

# Data-driven Gene Regulatory Network Inference based on Classification Algorithms

---

Sergio Peignier<sup>a</sup>, Pauline Schmitt, Federica Calevro

*INSA Lyon, INRA*

*BF2I, UMR0203, F-69621*

*Lyon, France*



---

<sup>a</sup> [sergio.peignier@insa-lyon.fr](mailto:sergio.peignier@insa-lyon.fr)

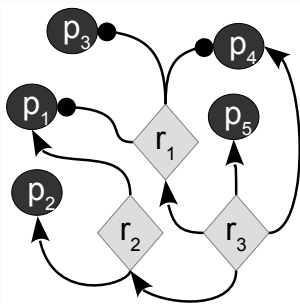
# Gene Regulatory Networks (GRNs)

## Central Dogma of Molecular Biology

*Gene*  $\xrightarrow{\text{Transcription}}$  *mRNA*  $\xrightarrow{\text{Translation}}$  *Protein*

## Definition

Set of interacting molecular regulators (e.g. transcription factors) controlling the creation of gene products (e.g. mRNA, proteins).



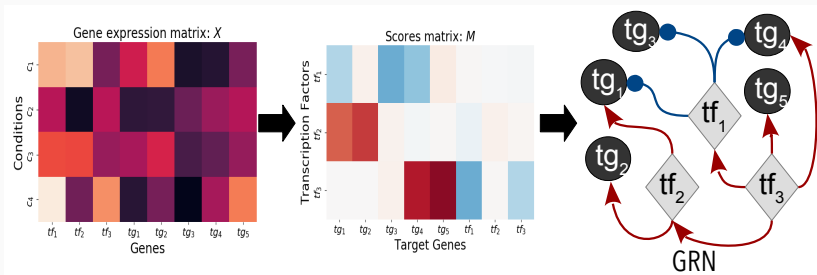
- Wide range of mechanisms:  
(e.g., epigenetic, transcriptional ...)
- Important biological role:
  - Adaptation
  - Differentiation
  - Versatility
  - Morphogenesis ...

# Data-driven GRN Inference

## General principle

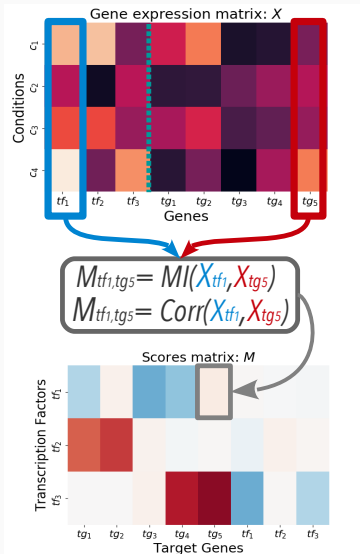
- Based on high-throughput gene expression data.
- Score possible links between:
  - Regulators, i.e. Transcription Factors (TFs)
  - Target Genes (TGs)
- Select most promising links.

Well-known paradigm: Simple, accurate, computationally efficient

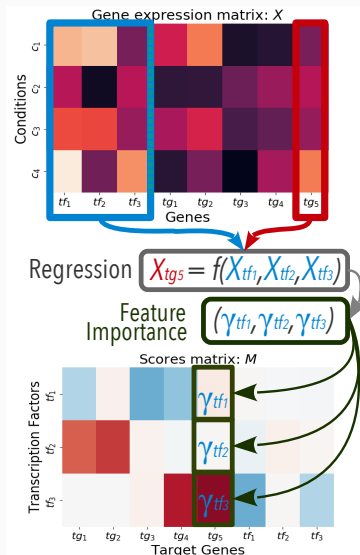


# Data-driven Inference Families

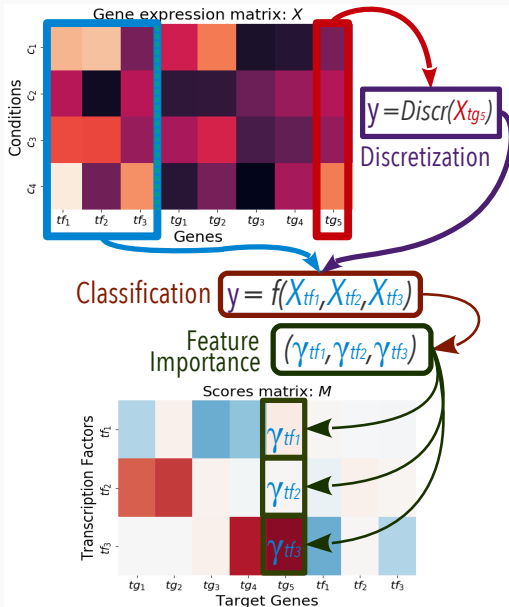
## Correlation | Mutual Information



## Regression Methods

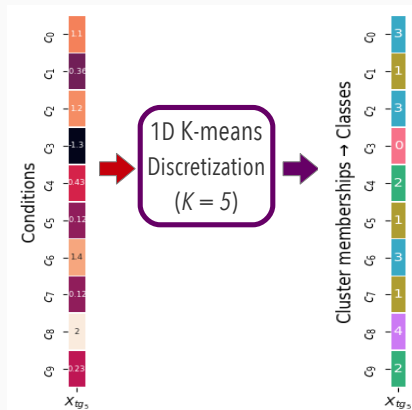


# Classification-based GRN Inference

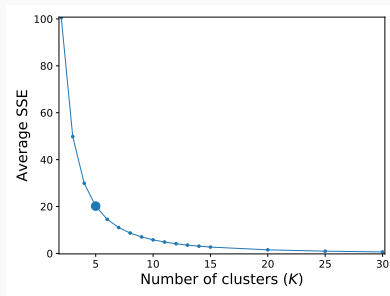


# Target gene expression discretization

- $K$ -means  $\rightarrow$  Discretize TG exp.
- Cluster membership  $\rightarrow$  Class

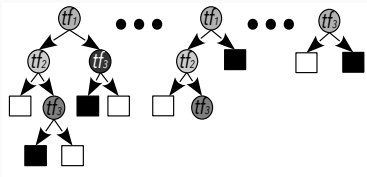


- Avg. SSE between **gene exp.** and **cluster centers** for different values of  $K$ .
- **Elbow** for  $k = 5$  clusters.



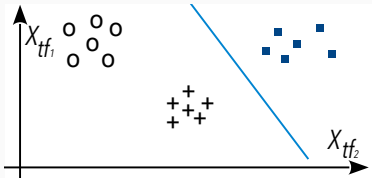
# Classification Algorithms

## Ensemble of decision trees<sup>1</sup>



- Random Forest (RF)
- Extremely Randomized Trees (XRT)
- AdaBoost (AB)
- Gradient Boosting (GB)

## Support Vector Machine (SVM)<sup>1</sup>



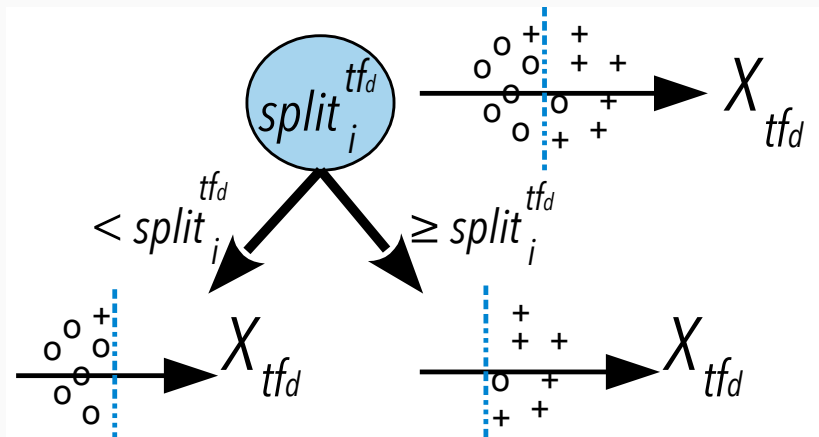
- One-vs-All linear multi-class SVM

<sup>1</sup>Implementations from **scikit-learn** (Python 3.7)

# Feature Importance | Decision Tree based

Impurity gain (e.g., GINI, Entropy) for the  $i$ -th split along feature  $X_{tf_d}$

$$\delta(\text{split}_i^{tf_d})$$



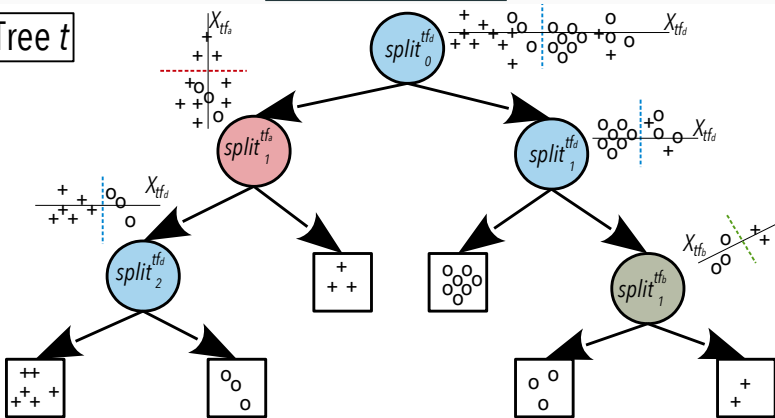


# Feature Importance | Decision Tree based

Importance of feature  $X_{tf_d}$  for tree  $t$

$$\gamma_t^{tf_d} = \sum_i \delta(split_i^{tf_d})$$

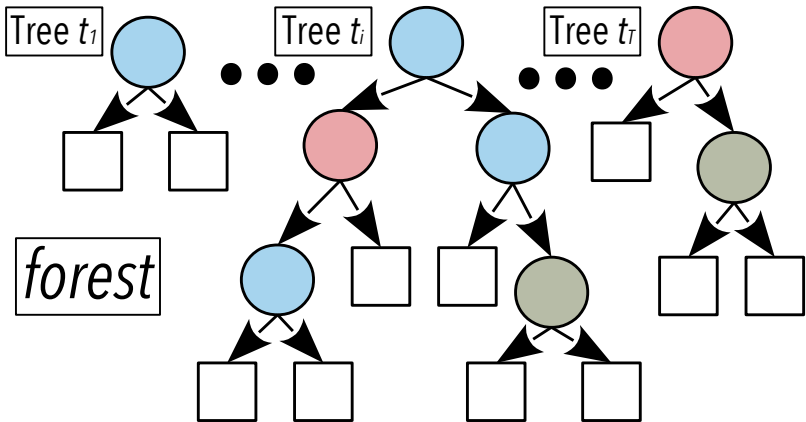
Tree  $t$



# Feature Importance | Decision Tree based

Importance of feature  $X_{tfd}$  for a forest, i.e., set of trees

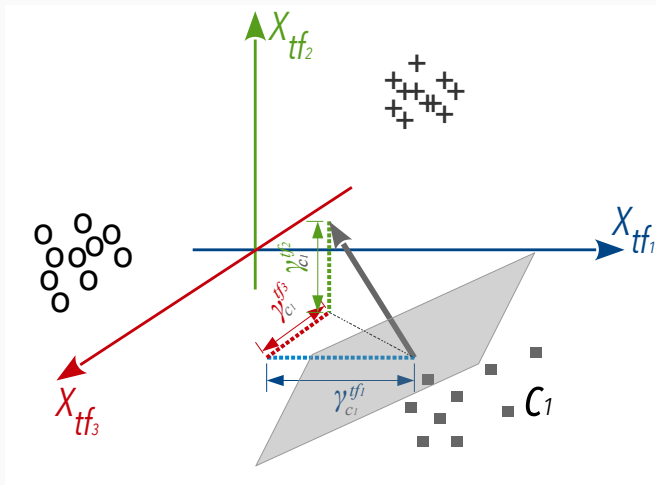
$$\gamma_{tfd}^{tfd} = \frac{\sum_{t \in \text{forest}} \gamma_t^{tfd}}{|\text{forest}|}$$



# Feature Importance | One-vs-all linear SVM

$\gamma_c^{tf_d}$ : Importance of feature  $X_{tf_d}$  for class  $c$  linear SVM  $\rightarrow$

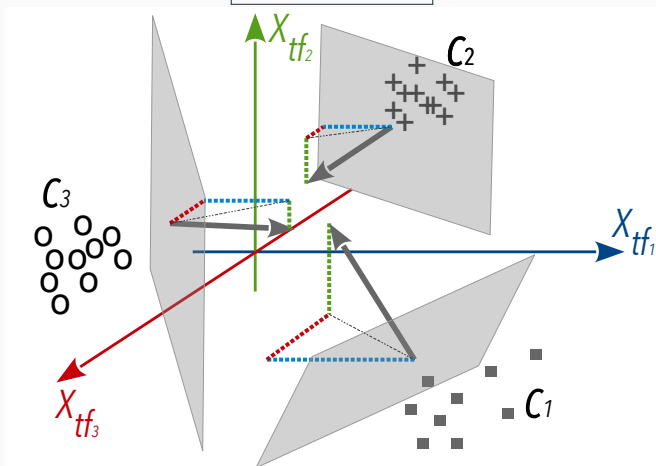
Norm of the separating hyperplane orthogonal vector



# Feature Importance | One-vs-all linear SVM

Importance of feature  $X_{tf_d}$  for a set  $C$  of one-vs-all linear SVM

$$\gamma^{tf_d} = \frac{\sum_{c \in C} \gamma_c^{tf_d}}{|C|}$$



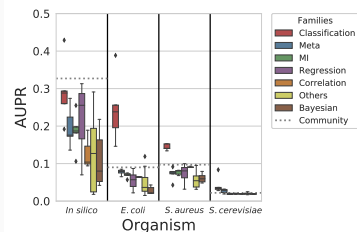
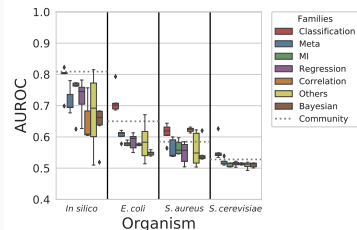
# Evaluation Setup

- Evaluation protocol designed by [Marbach et al. 2012]
- Standard evaluation measures: AUROC and AUPR
- Impact of 7 pre-processing techniques.
- DREAM5 Benchmark datasets
- Comparison w.r.t. 36 methods

Dataset	# condit.	TGs	TFs
<i>In Silico</i>	805	1,643	195
<i>S. aureus</i>	160	2,810	99
<i>E. coli</i>	805	4,511	334
<i>S. cerevisiae</i>	536	5,950	333

Paradigm	# Methods
Community	1
MI	5
Meta	5
Regression	8
Correlation	3
Bayesian	6
Others	8

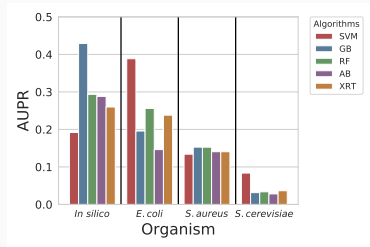
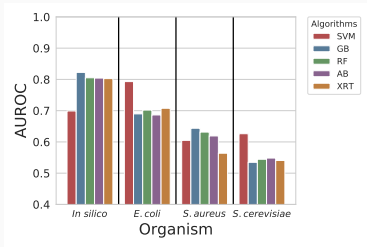
# Results: Comparison with other Paradigms



Paradigm	Avg. AUROC	Avg. AUPR
<b>Classification</b>	<b>0.67</b>	<b>0.18</b>
Community	0.64	0.13
Others	0.58	0.06
MI	0.60	0.09
Meta	0.60	0.09
Regression	0.59	0.09
Correlation	0.59	0.08
Bayesian	0.56	0.05

- **Best AUROC and AUPR on avg.**
- **Surpass community results**
- **Good results for all datasets**

# Results: Classification Methods Comparison



- No ever-winning method.
- Analogous phenomenon reported in [Marbach et al. 2012]

# Conclusion

## Summary

- **Classification** methods **outperform** other **families** on avg.
- Interesting **complementary tool** for the **community**

## Implementation

- **GReNaDIne Python package:**

Gene Regulatory Network Data-driven Inference

- **GitLab repository:**

[gitlab.com/bf2i/grenadine](https://gitlab.com/bf2i/grenadine)

- **Documentation:**

[grenadine.readthedocs.io](https://grenadine.readthedocs.io)



`pip install GReNaDIne`