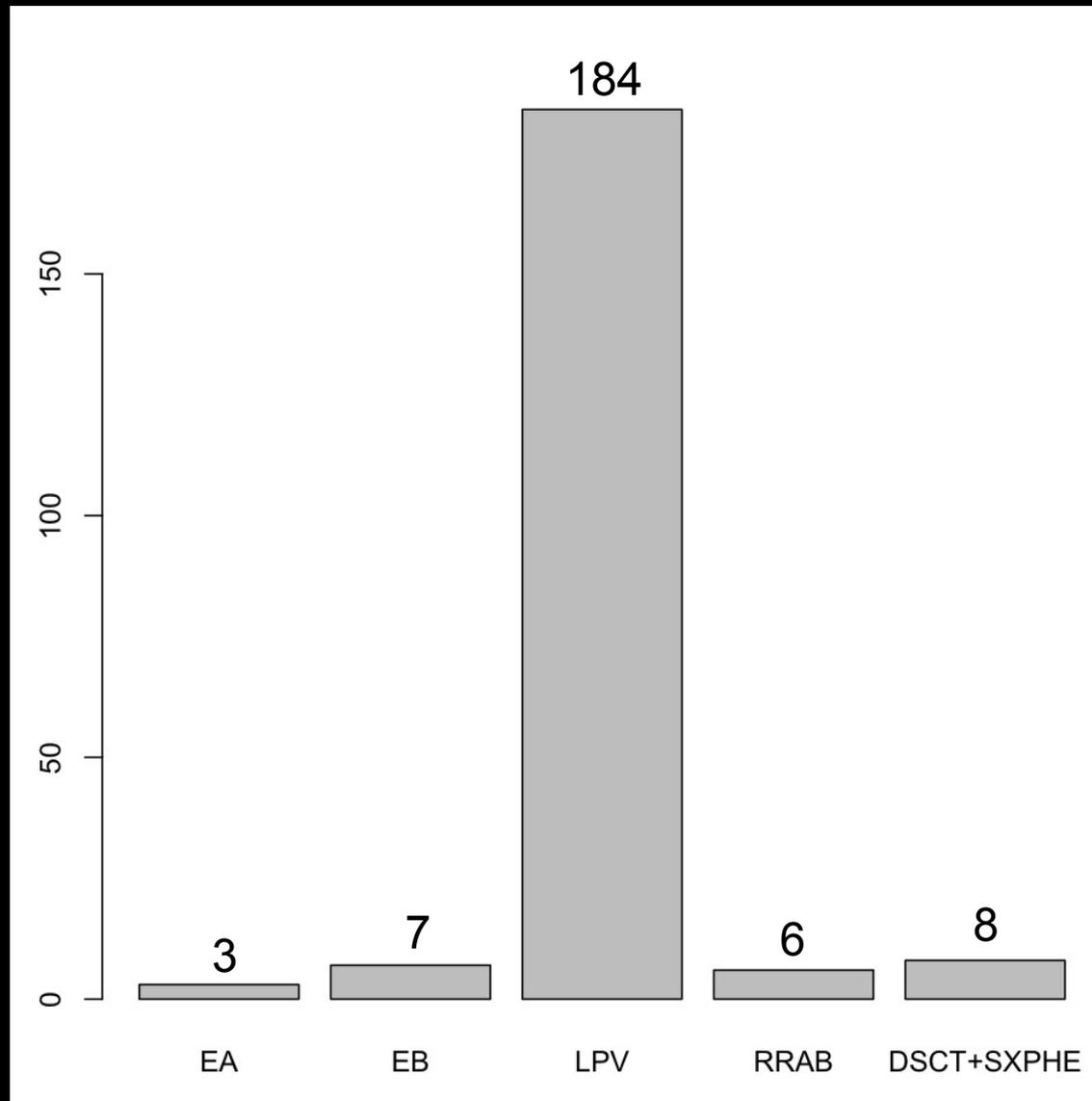- Attributes with Spearman correlation larger than 80% are trimmed

# Random forest CV error rates

# Random forest confusion matrix



|  | Constant | Microvariable | Unsolved | Periodic |  |
|---|---|---|---|---|---|
| 212 | 549 |  | 1 |  | Constant |
| 112 | 3738 | 94 | 405 | 11 |  |
|  | 319 | 313 | 152 | 10 | Microvariable |
|  | 523 | 68 | 6167 | 196 | Unsolved |
|  | 33 | 27 | 503 | 2094 | Periodic |

# Predicted types for non-periodic identified as periodic

# Conclusion

- We established a complete scheme for variable star classification
- Optimized for Hipparcos data
  - ➔ Training set must be representative of the test set
- Hipparcos classification relatively easy: clean sample and well known stars
- Can be completed with additional information
  - ➔ Color light curves
  - ➔ Radial velocity time series
- Next step: apply our scheme to other surveys