

How to build the metadata files

ODAM



EDTMS

Tutorial

Daniel Jacob

UMR 1332 BFP – Metabolism Group

Bordeaux Metabolomics Facility

May 2016

Explanation based on an example



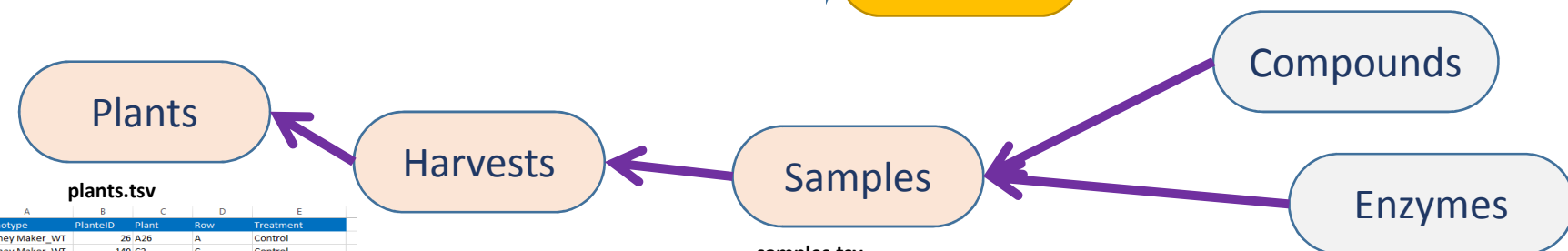
FRIM - Fruit Integrative Modelling

See <http://www.erasysbio.net/index.php?index=266> & <http://frim.brookes.ac.uk/Home>

Design of Experiment (DoE)

- 1 Genotype **MoneyMaker**, Species **Lyco. Sola**.
- ~600 plants, >1500 samples
- Tissues : **Fruit** (F)
- Treatment : **Control** (TC), **Shadow** (TS), **Water Shortage** (TWS)
- bouquet, weight, height, diameter,
- Dev. Stages, 9 levels: from march up to august
- Compounds: 12 quantified metabolites
- Enzymes: 38 quantified enzymes

- **5 object types of study**, namely: plants, harvests, samples, compounds and enzymes
- **2 factors**: Treatment, Development stages
- **53 quantitative variables**: compounds (12) + enzymes (38) + weight, height, diameter(3)



68	enzymes.tsv				
	DH	ND	Eno5	ISDH	NADiPEPC
	50.99	NA	NA		
	975.77	514.4	133.6		
	355.48	388.57	167.88		
	72.23	676.61	87.43		
	421.33	504.23	231.52		
	416.04	530.9	59.01		
	606.68	645.49	221.32		
	965.12	530.51	68.32		
	437.56	548.96	134.09		
	603.02	557.28	163.25		
	639.72	579.25	NA		
	69.02	518.03	91.02		
	210.84	581.39	118.38		
	247.36	587.58	20.21		
	190.41	545.06	58.85		
	301.7	NA	168.46		
	204.19	500.82	88.22		
	366.29	639.19	45.01		1
	256.53	766.52	NA		
	202.99	600.66	41.36		
	146.66	660.62	121.36		

Data subsets files must be compliant with the **TSV standard**
(Tab-Separator-Values)



FRIM - Preparation and cleaning of the data sub-sets of files

1	Genotype	PlantID	Plant	Row	Treatment
2	Money Maker_WT	26	A26	A	Control
3	Money Maker_WT	140	C2	C	Control
4	Money Maker_WT	222	D15	D	Control
5	Money Maker_WT	295	E19	E	Control
6	Money Maker_WT	310	E24	E	Control

plants.tsv

1	PlantID	Lot	Truss	FruitAge	HarvestDate	HarvestHour	FruitPosition	FruitFW	FruitDiameter	FruitHeight
8	18	5	Truss_5	00.08DPA	40379	0.51875	4	1.1	13.29	13.17
9	155	5	Truss_5	00.08DPA	40379	0.51875	4	1.01	12.38	12.29
10	164	5	Truss_5	00.08DPA	40379	0.51875	4	1.02	12.28	12.44
11	221	5	Truss_5	00.08DPA	40379	0.51875	2	1.02	12.8	12.28
12	226	5	Truss_5	00.08DPA	40379	0.51875	3	0.81		
13	322	5	Truss_5	00.08DPA	40379	0.51875	2	0.79	1	
14	343	5	Truss_5	00.08DPA	40379	0.51875	4	0.69	2	
15	361	5	Truss_5	00.08DPA	40379	0.51875	5	1.06	3	
16	369	5	Truss_5	00.08DPA	40379	0.51875	4	0.73	4	
17	372	5	Truss_5	00.08DPA	40379	0.51875	5	1.21	5	

harvests.tsv

1	Lot	SampleID	NbFruit	GellyFW	GellyFruit	BER
1	1	1	10	2.83	0.283	FALSE
2	1	2	10	2.83	0.283	FALSE
3	2	3	8	4.05	0.50625	FALSE
4	2	3	10	3.3	0.33	FALSE
5	3	5	10	3.3	0.33	FALSE

samples.tsv

Data subset files

compounds.tsv

- Whatever the kind of experiment, this assumes a design of experiment (DoE) involving individuals, samples or whatever things, as the main objects of study (e.g. plants, tissues, bacteria, ...)
- This also assumes the observation of dependent variables resulting of effects of some controlled experimental **factors**.
- Moreover, the objects of study have usually an **identifier** for each of them, and the variables can be **quantitative** or **qualitative**.
- We can have either one object type of study or several kinds, but in this latter case, it must exist a relationship between object types that we assume of "**obtainedFrom**" type.

enzymes.tsv

L	C	M	K	O
SODH_NAD	EnoS	ISODH_NADI	PEPC	
508.99	NA	NA		93
975.77	514.4		133.6	11
555.54	588.57		167.86	99
572.33	676.61		87.43	27
421.23	504.23		231.52	81
416.04	530.9		59.01	8
660.68	645.49		221.32	68
965.12	530.51		68.32	9
437.56	548.96		194.09	84
607.02	557.28		163.25	1
632.42	579.25	NA		90
659.72	518.63		91.02	87
210.84	581.39		118.38	61
247.36	587.58		20.21	86
190.41	545.06		58.85	57
301.7	NA		168.46	83
204.19	500.82		88.22	
366.29	639.19		45.01	127
256.53	766.52	NA		6
202.93	600.66		41.36	9
146.66	660.62		124.38	58



FRIM - Classification of each column within its right category

Data subset files

Genotype	PlantID	Plant	Row	Treatment
Money Maker_WT	26	A26	A	Control
Money Maker_WT	140	C2	C	Control
Money Maker_WT	222	D15	D	Control

plants.tsv

PlantID	Lot	Truss	FruitAge	HarvestDate	HarvestHour	FruitPositior	FruitFW	FruitDiameter	FruitHeight
18	5	Truss_5	00.08DPA	40379	0.51875	4	1.1	13.29	13.17
155	5	Truss_5	00.08DPA	40379	0.51875	4	1.01	12.38	12.29
164	5	Truss_5	00.08DPA						
221	5	Truss_5	00.08DPA						
226	5	Truss_5	00.08DPA						
322	5	Truss_5	00.08DPA						

harvests.tsv

Lot	SampleID	NbFruit	GellyFW	GellyFruit	BER
1	1	10	2.83	0.283	FALSE
1	1	10	2.83	0.283	FALSE
2	2	10	2.83	0.283	FALSE

samples.tsv

SampleID	DPA	MassBefore	MassMIA	RDT	Starch1	Starch2	RHAMNOSE	ARABINOSE	XYLOSE	MANNONE	GALACTOSE	GLUCOSE	OsesN	PoidsAU	Poids
363	15	0.192	0.0002	52.19	36.55	0.19074531	0.89	2.17	1.73	1.2	5.99	11.31	23.29	7.17	
369	15	0.107	0.0511	47.76	32.43	0.15487598	0.74	2.19	1.68	1.3	6.03	8.51	20.45	7.35	
373	15	0.166	0.0771	46.45	32.05	0.14885874	0.93	2.44	2.56	2.95	9.3	18.18	36.36	10.88	
375	15	0.14	0.0684	48.86	34.89	0.17046257	0.71	1.89	1.98	2.04	7.33	10.47	24.42	8.79	
429	28	0.104	0.0431	41.44	19.78	0.08197289	0.79	2.13	2.33	2.42	8.51	15.2	31.38	9.94	

compounds.tsv

SampleID	PGM	F16BP_Cyt	PyrK	CitS	PPI	AcoS	PFK	FruS	F16BP_Stron	Glus	ISODH_NAD	EnoS	ISODH_NADI	PEPC
1	NA	8.97	1599.53	64.53	2767.89	1172.28	192.05	876.13	523.57	722.19	508.99	NA	NA	936.41
3	6844.78	85.02	1839.39	NA	3373.46	2014.02	263.33	984.08	634.14	622.22	975.77	514.4	133.6	1107.8

enzymes.tsv

SampleID	DPA	MassBefore	MassMIA	RDT	Starch1	Starch2	RHAMNOSE	ARABINOSE	XYLOSE	MANNONE	GALACTOSE	GLUCOSE	OsesN	PoidsAU	Poids
363	15	0.192	0.0002	52.19	36.55	0.19074531	0.89	2.17	1.73	1.2	5.99	11.31	23.29	7.17	
369	15	0.107	0.0511	47.76	32.43	0.15487598	0.74	2.19	1.68	1.3	6.03	8.51	20.45	7.35	
373	15	0.166	0.0771	46.45	32.05	0.14885874	0.93	2.44	2.56	2.95	9.3	18.18	36.36	10.88	
375	15	0.14	0.0684	48.86	34.89	0.17046257	0.71	1.89	1.98	2.04	7.33	10.47	24.42	8.79	
429	28	0.104	0.0431	41.44	19.78	0.08197289	0.79	2.13	2.33	2.42	8.51	15.2	31.38	9.94	

categories

- link
- identifier
- factor
- qualitative
- quantitative

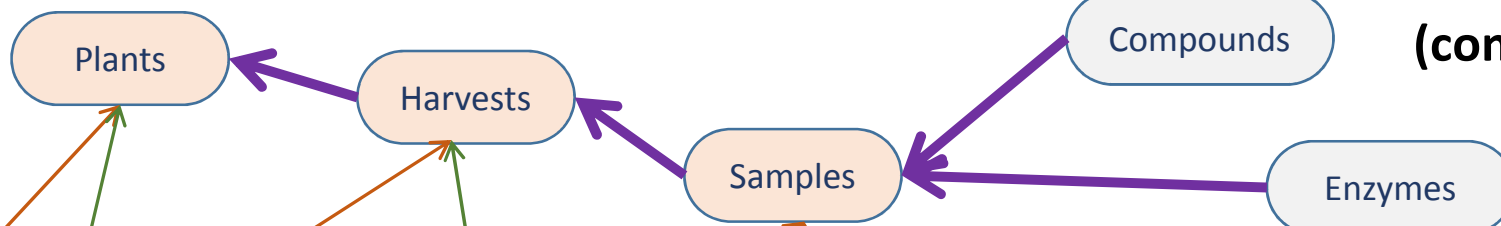
You have to organize your data subsets so that links could be established between them. In practical, it means to add a column containing the identifiers corresponding to the entity to which you want to connect the subset. It is to be noted that this duplication of identifiers must be the only redundant information, through all data subsets.

Money Maker_WT	24	1169	15	363	4753.73	17.77	728.5	NA	1869.46	450.13	171.89	593.3	357.74		
Money Maker_WT	25	1187	16	365	3836.19	4.16	603.74	31.55	1836.44	330.26	131.87	499.38	208.85		
	26	1343	17	367	4.11	4.87	507.67	7.39	1371.13	573.06	128.65	390.78	283.25		
	27	1345	18	369	3798.59	15.14	703.9	82.78	1473.25	286.02	120.11	438.09	260.44		
	28	1347	19	371	5610.69	14.69	1092.42	99.09	2384.35	799.55	182.18	812.09	425.56		
	29	1349	20	373	4220.68	31.7	592.57	14.81	1654.27	602.62	150.23	480.9	283.92		
	30	1351													
	31	1353													
	32	1355	8	0.093	0.057	61.2903226	29.63	0.18160323	0.64	2.25	2.57	3.2	3.65	15.73	28.04



FRIM - Connections between the dataset files based on identifiers

Entities
(concepts)



plants.tsv

harvests.tsv

samples.tsv

enzymes.tsv

compounds.tsv

categories

link

identifier

factor

qualitative

quantitative

Identifier of the central entity of the subset

Data subset files



FRIM - Creation of the metadata files

Supplementary files

In order to allow data to be explored and mined, we have to adjoin some minimal but relevant metadata:

For that, **2 metadata files** are required

- **s_subsets.tsv**: a file allowing to associate with each subset of data a key concept corresponding to the main entity of the subset and the relations of the type "obtainedFrom" between these concepts
- **a_attributes.tsv**: a metadata file allowing each attribute (concept/variable) to be annotated with some minimal but relevant metadata

Note: Data subsets files and their associated metadata files must be compliant with the **TSV standard** (Tab-Separator-Values)

TSV is an alternative to the common comma-separated values (CSV) format, which often causes difficulties because of the need to escape commas

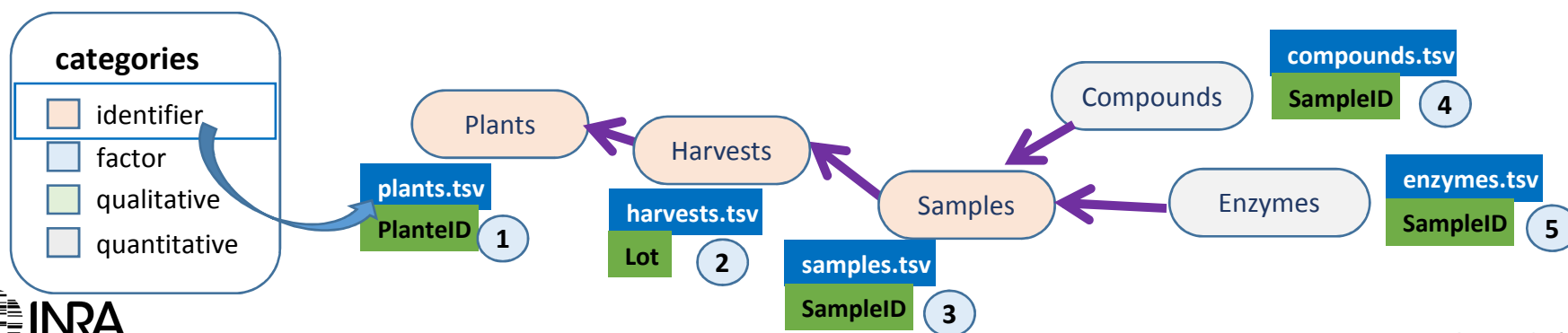
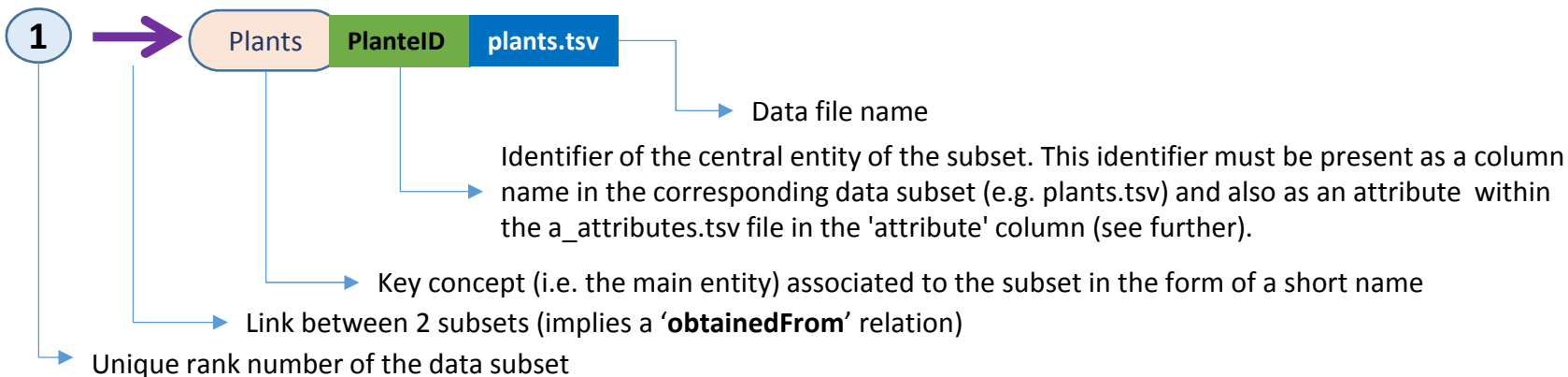


FRIM - Creation of the metadata files

s_subsets.tsv

This metadata file allows to associate a key concept to each **data subset file**

	A	B	C	D	E	F	G	H
1	rank	obtainedFrom	subset	identifier	file	description	CV_term_id	CV_term_name
2	1	0	plants	PlantID	plants.tsv	Plant features	http://purl.obolibrary.org/obo/PO_0000003	whole plant
3	2	1	harvests	Lot	harvests.tsv	Harvest features	http://purl.obolibrary.org/obo/OBI_1110046	organ harvesting
4	3	2	samples	SampleID	samples.tsv	Samples features	http://purl.obolibrary.org/obo/PO_0009001	fruit
5	4	3	compounds	SampleID	compounds.tsv	Compound quantifications	http://purl.obolibrary.org/obo/CHEBI_24431	chemical entity
6	5	3	enzymes	SampleID	enzymes.tsv	Enzyme Features	http://purl.obolibrary.org/obo/OBI_0000427	enzyme





FRIM - Creation of the metadata files

s_subsets.tsv

This metadata file allows to associate a key concept to each **data subset file**

	A	B	C	D	E	F	G	H
1	rank	obtainedFrom	subset	identifier	file	description	CV_term_id	CV_term_name
2	1	0 plants	PlantID	plants.tsv	Plant features	http://purl.obolibrary.org/obo/PO_0000003	whole plant	
3	2	1 harvests	Lot	harvests.tsv	Harvest features	http://purl.obolibrary.org/obo/OBI_1110046	organ harvesting	
4	3	2 samples	SampleID	samples.tsv	Samples features	http://purl.obolibrary.org/obo/PO_0009001	fruit	
5	4	3 compounds	SampleID	compounds.tsv	Compound quantifications	http://purl.obolibrary.org/obo/CHEBI_24431	chemical entity	
6	5	3 enzymes	SampleID	enzymes.tsv	Enzyme Features	http://purl.obolibrary.org/obo/OBI_0000427	enzyme	

a description for each subset

*Optional:
an annotation based on ontology*



FRIM - Creation of the metadata files : Be careful in the spelling

s_subsets.tsv

This metadata file allows to associate a key concept to each **data subset file**

	A	B	C	D	E	F	G	H
1	rank	obtainedFrom	subset	identifier	file	description	CV_term_id	CV_term_name
2	1	0 plants	PlantID	plants.csv	Plant features	http://purl.obolibrary.org/obo/PO_0000003	whole plant	
3	2	1 harvests	Lot	harvests.csv	Harvest features	http://purl.obolibrary.org/obo/OBI_1110046	organ harvesting	
4	3	2 samples	SampleID	samples.csv	Samples features	http://purl.obolibrary.org/obo/PO_0009001	fruit	
5	4	3 compounds	SampleID	compounds.csv	Compound quantifications	http://purl.obolibrary.org/obo/CHEBI_24431	chemical entity	
6	5	3 enzymes	SampleID	enzymes.csv	Enzyme Features	http://purl.obolibrary.org/obo/OBI_0000427	enzyme	

1. For short names of the data subsets (**C column**), the name of the identifier attributes (**D column**) and the names of the files (**E column**),
 - only the alphanumerical characters and the underscore are allowed (i.e. 'a-z', 'A-Z', '0-9' and '_').
 - Moreover, these names should not start with a digit!
2. For description (**F column**),
 - the allowed characters are : 0-9 a-z A-Z , : + * () [] { } - % ! | / . ?

The identifier attributes (**D column**)

1. should be the only attribute declared as '*identifier*' in the '**category**' column in the a_attributes.tsv file (**D column**)
2. should be available as a column item in the corresponding data subset file



FRIM - Creation of the metadata files

a_attributes.tsv

This metadata file allows each **attribute (variable)** to be annotated with some minimal but relevant metadata

	A	B	C	D	E	F	G	H
	subset	attribute	entry	category	type	Description	CV term id	CV term name
Plants	plants	PlantID	plantid	identifier	numeric	Plant identifier	http://purl.obolibrary.org	individual organism identifier
	plants	Row	row	qualitative	string	Row of the individual plant on the table	http://ncicb.nci.nih.gov	Row
	plants	Plant	plant		string	Code identifier of		
	plants	Treatment	treatment	factor	string	Treatment app		
	plants	Genotype		qualitative	string	Genotype		
Harvests	harvests	Lot	lot	identifier	numeric	Pool of several		
	harvests	PlantID			numeric	Plant identifier		
	harvests	Truss		qualitative	string	Position on the		
	harvests	HarvestDate			string	Harvest date		
	harvests	HarvestHour			string	Harvest hour		
	harvests	FruitAge	age	factor	string	Fruit developm		
	harvests	FruitPosition		qualitative	numeric	Position on the		
	harvests	FruitDiameter		quantitative	numeric	Fruit diameter		
Samples	harvests	FruitHeight		quantitative	numeric	Fruit height (m		
	harvests	FruitFW		quantitative	numeric	Fruit Fresh We		
	samples	SampleID	sampleid	identifier	numeric	Sample identifi		
	samples	Lot			numeric	Pool of several ha		
	samples	NbFruit			numeric	Fruit Number per sample		
	samples	GellyFW		quantitative	numeric			
Compounds	samples	GellyFruit			numeric			
	samples	BER			string			
	compounds	SampleID			numeric		http://purl.obolibrary.org	centrally registered identifier
	compounds	DPA			numeric			
	compounds	MassBefore			numeric			
	compounds	MassMIA			numeric			
	compounds	RDT			numeric			
	compounds	Starch1			numeric		http://purl.obolibrary.org	starch
	compounds	Starch2			numeric		http://purl.obolibrary.org	starch
	compounds	RHAMNOSE			numeric		http://purl.obolibrary.org	rhamnose

1. Short names of the data subsets (**A column**) must be declared in the s_subsets.tsv file (**C column**) and vice versa.
2. One and only one attribute (**B column**) must be declared as '*identifier*' in the '**category**' column (**D column**) per data subset (**A column**)

the
attribute
names

the
attribute
data types



FRIM - Creation of the metadata files

a_attributes.tsv

This metadata file allows each **attribute (variable)** to be annotated with some minimal but relevant metadata

plants.tsv

Plants

Harvests

Samples

Compounds

...

Transpose

Get attribute names simply from the data subsets

subset	attribute	entry	category
plants	PlantID	plantid	identifier
plants	Row	row	quality
plants	Plant	plant	quality
plants	Treatment	treatment	factor
plants	Genotype		quality
harvests	Lot	lot	identifier
harvests	PlantID		
harvests	Truss		quality
harvests	HarvestDate		quality
harvests	HarvestHour		quality
harvests	FruitAge	age	factor
harvests	Fruit		quality
harvests	FruitFW		quality
harvests	FruitFW		quality
harvests	FruitFW		quality
samples	SampleID	sampleid	identifier
samples	Lot		
samples	NbFruit		quality
samples	GellyFW		quality
samples	GellyFruit		quality
samples	BER		quality
compounds	SampleID	sampleid	identifier
compounds	DPA		factor
compounds	MassBefore		quality
compounds	MassMIA		quality
compounds	RDT		quality
compounds	Starch1		quality
compounds	Starch2		quality
compounds	RHAMNOSE		quality

Genotype	PlantID	Plant	Row	Treatment
Money Make	26	A26	A	Control
Money Make	140	C2	C	Control
Money Make	222	D15	D	Control
Money Make	295	E19	E	Control
Money Make	310	E34	E	Control
Money Make	314	E38	E	Control
Money Make	374	H29	H	Control
Money Make	379	H34	H	Control
Money Make	397	H52	H	Control
Money Make	406	H61	H	Control
Money Make	415	F1	F	WATER STRESS
Money Make	438	F24	F	WATER STRESS
Money Make	460	F46	F	WATER STRESS
Money Make	486	G3	G	WATER STRESS
Money Make	494	G11	G	WATER STRESS
Money Make	511	G28	G	WATER STRESS



FRIM - Creation of the metadata files

a_attributes.tsv

This metadata file allows each **attribute (variable)** to be annotated with some minimal but relevant metadata

harvests.tsv

Plants

Harvests

Samples

Compounds

subset	attribute
plants	PlantID
plants	Row
plants	Plant
plants	Treatment
plants	Genotype
harvests	Lot
harvests	PlantID
harvests	Truss
harvests	HarvestDate
harvests	HarvestHour
harvests	FruitAge
harvests	FruitPosition
harvests	FruitDiameter
harvests	FruitHeight
harvests	FruitFW
samples	SampleID
samples	Lot
samples	NbFruit
samples	GellyFW
samples	GellyFruit
samples	BER
compounds	SampleID
compounds	DPA
compounds	MassBefore
compounds	MassMIA
compounds	RDT
compounds	Starch1
compounds	Starch2
compounds	RHAMNOSE

harvests.csv - Excel

FICHIER ACCUEIL INSERTION MISE EN PAGE FORMULES DONNÉES RÉVISION AFFICHAGE PDF Architect Connexion

A1 : X ✓ fx PlantID

PlantID	Lot	Truss	FruitAge	HarvestDate	HarvestHour	FruitPosition	FruitFW	FruitDiameter	FruitHeight
30	18	5 Truss_5	00.08DPA	40379	0.51875	4	1.1	13.29	13.17
31	155	5 Truss_5	00.08DPA	40379	0.51875	4	1.01	12.38	12.29
32	164	5 Truss_5	00.08DPA	40379	0.51875	4	1.02	12.28	12.44
33	221	5 Truss_5	00.08DPA	40379	0.51875	2	1.02	12.8	12.28
34	226	5 Truss_5	00.08DPA	40379	0.51875	3	0.81	12.27	11.43
35	322	5 Truss_5	00.08DPA	40379					10.46
36	343	5 Truss_5	00.08DPA	40379					10.09
37	351	5 Truss_5	00.08DPA	40379					10.58
38	372	5 Truss_5	00.08DPA	40379					10.72
39	372	5 Truss_5	00.08DPA	40379					10.66
40	448	6 Truss_5	00.08DPA	40379	0.52430556	5	0.66	10.99	10.82
41	476	6 Truss_5	00.08DPA	40379	0.52430556	5	0.66	10.99	10.18
42	490	6 Truss_5	00.08DPA	40379	0.52430556	5	0.62	10.72	10.79
43	507	6 Truss_5	00.08DPA	40379	0.52430556	5	0.75	12.14	10.77
44	512	6 Truss_5	00.08DPA	40379	0.52430556	2	0.58	10.92	9.92
45	522	6 Truss_5	00.08DPA	40379	0.52430556	5	0.86	11.88	11.52

Transpose

Get attribute names simply from the data subsets



FRIM - Creation of the metadata files

a_attributes.tsv

This metadata file allows each **attribute (variable)** to be annotated with some minimal but relevant metadata

	A	B	C	D	E	F	G	H
	subset	attribute	entry	category	type	description	CV_term_id	CV_term_name
Plants	plants	PlantID	plantid	identifier	numeric	Plant identifier	http://purl.obolibrary.org	individual organism identifier
	plants	Row	row	qualitative	string	Row of the invidual plant on the table	http://ncicb.nci.nih.gov/	Row
	plants	Plant	plant		s			
	plants	Treatment	treatment	factor				
	plants	Genotype		qualitative	s			
	plants	Lot	lot	identifier	n			
Harvests	harvests	PlantID			n			
	harvests	Truss		qualitative	s			
	harvests	HarvestDate			s			
	harvests	HarvestHour			s			
	harvests	FruitAge	age	factor	s			
	harvests	FruitPosition		qualitative	n			
	harvests	FruitDiameter		quantitative	n			
	harvests	FruitHeight		quantitative	n			
Samples	samples	FruitFW		quantitative	n			
	samples	SampleID	sampleid	identifier	n	sample identifier	https://purl.obolibrary.org/	centrally registered identifier
	samples	Lot			numeric	Pool of several harvests	http://www.ebi.ac.uk/efo/	sample pooling
	samples	NbFruit			numeric	Fruit Number per sample		
	samples	GellyFW		quantitative	numeric	Gelly Fred Weight		
	samples	GellyFruit		quantitative	numeric	Gelly per Fruit		
Compounds	samples	BER			string	BER		
	compounds	SampleID	sampleid	identifier	numeric	Sample identifier	http://purl.obolibrary.org/	centrally registered identifier
	compounds	DPA		factor	numeric	Day Per Anthesis		
	compounds	MassBefore		quantitative	numeric	m av.extracted		
	compounds	MassMIA		quantitative	numeric	masse MIA (g)		
	compounds	RDT		quantitative	numeric	Rdt (% MIA/DV)		
	compounds	Starch1		quantitative	numeric	Dosage amido	http://purl.obolibrary.org/	starch
	compounds	Starch2		quantitative	numeric	amidon (g/gD)	http://purl.obolibrary.org/	starch
	compounds	RHAMNOSE		quantitative	numeric	RHAMNOSE	http://purl.obolibrary.org/	rhannose

Plants

Harvests

Samples

Compounds

Dependent variables resulting of effects of some controlled experimental **factors**.

Objects of study have usually an **identifier** for each of them, and the variables can be **quantitative** or **qualitative**.

categories

identifier

factor

qualitative

quantitative



FRIM - Creation of the metadata files

a_attributes.tsv

This metadata file allows each **attribute (variable)** to be annotated with some minimal but relevant metadata

	A	B	C	D	E	F	G	H
	subset	attribute	entry	category	type	description	CV_term_id	CV_term_name
Plants	plants	PlantID	plantid	identifier	numeric	Plant identifier	http://purl.obolibrary.org	individual organism identifier
	plants	Row	row	qualitative	string	Row of the individual plant on the table	http://ncicb.nci.nih.gov/	Row
	plants	Plant	plant		string	Code identifier of the individual plant	http://ncicb.nci.nih.gov/	Discrete Set Coded String Data Type
	plants	Treatment	treatment	factor	string	Treatment applied on plants	http://www.ebi.ac.uk/ef/	environmental factor
	plants	Genotype		qualitative	string	Genotype		
Harvests	harvests	Lot	lot	identifier	numeric	Pool of several harvests	http://www.ebi.ac.uk/ef/	sample pooling
	harvests	PlantID			numeric	Plant identifier	http://purl.obolibrary.org	individual organism identifier
	harvests	Truss		qualitative	string	Position on the stem of the truss	http://purl.obolibrary.org	stem node
	harvests	HarvestDate			string	Harvest date		
	harvests	HarvestHour						
	harvests	FruitAge	age	factor			http://purl.obolibrary.org	fruit development stage
	harvests	FruitPosition		qualitative			http://ncicb.nci.nih.gov/	Position Number
	harvests	FruitDiameter		quantitative			http://ncicb.nci.nih.gov/	Diameter
Samples	samples	SampleID	sampleid	identifier			http://ncicb.nci.nih.gov/	Height
	samples	Lot					http://ncicb.nci.nih.gov/	Weight
	samples	NbFruit		quantitative			http://purl.obolibrary.org	centrally registered identifier
	samples	GellyFW		quantitative			http://www.ebi.ac.uk/ef/	sample pooling
	samples	GellyFruit		quantitative				
Compounds	compounds	BER						
	compounds	SampleID	sampleid	identifier			http://purl.obolibrary.org	centrally registered identifier
	compounds	DPA		factor				
	compounds	MassBefore		quantitative				
	compounds	MassMIA		quantitative				
	compounds	RDT		quantitative				
	compounds	Starch1		quantitative			http://purl.obolibrary.org	starch
	compounds	Starch2		quantitative	numeric	amidon (g/gDW)	http://purl.obolibrary.org	starch
	compounds	RHAMNOSE		quantitative	numeric	RHAMNOSE	http://purl.obolibrary.org	rhamnose

Give opportunity to make a selection on the attributes via web-services by associating them an alias name (called an "entry")

[http://myhost.org/getdata/xml/frim1/\(samples\)/treatment/Control](http://myhost.org/getdata/xml/frim1/(samples)/treatment/Control)



FRIM - Creation of the metadata files

a_attributes.tsv

This metadata file allows each **attribute (variable)** to be annotated with some minimal but relevant metadata

	A	B	C	D	E	F	G	H
	subset	attribute	entry	category	type	description	CV_term_id	CV_term_name
Plants	plants	PlantID	plantid	identifier	numeric	Plant identifier	http://purl.obolibrary.org	individual organism identifier
	plants	Row	row	qualitative	string	Row of the invidual plant on the table	http://ncicb.nci.nih.gov/	Row
	plants	Plant	plant		string	Code identifier of the individual plant	http://ncicb.nci.nih.gov/	Discrete Set Coded String Data Type
	plants	Treatment	treatment	factor	string	Treatment applied on plants	http://www.ebi.ac.uk/efo	environmental factor
	plants	Genotype		qualitative	string	Genotype		
	plants	Lot	lot	identifier	numeric	Pool of several harvests	http://www.ebi.ac.uk/efo	sample pooling
Harvests	harvests	PlantID			numeric	Plant identifier	http://purl.obolibrary.org	individual organism identifier
	harvests	Truss		qualitative	string	Position on the stem of the truss	http://purl.obolibrary.org	stem node
	harvests	HarvestDate			string	Harvest date		
	harvests	HarvestHour			string	Harvest hour		
	harvests	FruitAge	age	factor	string	fruit development stage	http://purl.obolibrary.org	fruit development stage
	harvests	FruitPosition		qualitative	numeric	Poistion on the truss of the fruit	http://ncicb.nci.nih.gov/	Position Number
	harvests	FruitDiameter		quantitative	numeric	Fruit diameter (mm)	http://ncicb.nci.nih.gov/	Diameter
	harvests	FruitHeight		quantitative	numeric	Fruit height (mm)	http://ncicb.nci.nih.gov/	Height
	harvests	FruitFW		quantitative	numeric	Fruit Fresh Weight(g)	http://ncicb.nci.nih.gov/	Weight
	samples	SampleID	sampleid	identifier	numeric	Sample identifier	http://purl.obolibrary.org	centrally registered identifier
Samples	samples	Lot			numeric	Pool of several harvests	http://www.ebi.ac.uk/efo	sample pooling
	samples	NbFruit			numeric	Fruit Number per sample		
	samples	GellyFW		quantitative	numeric	Gelly Fred Weight(g) per sample		
	samples	GellyFruit						
	samples	BER						
Compounds	compounds	SampleID	sampleid	ic			h	
	compounds	DPA		fa				
	compounds	MassBefore		q				
	compounds	MassMIA		q				
	compounds	RDT		q				
	compounds	Starch1		q			h	
	compounds	Starch2		q			h	
	compounds	RHAMNOSE		q			h	

Plants

Harvests

Samples

Compounds

a description for each attribute

Optional:
an annotation based on ontology

a description for each attribute

Optional:
an annotation based on ontology



FRIM - Creation of the metadata files : Be careful in the spelling

a_attributes.tsv

This metadata file allows each **attribute (variable)** to be annotated with some minimal but relevant metadata

	A	B	C	D	E	F	G	H
1	subset	attribute	entry	category	type	description	CV_term_id	CV_term_name
2	plants	PlantID	plantid	identifier	numeric	Plant identifier	http://purl.obolibrary.org	individual organism identifier
3	plants	Row	row	qualitative	string	Row of the invidual plant on the table	http://ncicb.nci.nih.gov/	Row
4	plants	Plant	plant		string	Code identifier of the individual plant	http://ncicb.nci.nih.gov/	Discrete Set Coded String Data Type
5	plants	Treatment	treatment	factor	string	Treatment applied on plants	http://www.ebi.ac.uk/ef/	environmental factor
6	plants	Genotype		qualitative	string	Genotype		
7	harvests	Lot	lot	identifier	numeric	Pool of several harvests	http://www.ebi.ac.uk/ef/	sample pooling
8	harvests	PlantID			numeric	Plant identifier	http://purl.obolibrary.org	individual organism identifier
9	ha							
10								

- For short names of the data subsets (**A column**), the name of the attributes (**B column**) and the entry names (**C column**):
 - only the alphanumeric characters and the underscore are allowed (i.e. 'a-z', 'A-Z', '0-9' and '_').
 - Moreover, these names should not start with a digit!
- For categorical names (**D column**):
 - the set of terms are fixed, namely: '*identifier*', '*factor*', '*quantitative*', '*qualitative*'. Leave as blank otherwise.
- For data types (**E column**):
 - the allowed names are restricted to '*numeric*' or '*string*'.
- For description (**F column**):
 - the allowed characters are : 0-9 a-z A-Z , : + * () [] { } - % ! | / . ?



FRIM - Updating the metadata files

s_subsets.tsv

	A	B	C	D	E	F	G	H
1	rank	obtainedFrom	subset	identifier	file	description	CV_term_id	CV_term_name
2	1		0 plants	PlantID	plants.tsv	Plant features	http://purl.obolibrary.org/obo/PO_0000003	whole plant
3	2		1 harvests	Lot	harvests.tsv	Harvest features	http://purl.obolibrary.org/obo/OBI_1110046	organ harvesting
4	3		2 samples	SampleID	samples.tsv	Samples features	http://purl.obolibrary.org/obo/PO_0009001	fruit
5	4		3 compounds	SampleID	compounds.tsv	Compound quantifications	http://purl.obolibrary.org/obo/CHEBI_24431	chemical entity
6	5		3 enzymes	SampleID	enzymes.tsv	Enzyme Features	http://purl.obolibrary.org/obo/OBI_0000427	enzyme
	...							

a_attributes.tsv

	A	B	C	D	E	F
1	subset	attribute	entry	category	type	description
2	plants	PlantID	plantid	identifier	numeric	Plant identifier
3	plants	Row	row	qualitative	string	Row of the individual plant or
4	plants	Plant	plant		string	Code identifier of the indiv
5	plants	Treatment	treatment	factor	string	Treatment applied on plant
6	plants	Genotype		qualitative	string	Genotype
7	harvests	Lot	lot	identifier	numeric	Pool of several harvests
8	harvests	PlantID			numeric	Plant identifier
	...					
24	compounds	DPA		factor	numeric	Day Per Anthesis
25	compounds	MassBefore		quantitative	numeric	m av.extraction (g)
26	compounds	MassMIA		quantitative	numeric	masse MIA (g)
27	compounds	RDT		quantitative	numeric	Rdt (% MIA/DW)
28	compounds	Starch1		quantitative	numeric	Dosage amidon (%poids/MIA) http://purl.obolibrary.org starch
29	compounds	Starch2		quantitative	numeric	amidon (g/gDW) http://purl.obolibrary.org starch
30	compounds	RHAMNOSE		quantitative	numeric	RHAMNOSE http://purl.obolibrary.org rhamnose
	...					

Additional subsets/ attributes can be added step by step, as soon as data are produced.









FRIM - Uploading your datasets in the data repository

➡ Merely dropping data files on the data repository (e.g. NAS) should allow users to access them by web services

Data capture

Minimal effort (PUT)

 samples.tsv
 plants.tsv
 harvests.tsv
 compounds.tsv
 a_attributes.tsv
 s_subsets.tsv

Your data subset files

Data subsets files and their associated metadata files must be compliant with the **TSV standard** (Tab-Separator-Values)

PUT


Data repository

 Z: (\\Storage)

 DataRepos

 frim1

Your dataset entry (named 'frim1' as example) within the data repository

mount



myhost.org

http://myhost.org/

GET



Data analysis/mining

Maximum efficiency (GET)

No database schema, no programming code and no additional configuration on the server side.



FRIM - *Checking online if your the data subset files are consistent*

Before dropping your data set in the data repository, check all points below:

Ref. Note	Note Description
1	A directory named as the dataset name should be actually created in the data repository; Be careful in the spelling, see note 6;
2	The s_subsets.tsv and a_attributes.tsv files should be present in the data repository
3	All data subset files declared in the s_subsets.tsv (col 5) should be available in the data repository
5	To be sure to have the right format, do a 'copy' of data from the spreadsheet then 'paste' them into a new file, then 'save as TSV format (separator: a tab character)'
4	1) all subsets in the a_attributes.tsv file (col 1) should be declared in the s_subsets.tsv file (col 3) 2) all subsets in the s_subsets.tsv file (col 3) should be declared in the a_attributes.tsv file (col 1) 3) all attribute names in the a_attributes.tsv file (col 2) should be available as a column in the corresponding data subset file declared in the s_subsets.tsv file (col 5)
6	Be careful in the spelling: 1) for data subset file names (col 5 in s_subsets.tsv), identifier name (col 2 in s_subsets.tsv), attribute names (col 2 in a_attributes.tsv), subset shortnames (col 3 in s_subsets.tsv and col 1 in a_attributes.tsv) and entry names (col 3 in a_attributes.tsv), only the alphanumerical characters and the underscore are allowed (i.e 'a-z', 'A-Z', '0-9' and '_'). Moreover, these names should not start with a digit! 2) for categorical names (col 4 in a_attributes.tsv), the number of terms and their spelling are fixed, namely: 'identifier', 'factor', 'quantitative', 'qualitative'. 3) for type (col 5 in a_attributes.tsv), the allowed names are restrited to 'numeric' or 'string'. 4) for descritpion, the allowed characters are : 0-9 a-z A-Z , : + * () [] { } - % ! / . ?
7	Identifiers declared in the s_subsets.tsv file (col 2) 1) should be declared as 'identifier' in the 'category' column in the a_attributes.tsv file (col 4) 2) should be available as an column item in the corresponding data subset file 3) should be the only one attribute declared as identifier for the corresponding data subset file in the a_attributes.tsv file (col 4)
8	Each subset having a 'father_rank' greater than 0 in the s_subsets.tsv file (col 2) 1) should have a column label in the corresponding data file that must be identical to the identifier label of the linked subset (i.e. corresponding to the 'father_rank' in col 1) 2) should have the linked subset identifier with no category (i.e. void) in the a-attributes.tsv file (col 4), except if the subset and the linked subset have the same identifier



FRIM - *Checking online if your the data subset files are consistent*

Fortunately, most of checking can be automatically done for you

Checklist

<http://myhost.org/check/frim1>

Type	Description	Information	Status
General	the dataset directory	Available	ok
General	the subset definition file	Available	ok
General	the attribute definition file	Available	ok
General	data subset files	all availables (11)	ok
s_subsets	definition file format	CSV-compliant	ok
a_attributes	definition file format	CSV-compliant	ok
a_attributes	Check if spelling of names are proper in a_attributes.tsv	all attributes seem ok	ok
plants.txt	Check if column names are the same that those declared in a_attributes.tsv	all attributes seem ok	ok
samples.txt	Check if column names are the same that those declared in a_attributes.tsv	all attributes seem ok	ok
aliquots.txt	Check if column names are the same that those declared in a_attributes.tsv	all attributes seem ok	ok
compounds.txt	Check if column names are the same that those declared in a_attributes.tsv	all attributes seem ok	ok
enzymes.txt	Check if column names are the same that those declared in a_attributes.tsv	all attributes seem ok	ok
pools.txt	Check if column names are the same that those declared in a_attributes.tsv	all attributes seem ok	ok
FRIM1Quantities.txt	Check if column names are the same that those declared in a_attributes.tsv	all attributes seem ok	ok
qnmr_metabo.txt	Check if column names are the same that those declared in a_attributes.tsv	all attributes seem ok	ok
plato_HexosesP.txt	Check if column names are the same that those declared in a_attributes.tsv	all attributes seem ok	ok
lipids_ag.txt	Check if column names are the same that those declared in a_attributes.tsv	all attributes seem ok	ok
aminoacids.txt	Check if column names are the same that those declared in a_attributes.tsv	all attributes seem ok	ok
identifiers	Check if identifiers are consistent	all identifiers seem consistent	ok
subsets	Check SQL on each subset	all SQL success	ok
merged_subsets	Check SQL on each merged subset	all SQL success	ok



FRIM - Testing online if your the data subset files are OK

<http://myhost.org/xml/frim1>

Subset	Description	Identifier	WSEntry	SetID	LinkID	CV_Term_ID	CV_Term_Name
plants	Plant features	PlantID	plant	1	0	http://purl.obolibrary.org/obo/PO_0000003	whole plant
samples	Sample features	SampleID	sample	2	1	http://purl.obolibrary.org/obo/OBI_1110046	organ harvesting
aliquots	Aliquots features	AliquotID	aliquot	3	2	http://purl.obolibrary.org/obo/PO_0009001	fruit
cellwall_metabo	Cell wall Compound quantifications	AliquotID	aliquot	4	3	http://purl.obolibrary.org/obo/CHEBI_24431	chemical entity
cellwall_metaboFW	Cell Wall Compound quantifications (FW)	AliquotID	aliquot	5	3	http://purl.obolibrary.org/obo/CHEBI_24431	chemical entity
activome	Activome Features	AliquotID	aliquot	6	3	http://purl.obolibrary.org/obo/CHEBI_24431	chemical entity
pools	Pools of remaining pools	PoolID	pool	7	2	http://purl.obolibrary.org/obo/OBI_1110046	organ harvesting
qMS_metabo	MS Components quantification	PoolID	pool	8	7	http://purl.obolibrary.org/obo/CHEBI_24431	chemical entity
qNMR_metabo	Pools of remaining pools	PoolID	pool	9	7	http://purl.obolibrary.org/obo/OBI_1110046	organ harvesting
plato_hexosesP	Hexoses Phosphate	AliquotID	aliquot	10	3	http://purl.obolibrary.org/obo/CHEBI_24431	chemical entity
lipids_AG	Lipids AG	AliquotID	aliquot	11	3	http://purl.obolibrary.org/obo/CHEBI_24431	chemical entity
AminoAcid	Amino Acids	AliquotID	aliquot	12	3	http://purl.obolibrary.org/obo/CHEBI_24431	chemical entity

To summarize

1. *Preparation and cleaning of the data sub-sets of files*
2. *Classification of each column within its right category*
3. *Connections between the dataset files based on identifiers*
4. *Creation of the metadata files: **s_subsets.tsv** and **a_attributes.tsv***
5. *Deposit of the dataset files in the data repository*
6. *Checking online if your the data subset files are consistent*
7. *Testing online the web-services on your dataset*

Note: Data subsets files and their associated metadata files must be compliant with the **TSV standard** (Tab-Separator-Values)

TSV is an alternative to the common comma-separated values (CSV) format, which often causes difficulties because of the need to escape commas