



# Discussing gene selection methods based on gene expression data

GT LEGO, 23 november 2023

Myriam Bontonou, Anaïs Haget, Maria Boulougouri, Jean-Michel Arbona, Benjamin Audit, Pierre Borgnat

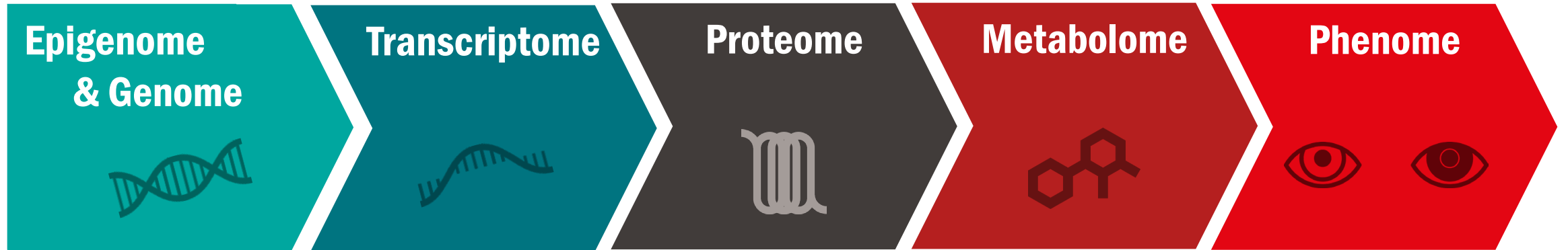


# Overview

- Functional genomics: objective and available data
- Selecting relevant genes with statistical and machine learning methods
- Application to cancer classification on bulk RNA-seq data

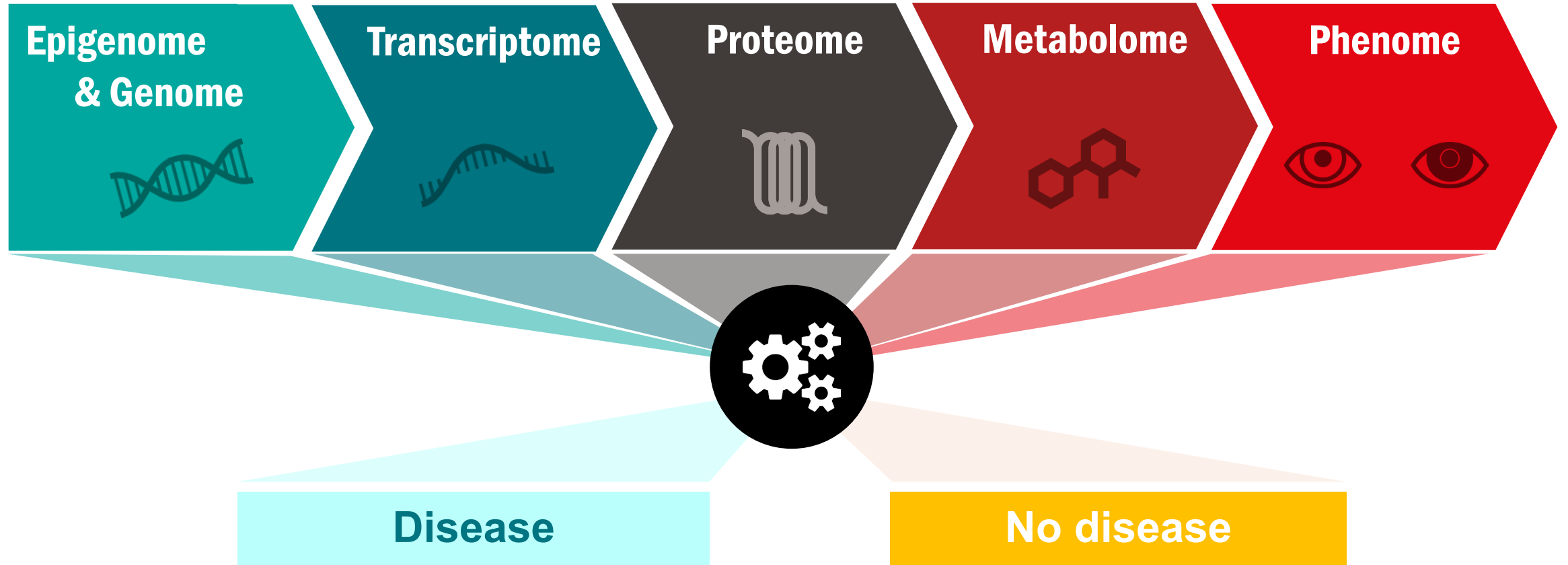
- **Functional genomics: objective and available data**
- Selecting relevant genes with statistical and machine learning methods
- Application to cancer classification on bulk RNA-seq data

# Introduction to functional genomics



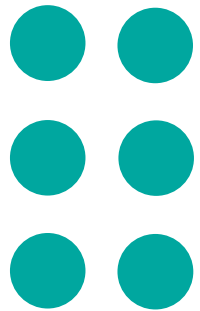
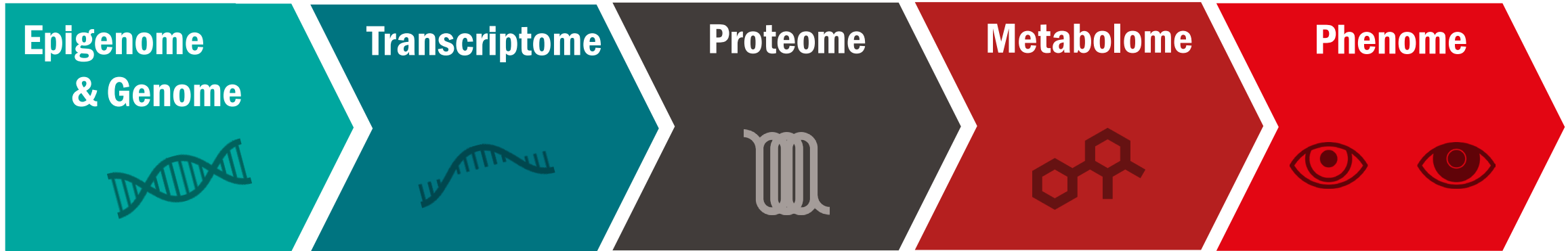
Adapted from Zitnik et al. Machine learning for integrating data in biology and medicine: Principles, practice, and opportunities. Inf Fusion (2019).

# Introduction to functional genomics

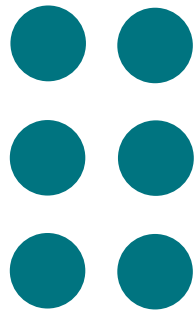


Adapted from Zitnik et al. Machine learning for integrating data in biology and medicine: Principles, practice, and opportunities. Inf Fusion (2019).

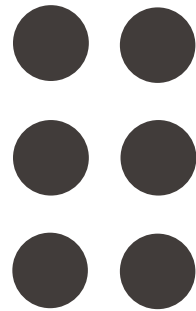
# Introduction to functional genomics



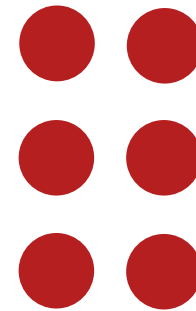
Epigenetic factors / Mutation



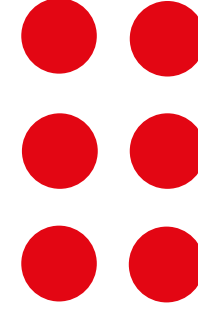
Genes



Proteins

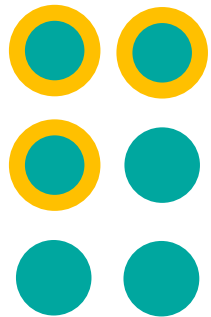
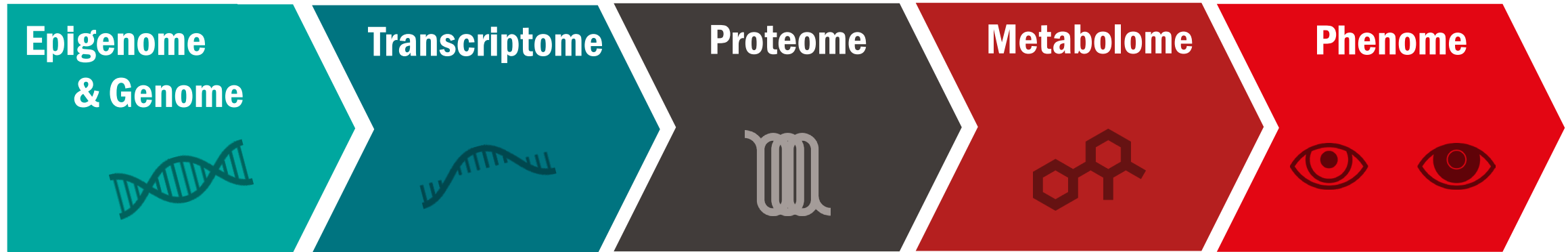


Metabolites

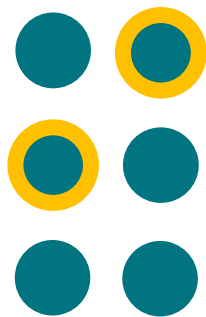


Phenotypes

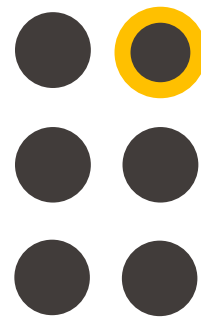
# Introduction to functional genomics



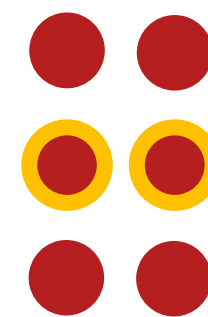
Epigenetic factors / Mutation



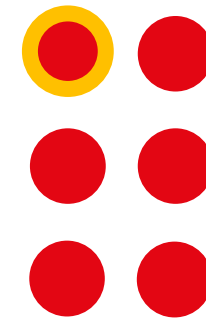
Genes



Proteins

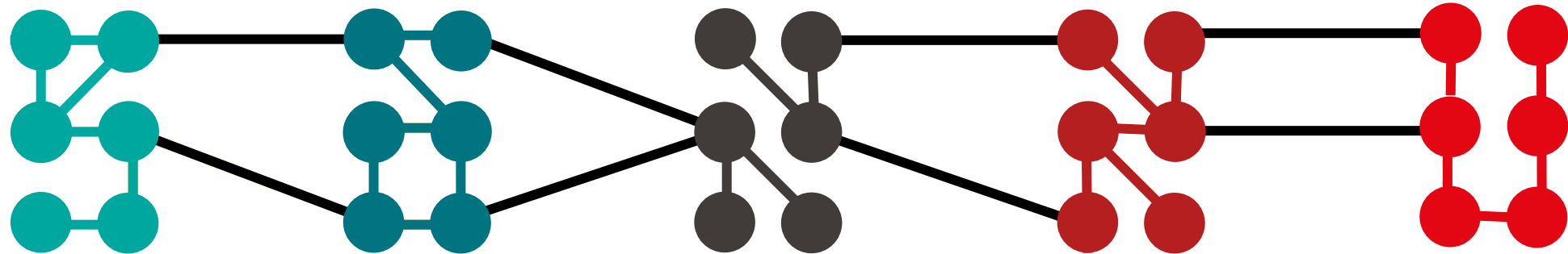
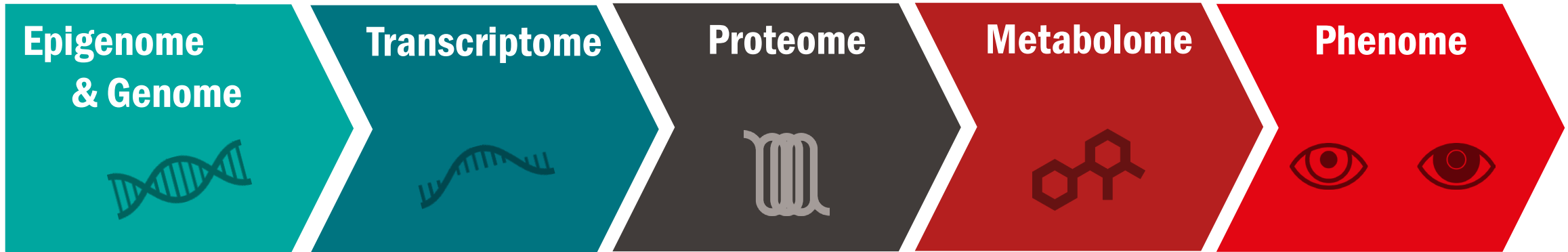


Metabolites



Phenotypes

# Introduction to functional genomics



Epigenetic factors / Mutation

Genes

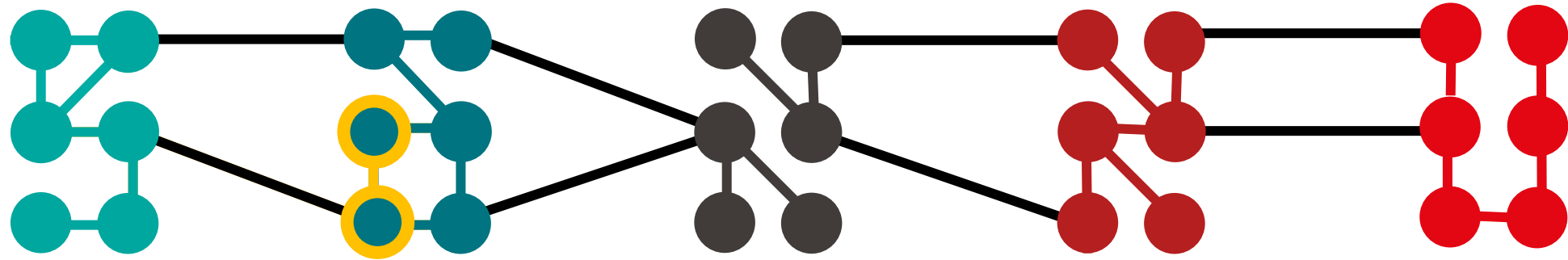
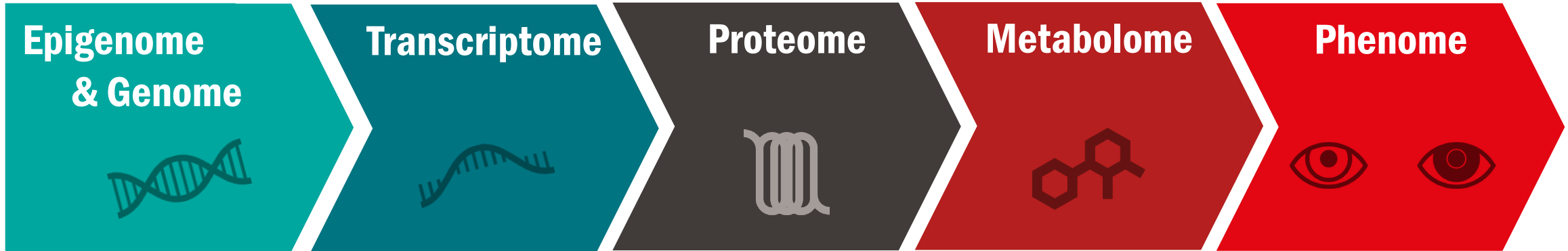
Proteins

Metabolites

Phenotypes



# Introduction to functional genomics



Epigenetic factors / Mutation

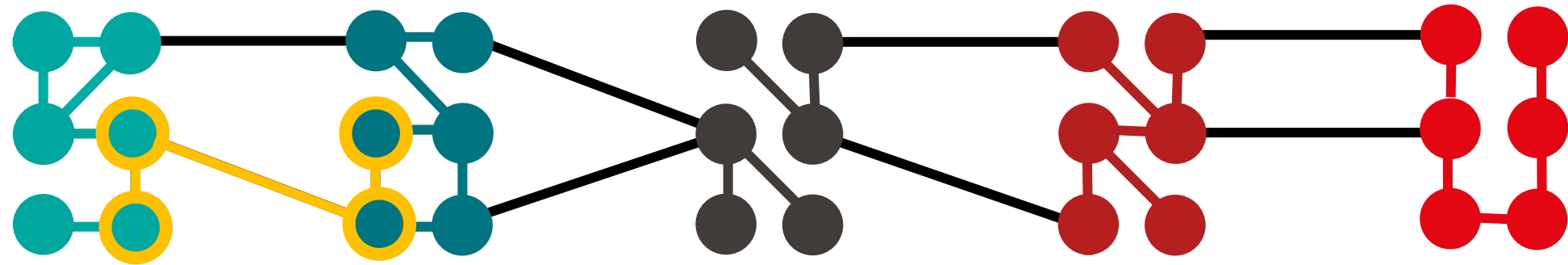
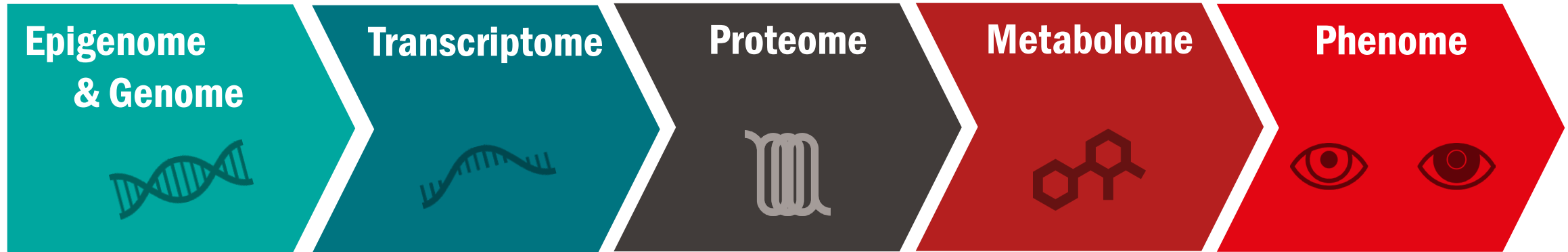
Genes

Proteins

Metabolites

Phenotypes

# Introduction to functional genomics



Epigenetic factors / Mutation

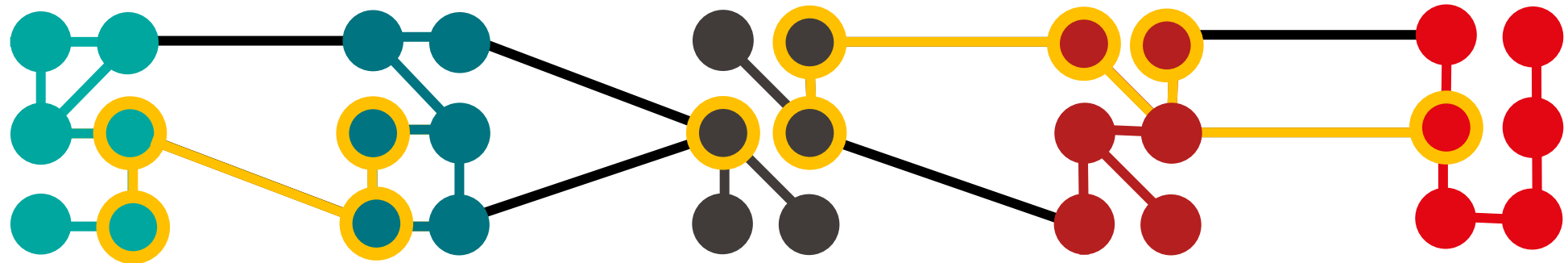
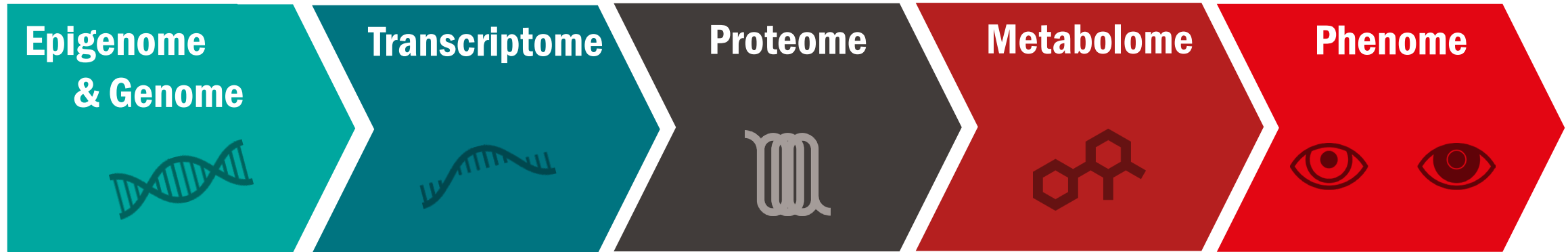
Genes

Proteins

Metabolites

Phenotypes

# Introduction to functional genomics



Epigenetic factors / Mutation

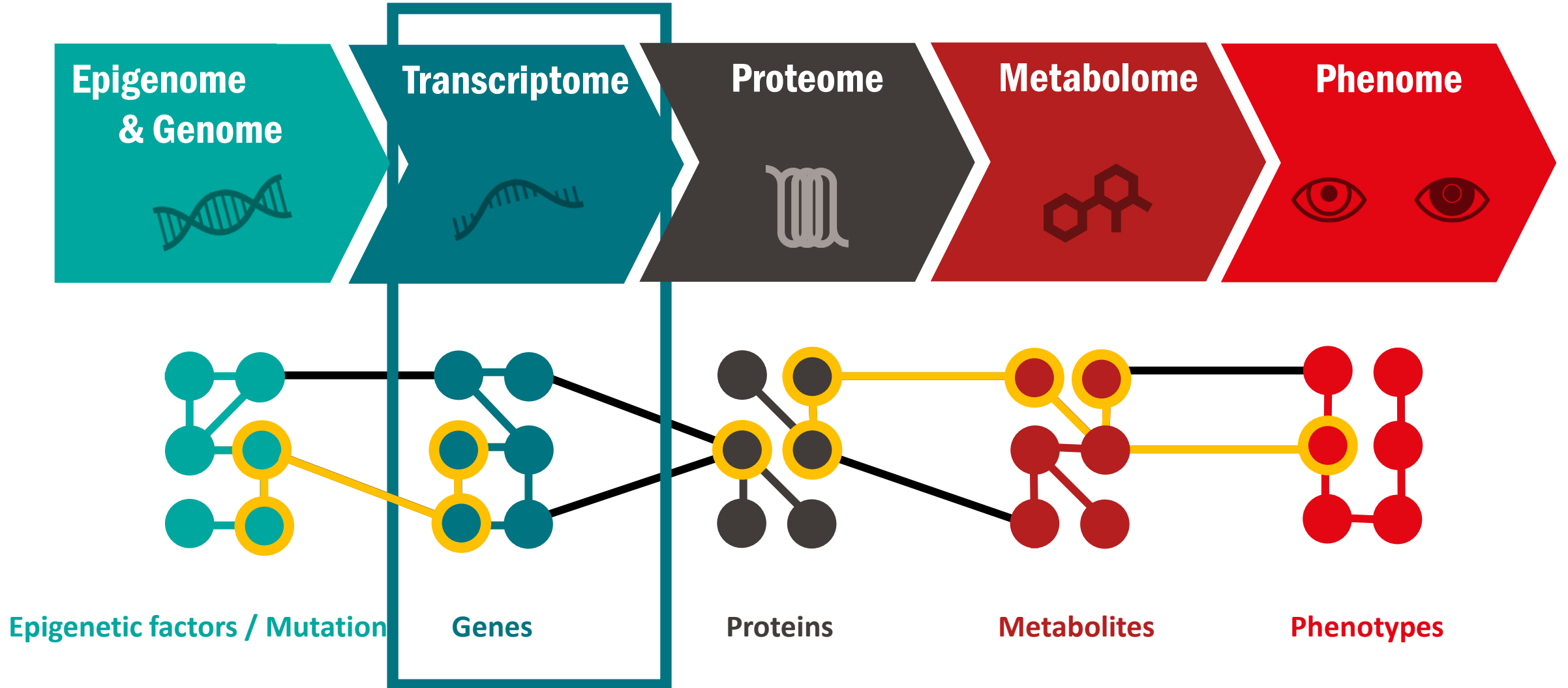
Genes

Proteins

Metabolites

Phenotypes

# Introduction to functional genomics



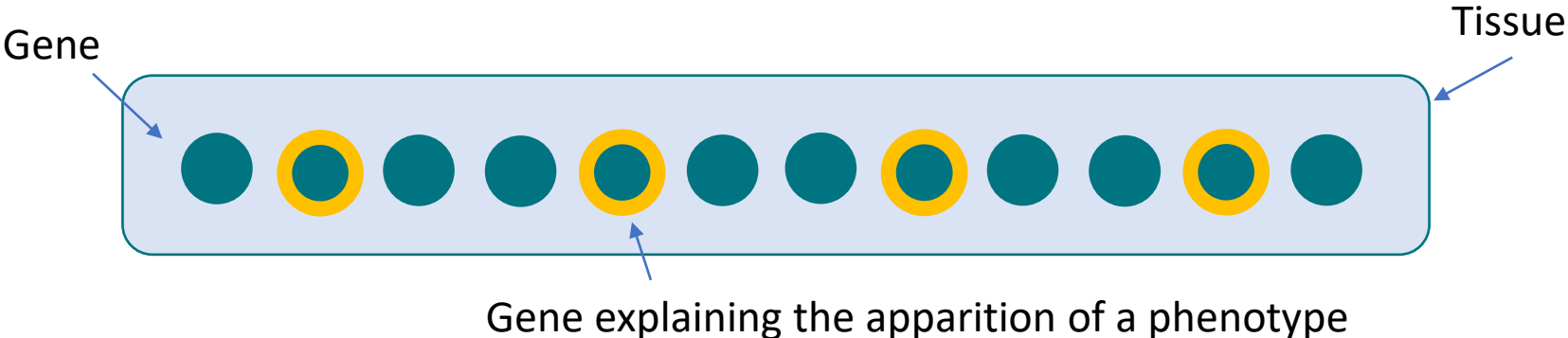
- Functional genomics: objective and available data
- **Selecting relevant genes with statistical and machine learning methods**
- Application to cancer classification on bulk RNA-seq data

# Statistics and machine learning on gene expression data

**Objective** Understand why certain phenotypes appear based on the genes expressed in a tissue.

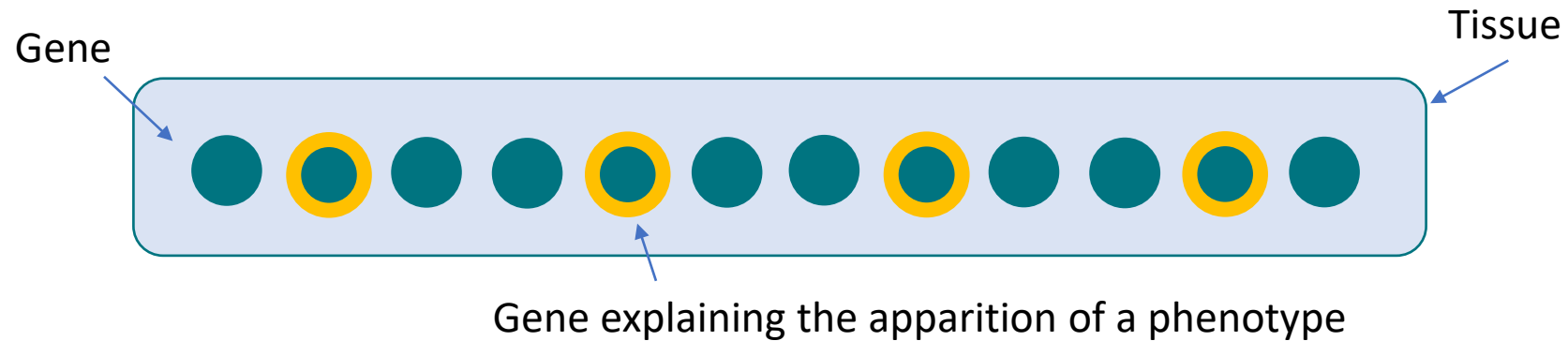
# Statistics and machine learning on gene expression data

**Objective** Understand why certain phenotypes appear based on the genes expressed in a tissue.



# Statistics and machine learning on gene expression data

**Objective** Understand why certain phenotypes appear based on the genes expressed in a tissue.

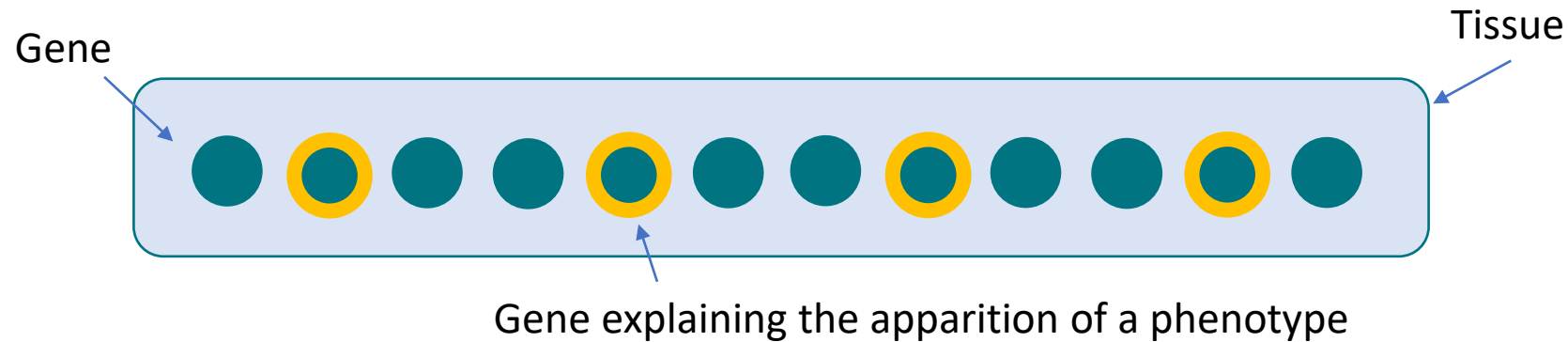


**Datasets** The Cancer Genome Atlas (TCGA), TARGET, Genotype-Tissue Expression project (GTEx).



# Statistics and machine learning on gene expression data

**Objective** Understand why certain phenotypes appear based on the genes expressed in a tissue.



**Datasets** The Cancer Genome Atlas (TCGA), TARGET, Genotype-Tissue Expression project (GTEx).

Dataset	Task	Source	# classes	# samples (max / min per class)	# genes
ttg-breast	Healthy vs Tumor	TCGA TARGET GTEx	2	1384 (1092 / 292)	14373
BRCA-pam	PAM50 classes	Legacy TCGA	5	916 (421 / 67)	13896

# Statistics and machine learning on gene expression data

How genes are selected with statistical and machine learning methods?

# Statistics and machine learning on gene expression data

How genes are selected with statistical and machine learning methods?

**Feature  
selection**

Filters

Embedded methods

Post-hoc explainability

# Statistics and machine learning on gene expression data

How genes are selected with statistical and machine learning methods?

**Feature  
selection**

Filters

Embedded methods

Post-hoc explainability

Genes are ranked according to...

# Statistics and machine learning on gene expression data

How genes are selected with statistical and machine learning methods?

**Feature  
selection**

Filters

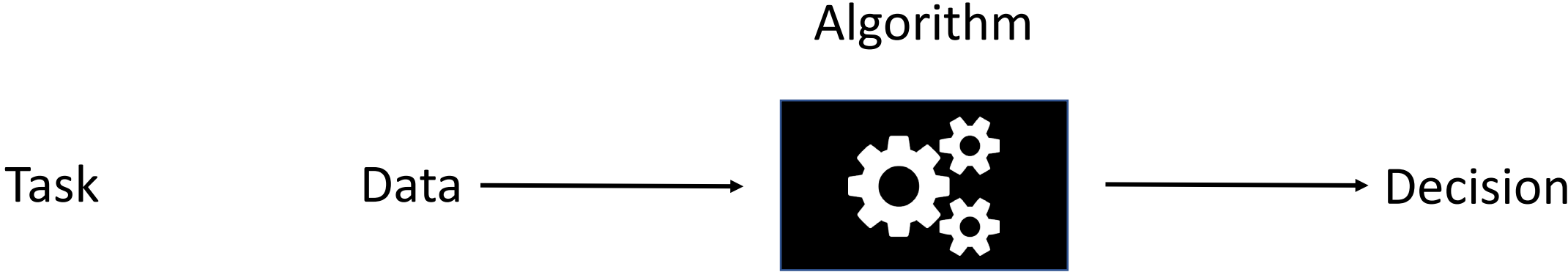
Embedded methods

Post-hoc explainability

Genes are ranked according to...

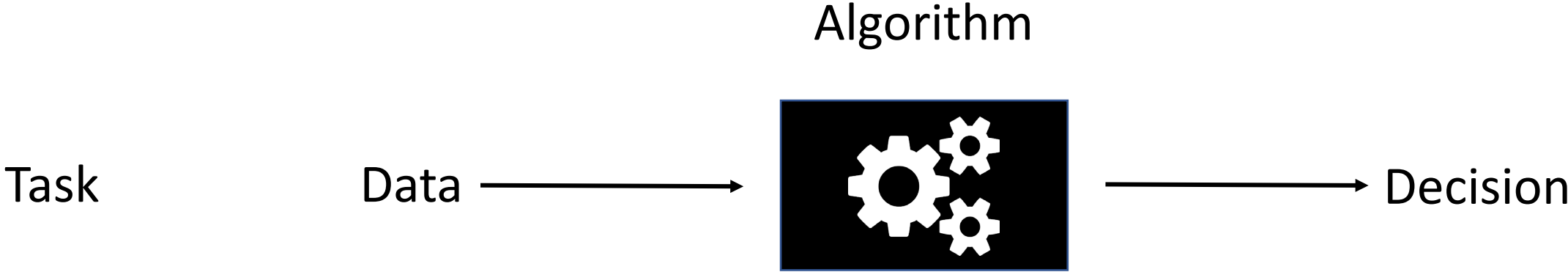
Intrinsic score related to a gene

# Statistics and machine learning on gene expression data



User

# Statistics and machine learning on gene expression data

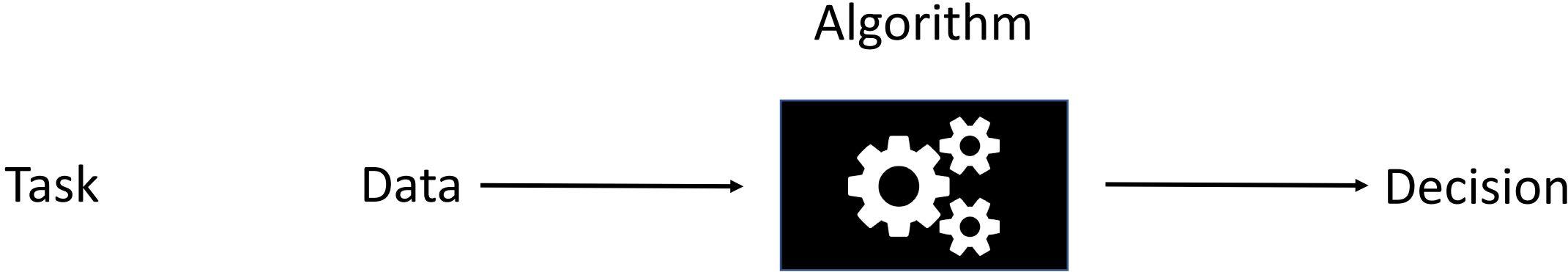


Medical diagnosis




User

# Statistics and machine learning on gene expression data



Medical diagnosis

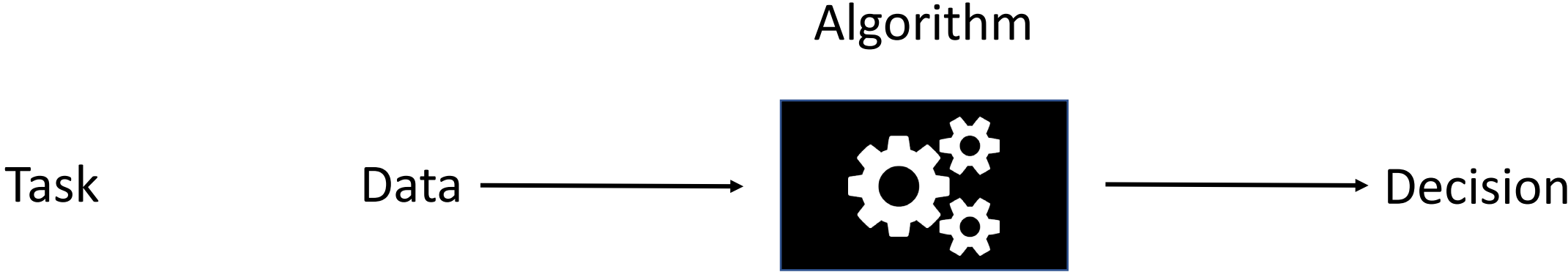
- ✓ Genome 
- ✓ Transcriptome 
- ✓ Proteome 
- ✓ Phenome... 



User



# Statistics and machine learning on gene expression data



Medical diagnosis

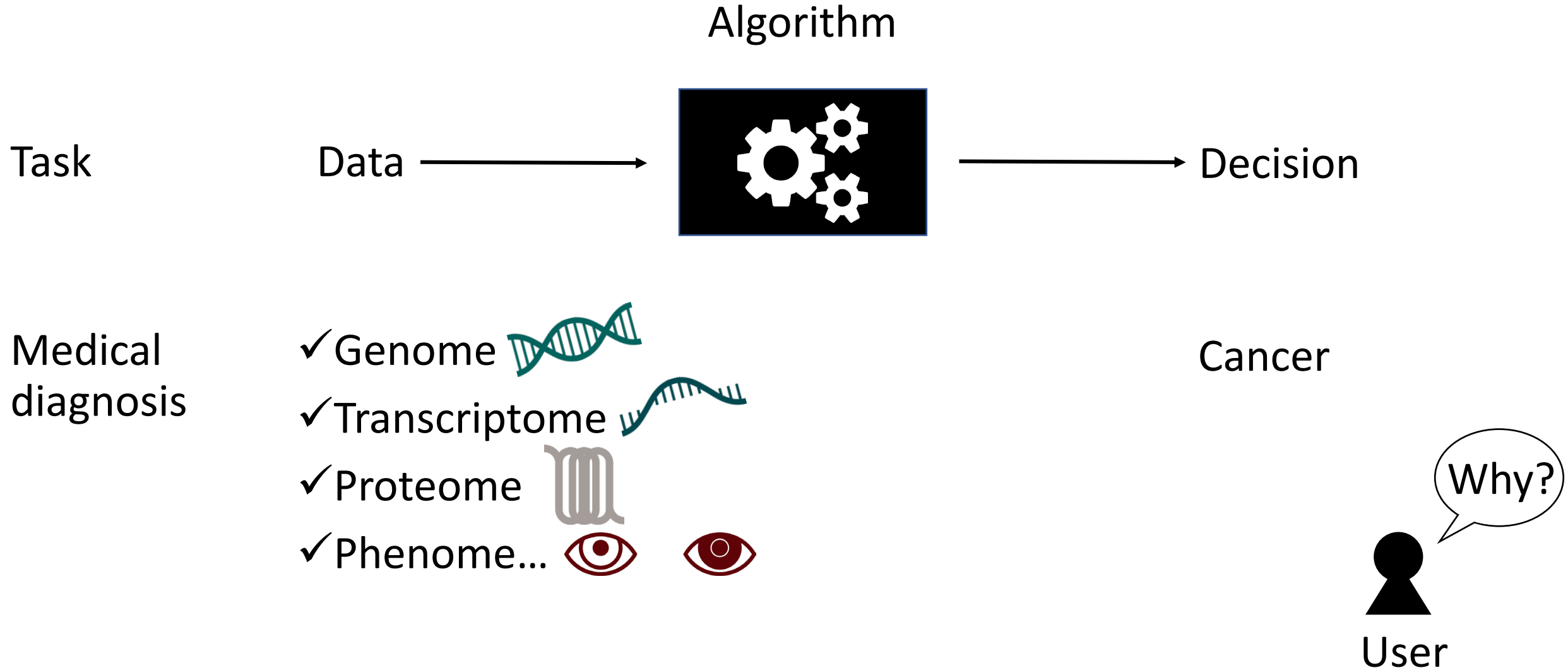
- ✓ Genome 
- ✓ Transcriptome 
- ✓ Proteome 
- ✓ Phenome... 

Cancer



User

# Statistics and machine learning on gene expression data



# Statistics and machine learning on gene expression data

How genes are selected with statistical and machine learning methods?

**Feature  
selection**

Filters

Embedded methods

Post-hoc explainability

Genes are ranked according to...

Intrinsic score related to a gene

# Statistics and machine learning on gene expression data

How genes are selected with statistical and machine learning methods?

**Feature  
selection**

Filters

Embedded methods

Post-hoc explainability

Genes are ranked according to...

Intrinsic score related to a gene

Intrinsic classifier score

# Statistics and machine learning on gene expression data

How genes are selected with statistical and machine learning methods?

## Feature selection



Genes are ranked according to...

Intrinsic score related to a gene

Intrinsic classifier score

Explainability score computed after training a classifier

# Statistics and machine learning on gene expression data

How genes are selected with statistical and machine learning methods?

## Feature selection

Filters

Embedded methods

Post-hoc explainability

Genes are ranked according to...

Intrinsic score related to a gene

Intrinsic classifier score

Explainability score computed after training a classifier

Examples

# Statistics and machine learning on gene expression data

How genes are selected with statistical and machine learning methods?

## Feature selection

Filters

Embedded methods

Post-hoc explainability

Genes are ranked according to...

Intrinsic score related to a gene

Intrinsic classifier score

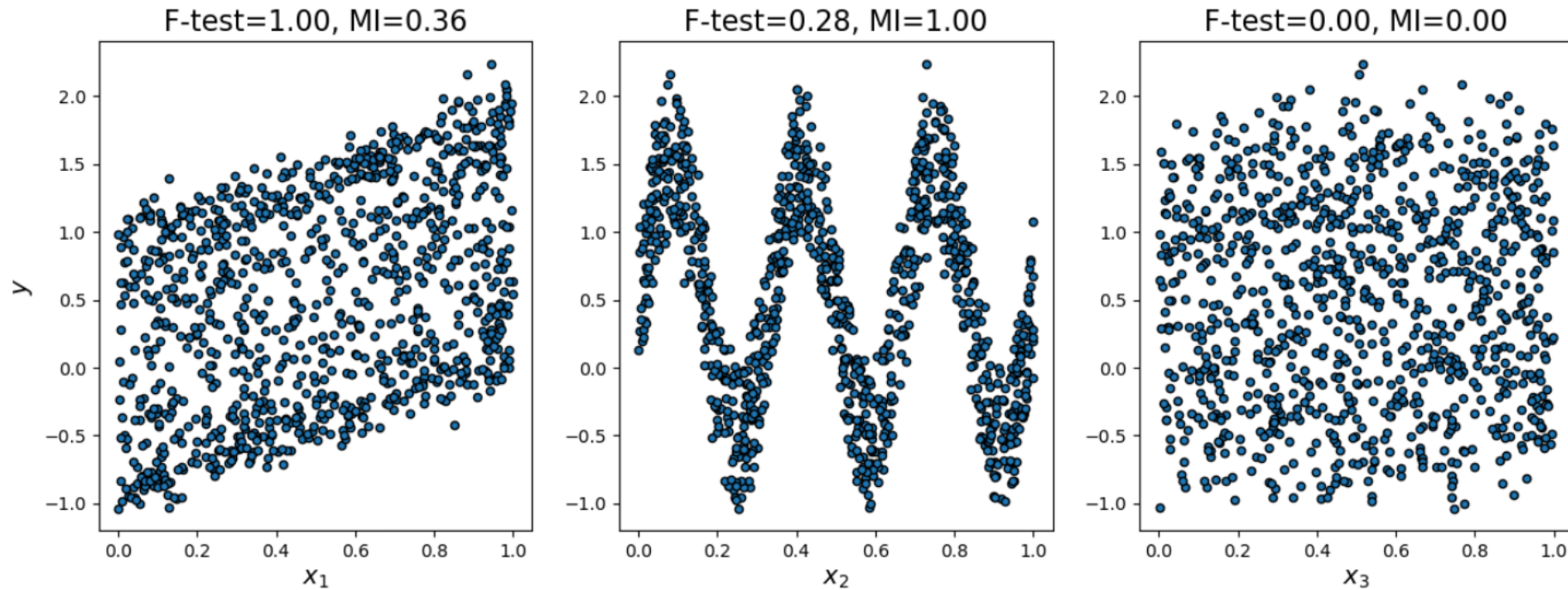
Explainability score computed after training a classifier

Examples

Mutual information

# Statistics and machine learning on gene expression data

## Mutual information (MI)



Example on scikit-learn.



# Statistics and machine learning on gene expression data

How genes are selected with statistical and machine learning methods?

## Feature selection

Filters

Embedded methods

Post-hoc explainability

Genes are ranked according to...

Intrinsic score related to a gene

Intrinsic classifier score

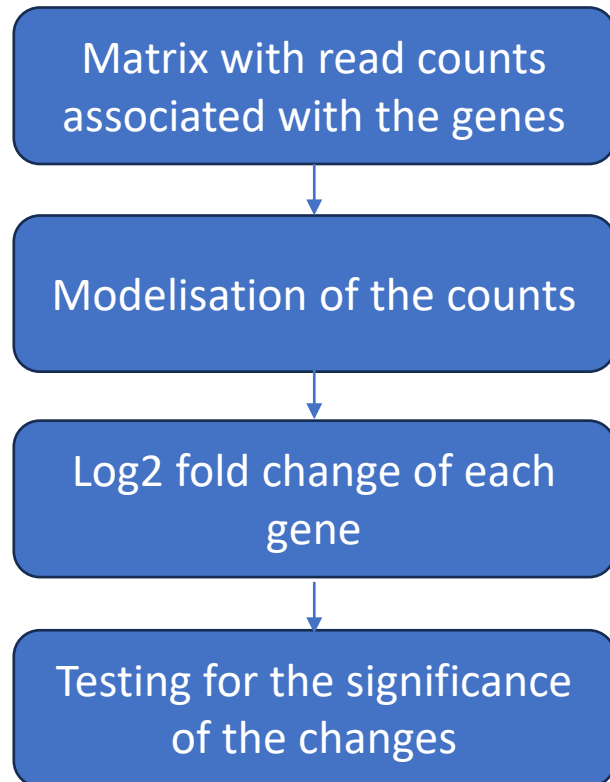
Explainability score computed after training a classifier

Examples

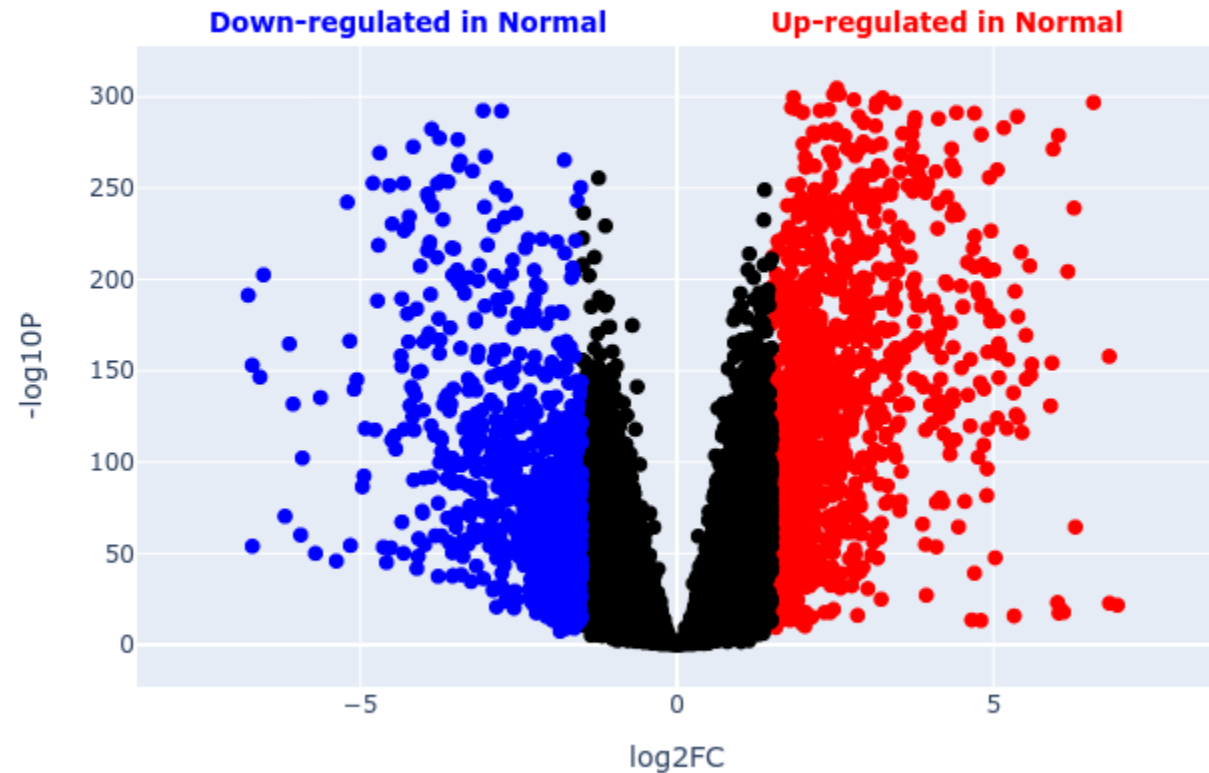
Mutual information  
DESeq2

# Statistics and machine learning on gene expression data

## DEseq2 (DE)



Tumor vs. Normal Signature | Volcano Plot



Love et al. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. Genome biology (2014).

# Statistics and machine learning on gene expression data

How genes are selected with statistical and machine learning methods?

## Feature selection

Filters

Embedded methods

Post-hoc explainability

Genes are ranked according to...

Intrinsic score related to a gene

Intrinsic classifier score

Explainability score computed after training a classifier

Examples

Mutual information  
DESeq2

Logistic regression with L1 penalty

# Statistics and machine learning on gene expression data

Logistic regression trained with a L1 penalty (LR+L1)

$$f(\mathbf{x}) = \textit{softmax}(\mathbf{w}_0 + \mathbf{w}_1\mathbf{x}_1 + \mathbf{w}_2\mathbf{x}_2 + \mathbf{w}_3\mathbf{x}_3 + \mathbf{w}_4\mathbf{x}_4 + \dots + \mathbf{w}_{10000}\mathbf{x}_{10000})$$

# Statistics and machine learning on gene expression data

Logistic regression trained with a L1 penalty (LR+L1)

$$f(\mathbf{x}) = \textit{softmax}(\mathbf{w}_0 + \mathbf{w}_1\mathbf{x}_1 + \mathbf{w}_2\mathbf{x}_2 + \mathbf{w}_3\mathbf{x}_3 + \mathbf{w}_4\mathbf{x}_4 + \dots + \mathbf{w}_{10000}\mathbf{x}_{10000})$$

↓ after training

$$f(\mathbf{x}) = \textit{softmax}(\mathbf{w}_0 + 0 + \mathbf{w}_2\mathbf{x}_2 + 0 + 0 + \dots + \mathbf{w}_{10000}\mathbf{x}_{10000})$$

# Statistics and machine learning on gene expression data

How genes are selected with statistical and machine learning methods?

## Feature selection

Filters

Embedded methods

Post-hoc explainability

Genes are ranked according to...

Intrinsic score related to a gene

Intrinsic classifier score

Explainability score computed after training a classifier

Examples

Mutual information  
DESeq2

Logistic regression with L1 penalty

Logistic regression  
Multilayer perceptron  
Graph neural network

# Statistics and machine learning on gene expression data

Logistic regression (LR+L2)

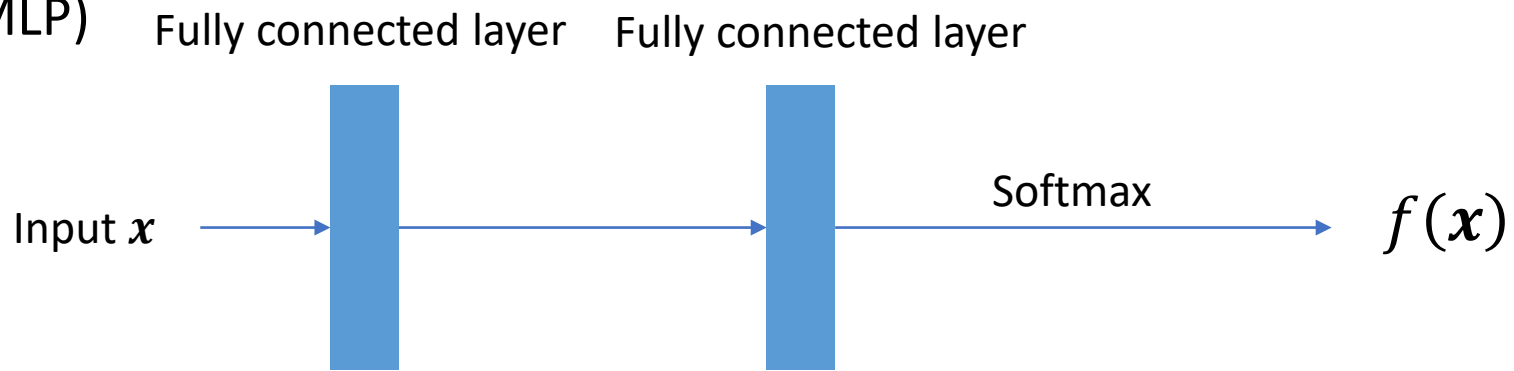
$$f(\mathbf{x}) = \textit{softmax}(\mathbf{w}_0 + \mathbf{w}_1x_1 + \mathbf{w}_2x_2 + \mathbf{w}_3x_3 + \mathbf{w}_4x_4 + \dots + \mathbf{w}_{10000}x_{10000})$$

# Statistics and machine learning on gene expression data

Logistic regression (LR+L2)

$$f(\mathbf{x}) = \text{softmax}(\mathbf{w}_0 + \mathbf{w}_1x_1 + \mathbf{w}_2x_2 + \mathbf{w}_3x_3 + \mathbf{w}_4x_4 + \dots + \mathbf{w}_{10000}x_{10000})$$

Multilayer perceptron (MLP)



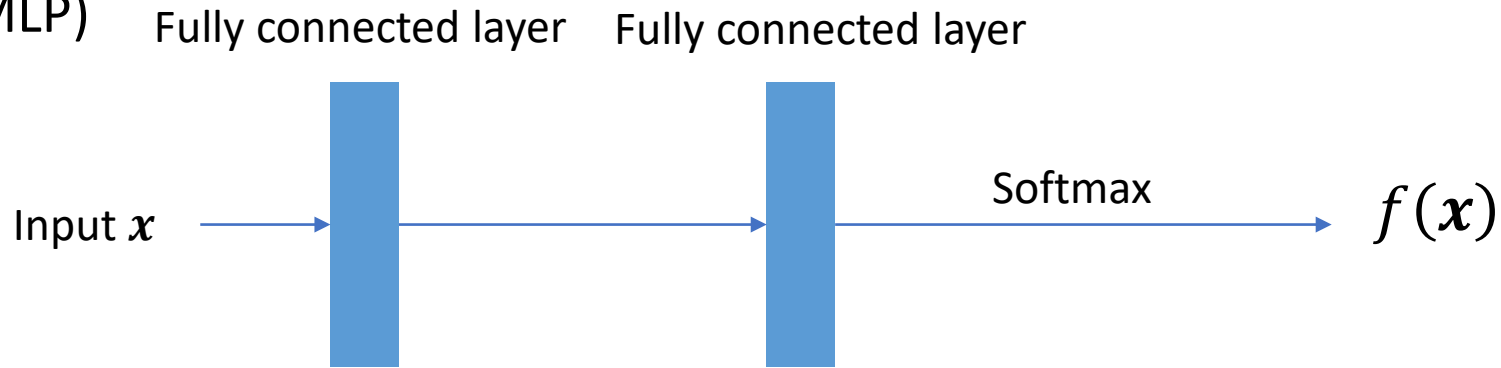


# Statistics and machine learning on gene expression data

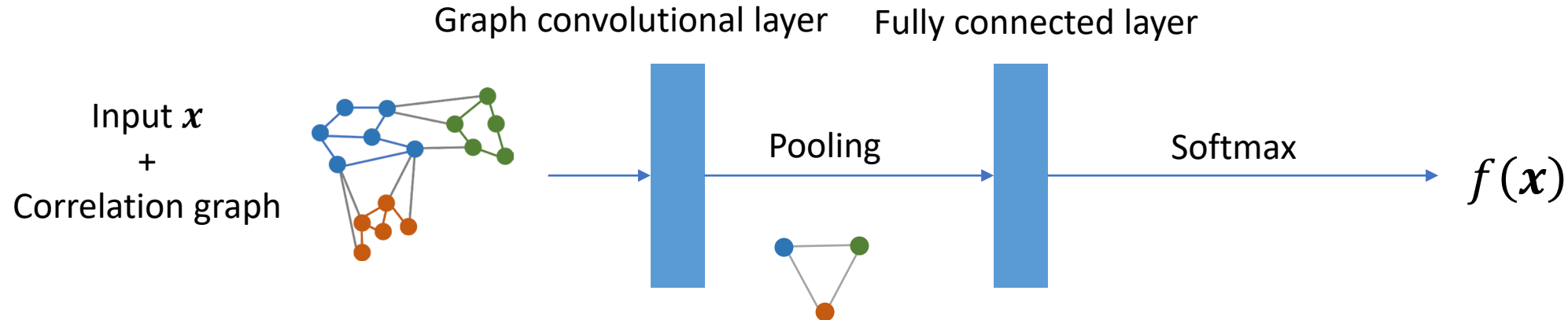
Logistic regression (LR+L2)

$$f(\mathbf{x}) = \text{softmax}(\mathbf{w}_0 + \mathbf{w}_1x_1 + \mathbf{w}_2x_2 + \mathbf{w}_3x_3 + \mathbf{w}_4x_4 + \dots + \mathbf{w}_{10000}x_{10000})$$

Multilayer perceptron (MLP)



Graph convolutional network (GCN)



# Statistics and machine learning on gene expression data

Post-hoc explainability

Input  $x \in \mathbb{R}^F$  (tumor tissue)

Reference  $x'$  (average of healthy tissues)

Model  $f: \mathbb{R}^F \rightarrow \mathbb{R}^C$

# Statistics and machine learning on gene expression data

## Post-hoc explainability

Input  $\mathbf{x} \in \mathbb{R}^F$  (tumor tissue)      Reference  $\mathbf{x}'$  (average of healthy tissues)

Model  $f: \mathbb{R}^F \rightarrow \mathbb{R}^C$

Explainability method attributing to each gene  $i$  of  $\mathbf{x}$  a score  $\phi_i(\mathbf{x})$ . *Here, integrated gradients.*

$$\phi_i(\mathbf{x}) = (\mathbf{x}_i - \mathbf{x}'_i) \int_{\alpha=0}^1 \frac{\partial f_c(\mathbf{z})}{\partial x_i} d\alpha \quad \text{with } \mathbf{z} = (\mathbf{1} - \alpha)\mathbf{x}' + \alpha\mathbf{x}$$

# Statistics and machine learning on gene expression data

## Post-hoc explainability

Input  $\mathbf{x} \in \mathbb{R}^F$  (tumor tissue)      Reference  $\mathbf{x}'$  (average of healthy tissues)

Model  $f: \mathbb{R}^F \rightarrow \mathbb{R}^C$

Explainability method attributing to each gene  $i$  of  $\mathbf{x}$  a score  $\phi_i(\mathbf{x})$ . *Here, integrated gradients.*

$$\phi_i(\mathbf{x}) = (\mathbf{x}_i - \mathbf{x}'_i) \int_{\alpha=0}^1 \frac{\partial f_c(\mathbf{z})}{\partial x_i} d\alpha \quad \text{with } \mathbf{z} = (\mathbf{1} - \alpha)\mathbf{x}' + \alpha\mathbf{x}$$

Genes ranked in decreasing order of average scores.

- Functional genomics: objective and available data
- Selecting relevant genes with statistical and machine learning methods
- **Application to cancer classification on bulk RNA-seq data**

# Application

Dataset	Task	# classes	# samples (max / min per class)	# genes
ttg-breast	Healthy vs Tumor	2	1384 (1092 / 292)	14373
BRCA-pam	PAM50 classes	5	916 (421 / 67)	13896

# Application

Dataset	Task	# classes	# samples (max / min per class)	# genes
ttg-breast	Healthy vs Tumor	2	1384 (1092 / 292)	14373
BRCA-pam	PAM50 classes	5	916 (421 / 67)	13896

Model	Hyperparameter	ttg-breast	BRCA-pam
MLP	# layers	1	1
	# hidden feat.	20	20
GNN	# layers	1	1
	# hidden feat.	1	2

# Application

Dataset	Task	# classes	# samples (max / min per class)	# genes
ttg-breast	Healthy vs Tumor	2	1384 (1092 / 292)	14373
BRCA-pam	PAM50 classes	5	916 (421 / 67)	13896

Dataset	LR+L1	LR+L2	MLP	GNN
ttg-breast	14374	14374	287541	12497
BRCA-pam	69485	69485	278085	118568

Number of parameters

Dataset	LR+L1	LR+L2	MLP	GNN
ttg-breast	262	10	10	24
BRCA-pam	673	6	7	15

Average training duration



# Application

Dataset	LR+L1	LR+L2	MLP	GNN
ttg-breast	$99.7 \pm 0.0$	$98.9 \pm 0.2$	$99.4 \pm 0.3$	$99.0 \pm 0.3$
BRCA-pam	$92.3 \pm 0.0$	$88.7 \pm 1.0$	$87.4 \pm 1.8$	$87.9 \pm 0.9$

Average balanced accuracy (%)

Healthy vs Tumor  
PAM50 breast cancer subtypes

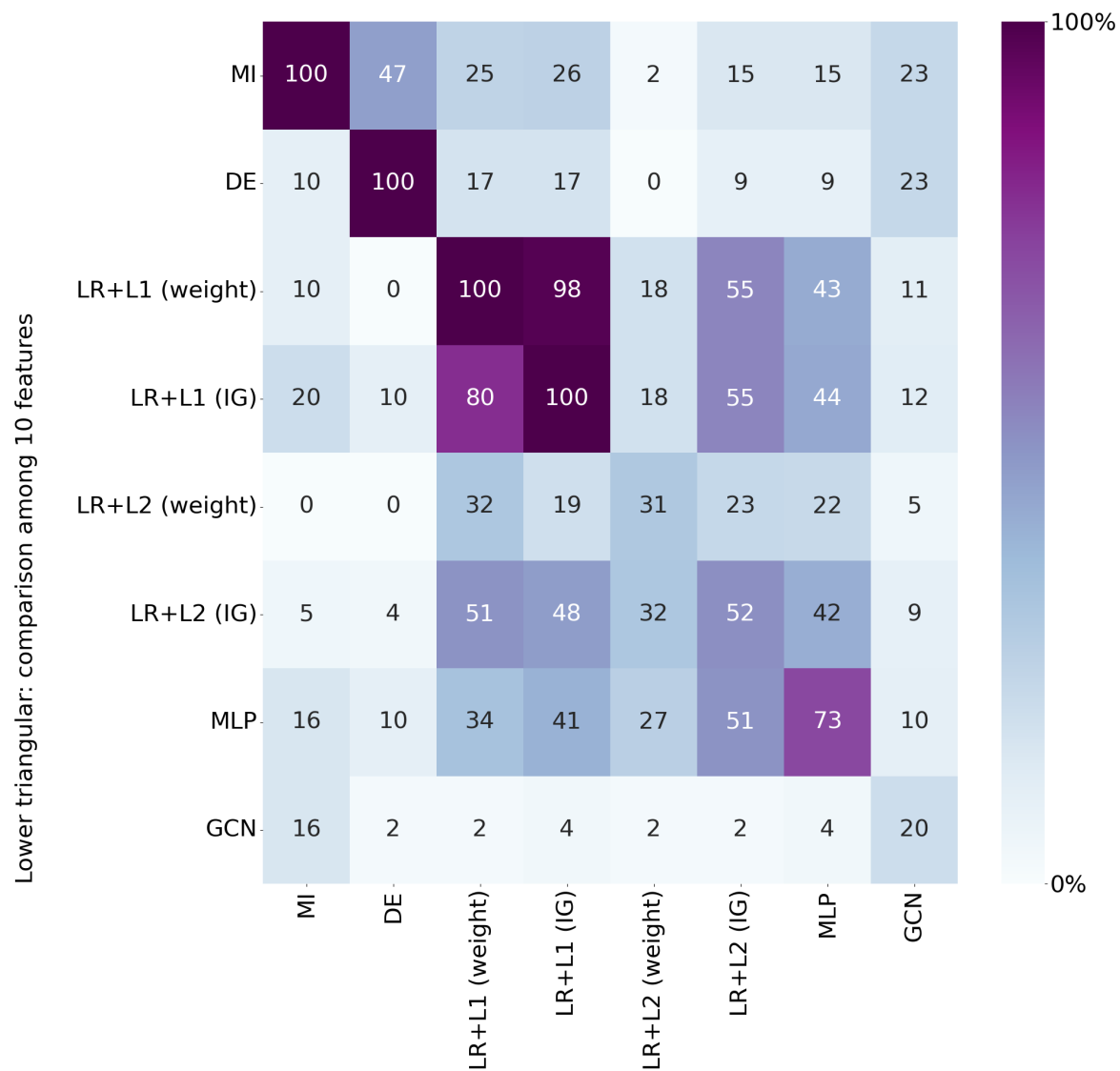
Dataset	LR+L1	LR+L2	MLP	GNN
ttg-breast	$99.5 \pm 0.0$	$98.5 \pm 0.1$	$99.3 \pm 0.2$	$99.1 \pm 0.2$
BRCA-pam	$91.8 \pm 0.0$	$87.4 \pm 1.0$	$87.7 \pm 1.3$	$87.8 \pm 0.7$

Average accuracy (%)

Healthy vs Tumor  
PAM50 breast cancer subtypes

# Application

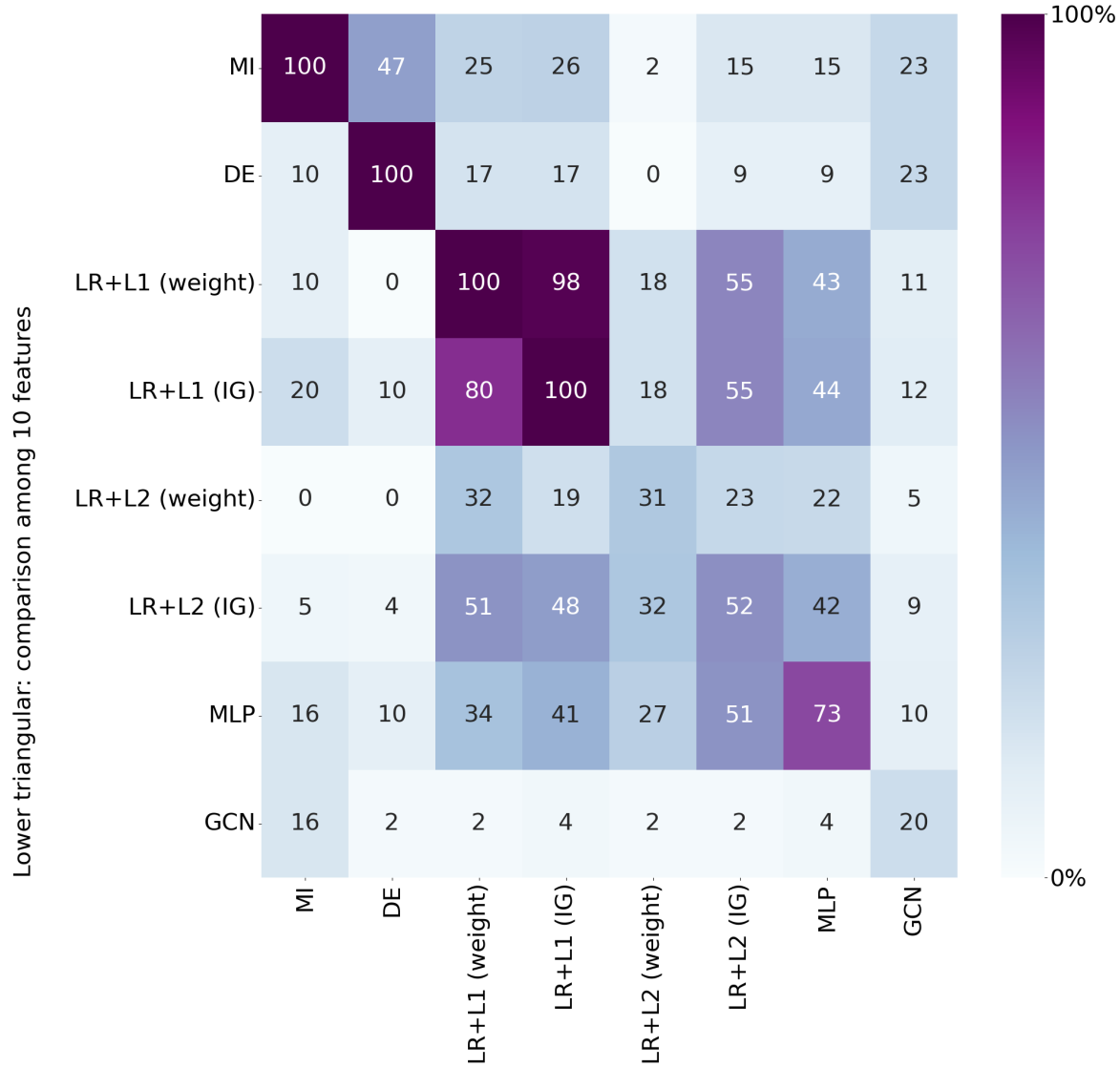
Upper triangular: comparison among 100 features



ttg-breast (healthy vs tumor)

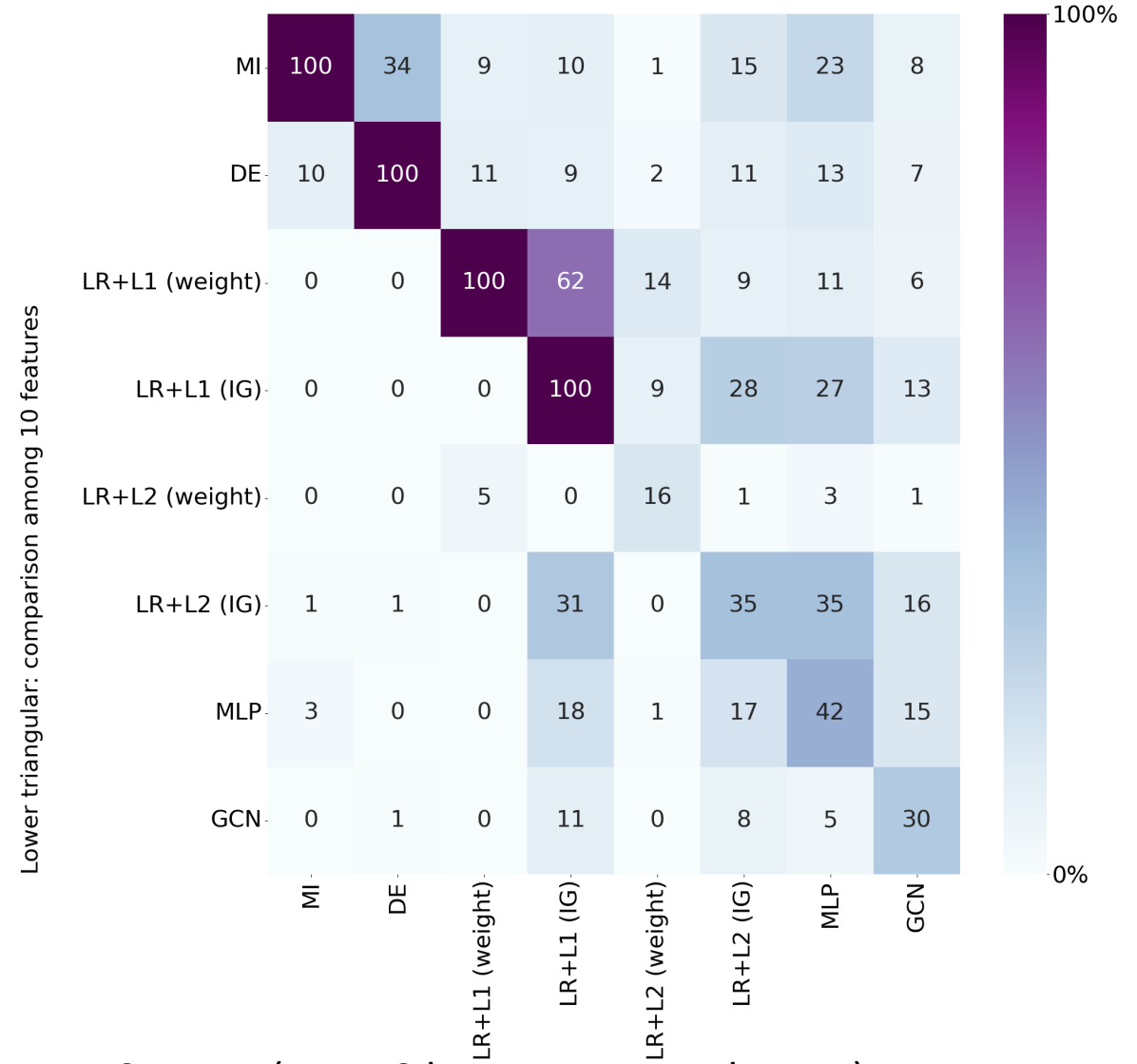
# Application

Upper triangular: comparison among 100 features



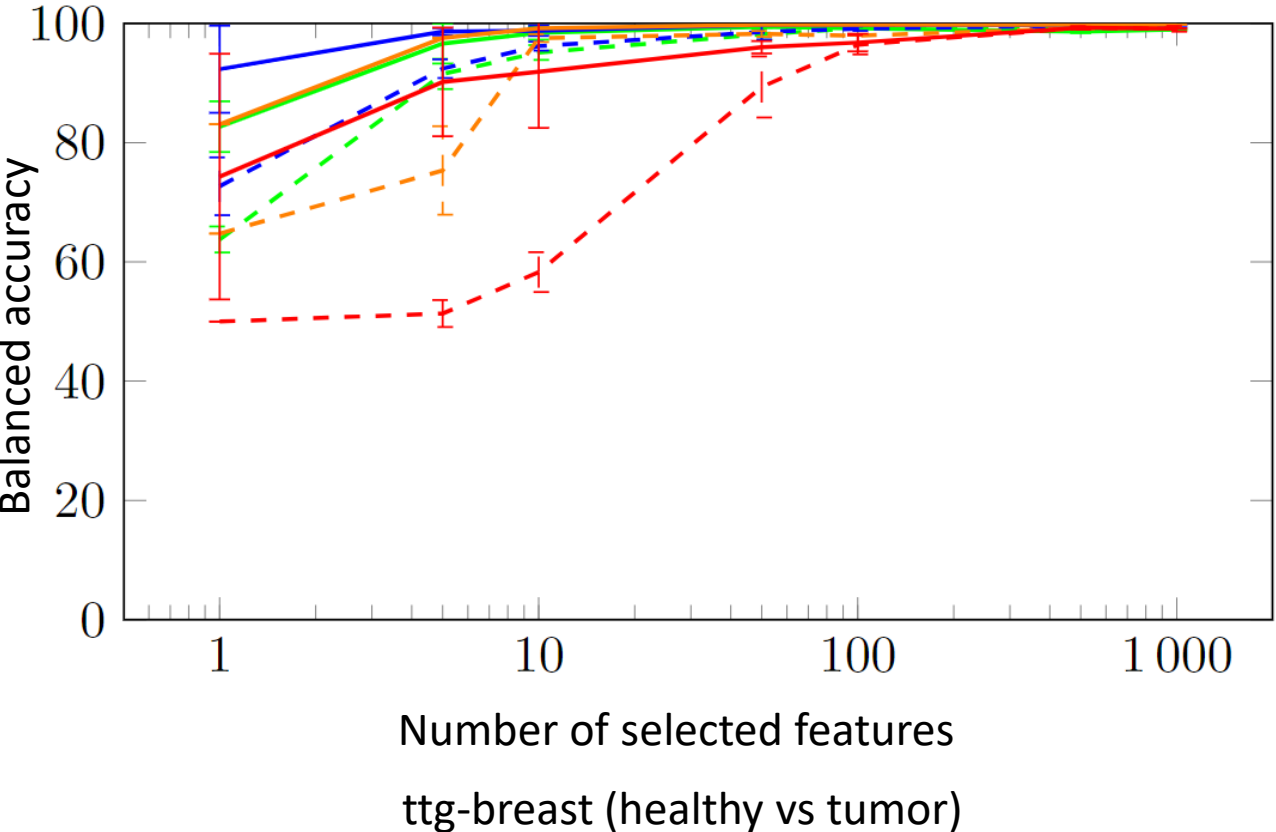
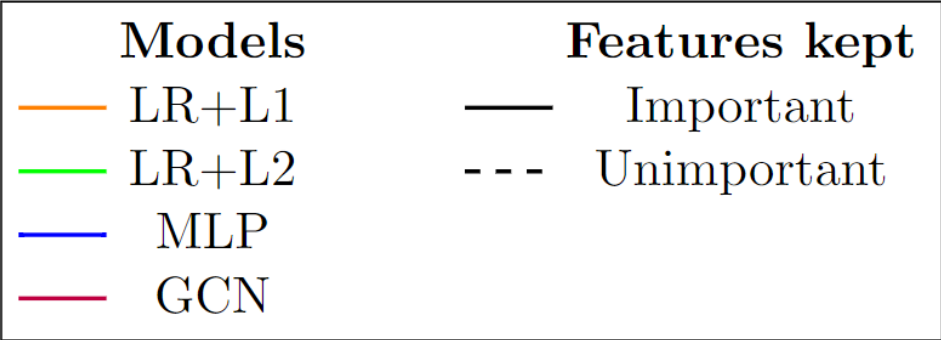
ttg-breast (healthy vs tumor)

Upper triangular: comparison among 100 features

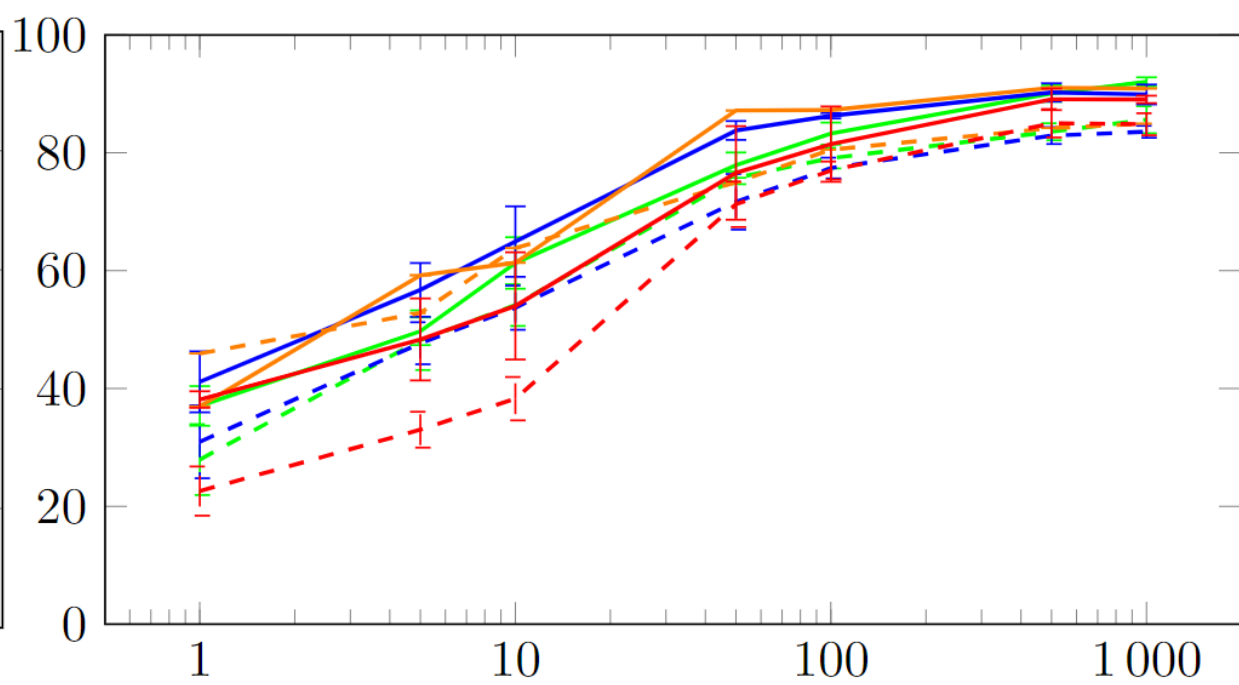
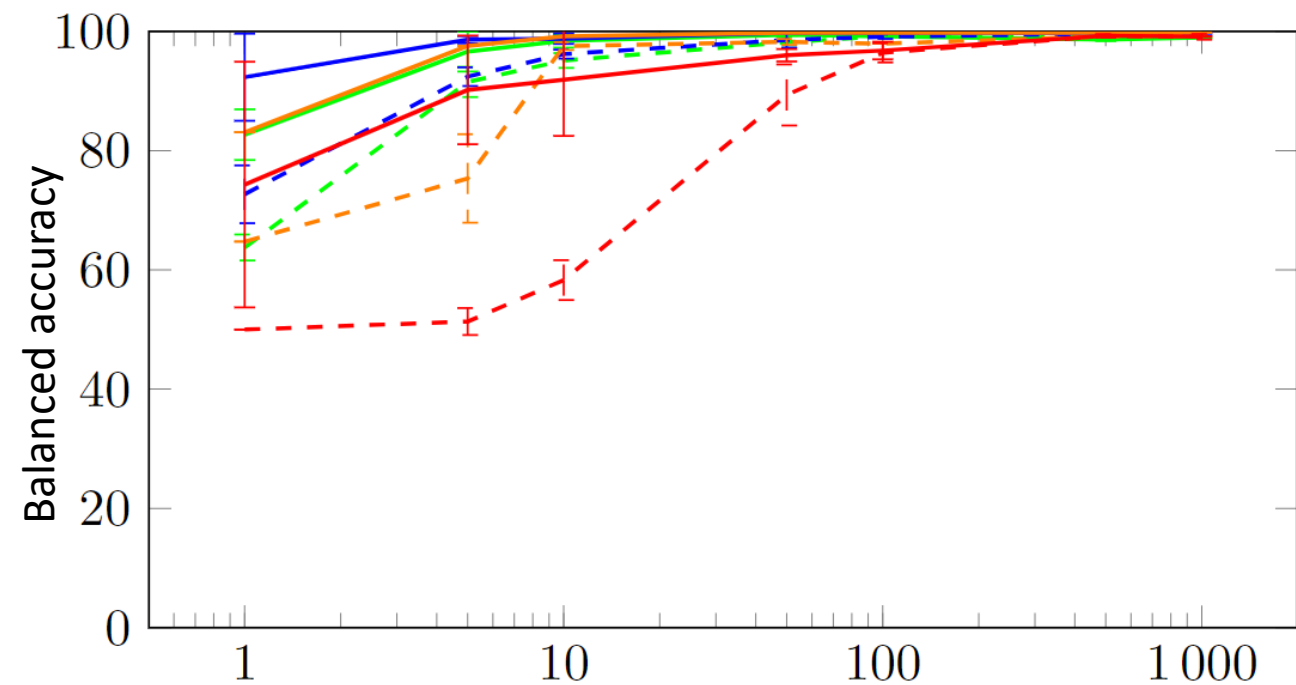
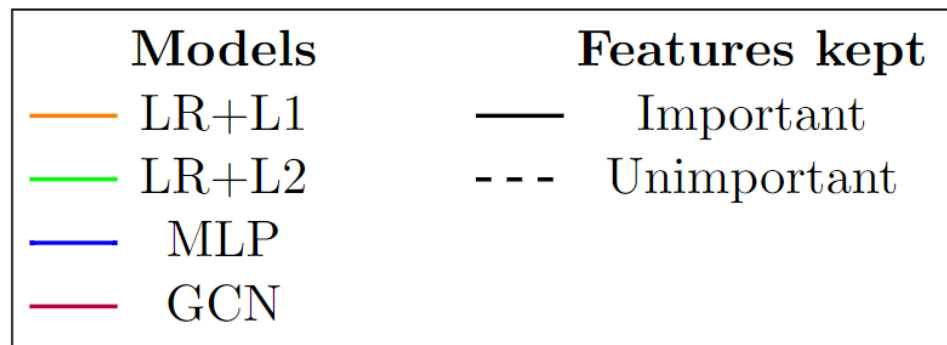


BRCA-pam (PAM50 breast cancer subtypes)

# Application



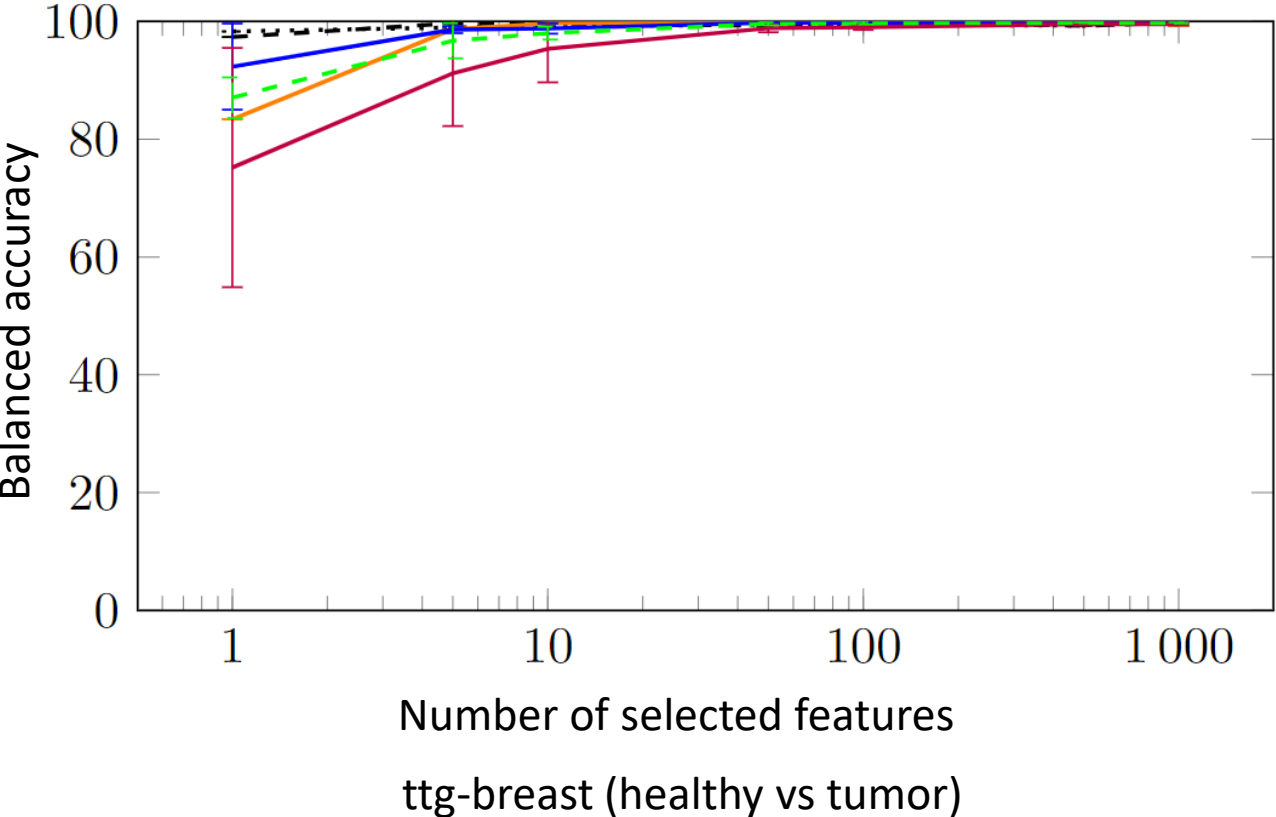
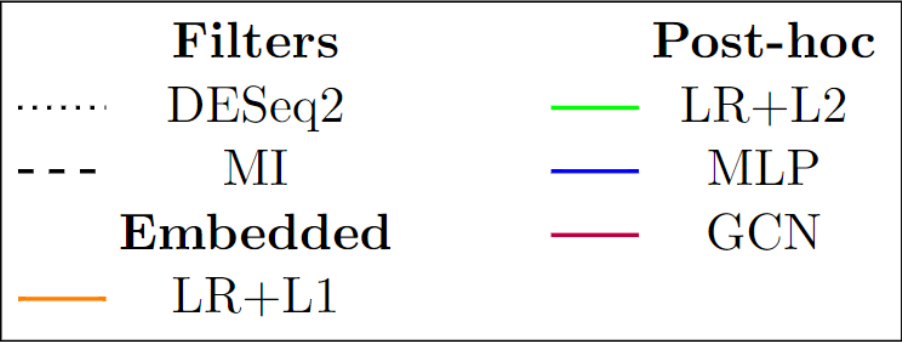
# Application



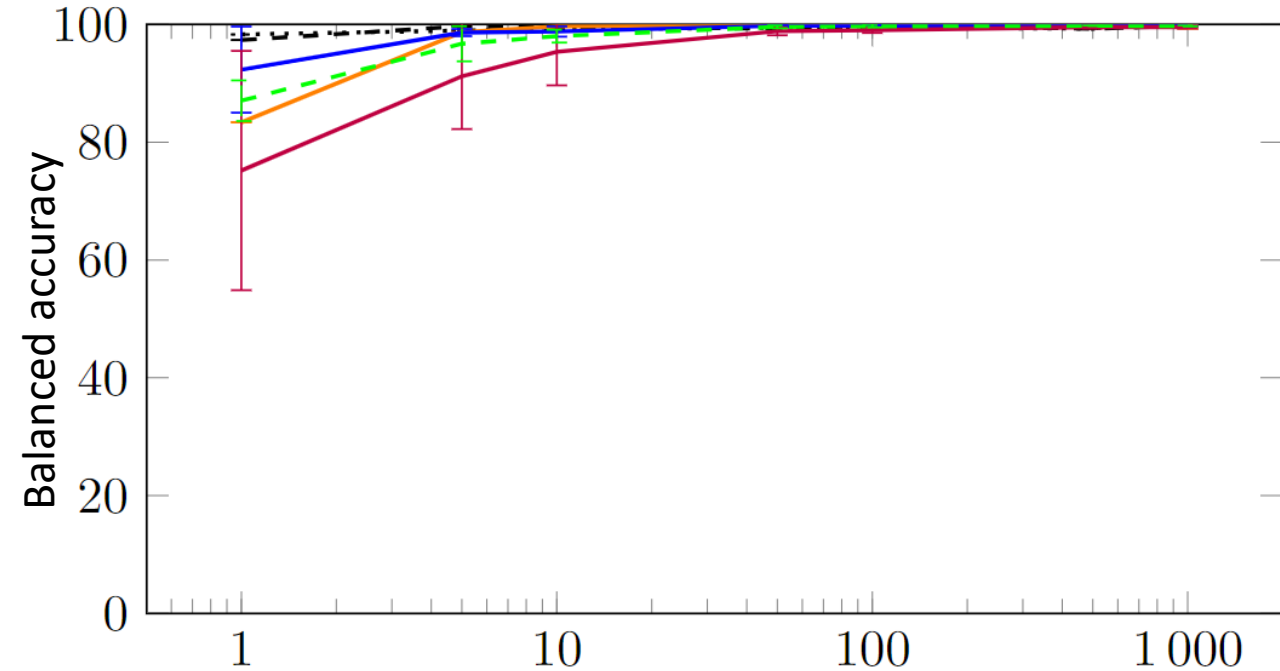
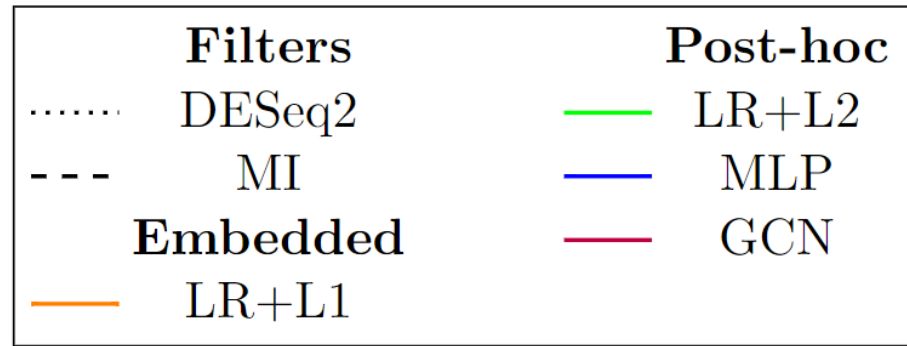
Number of selected features  
ttg-breast (healthy vs tumor)

Number of selected features  
BRCA-pam (PAM50 breast cancer subtypes)

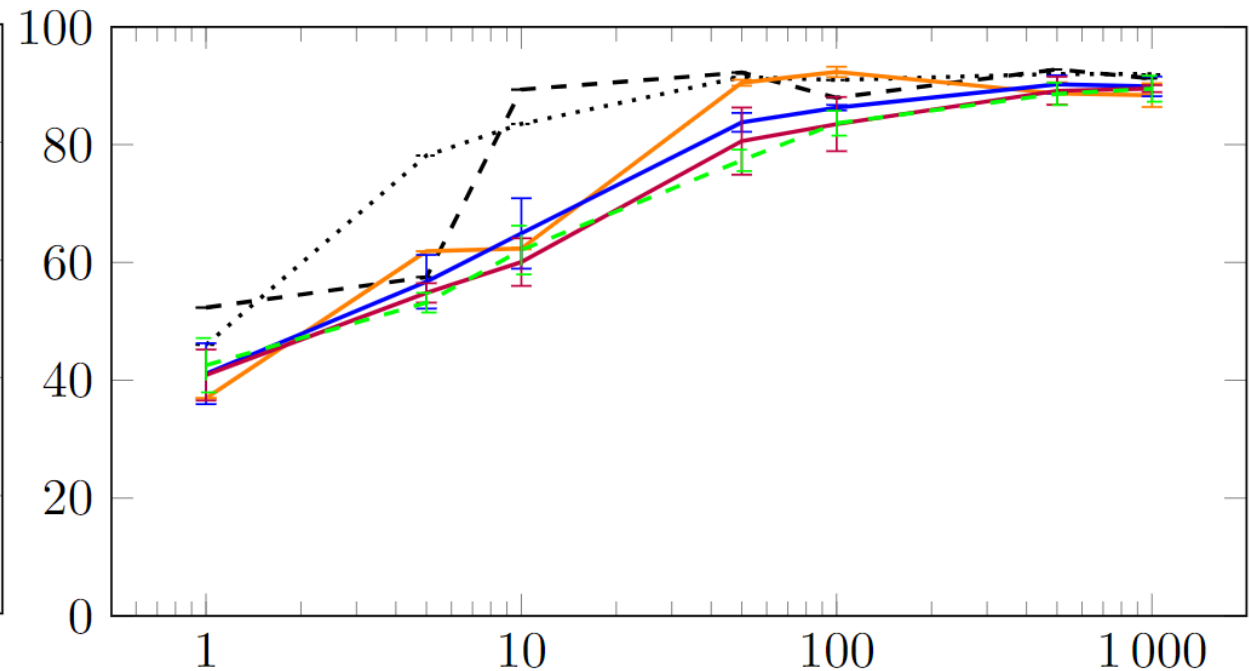
# Application



# Application



Number of selected features  
ttg-breast (healthy vs tumor)



Number of selected features  
BRCA-pam (PAM50 breast cancer subtypes)

# Conclusion

## Summary

- Challenges of functional genomics
- Machine learning methods selecting relevant genes
- Discussion on 2 datasets
  - ❖ Good classification performance, even with linear models.
  - ❖ Small sets containing the best-ranked genes are sufficient to achieve a good classification.
  - ❖ Sets generated by the different methods differ significantly.
  - ❖ Similar classification performance can be achieved with lower ranked genes.



# Conclusion

## Summary

- Challenges of functional genomics
- Machine learning methods selecting relevant genes
- Discussion on 2 datasets
  - ❖ Good classification performance, even with linear models.
  - ❖ Small sets containing the best-ranked genes are sufficient to achieve a good classification.
  - ❖ Sets generated by the different methods differ significantly.
  - ❖ Similar classification performance can be achieved with lower ranked genes.

## Next steps

- Extension to other classification problems
- Extension to multiomics datasets
- Automatically extract explanations at a relevant scale