



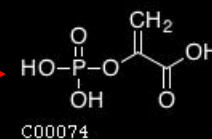
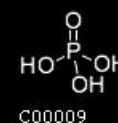
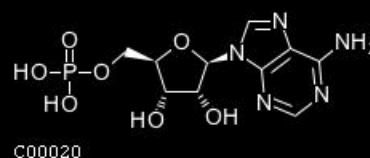
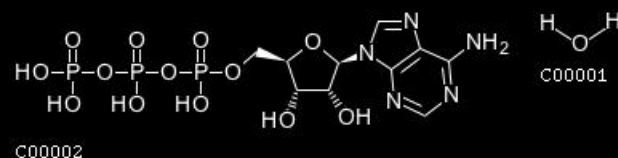
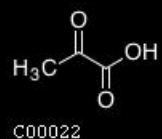
NIH:
West Coast Metabolomics Center



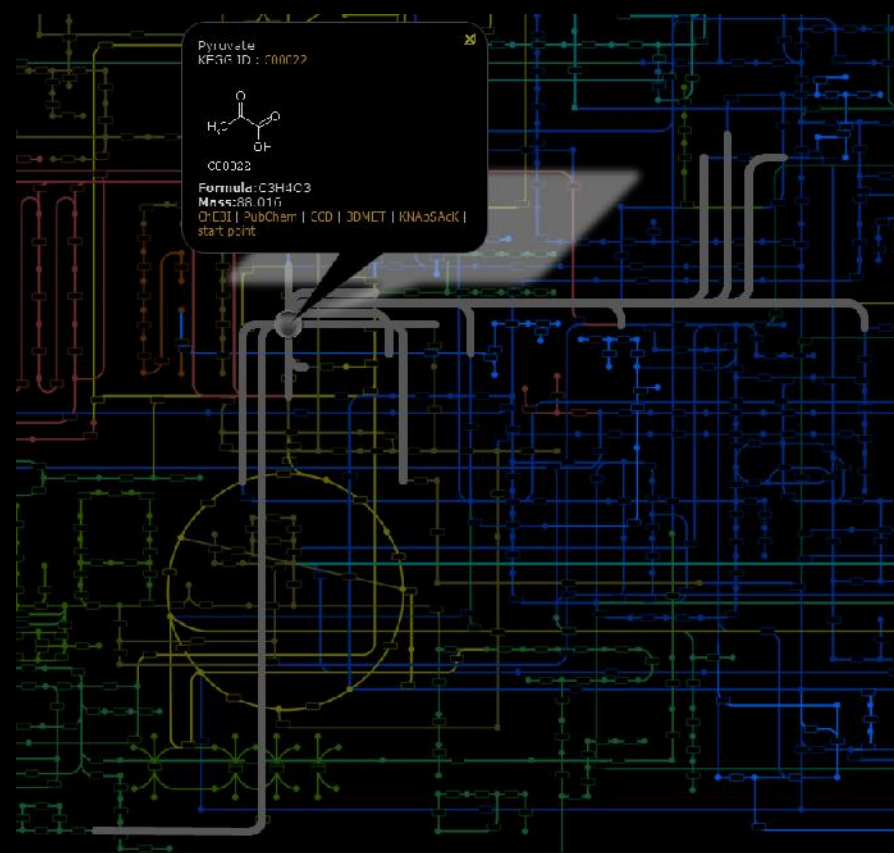
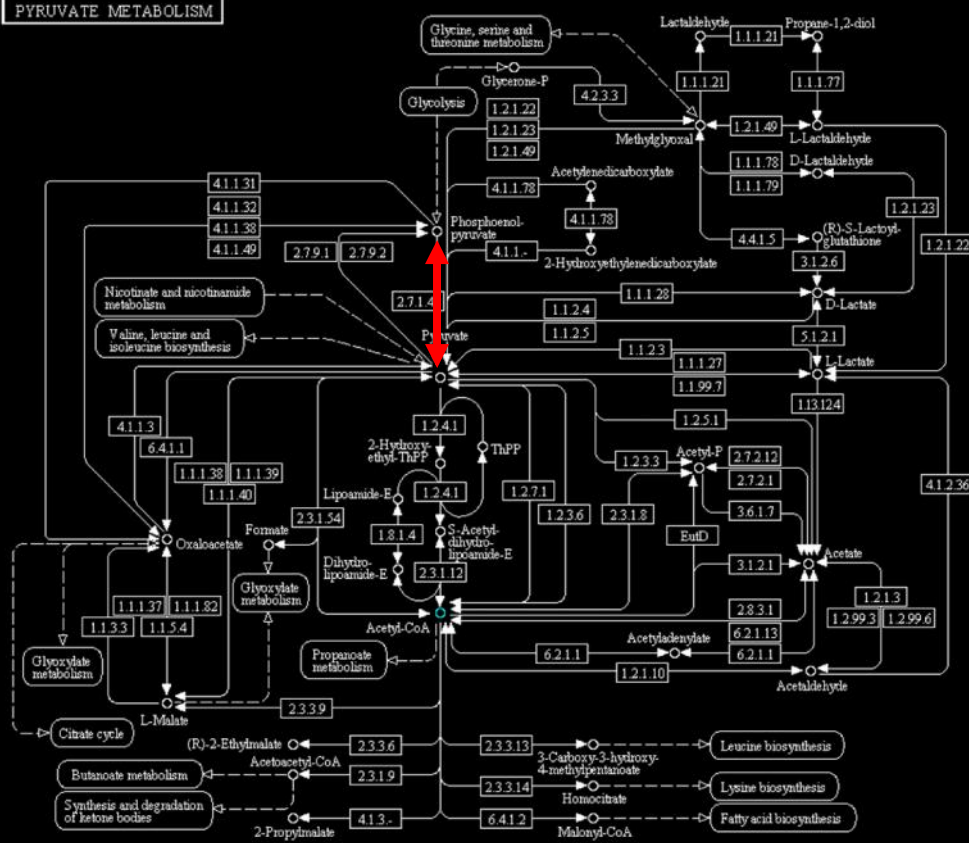
Strategies for Metabolomic Data Analysis

Dmitry Grapov, PhD

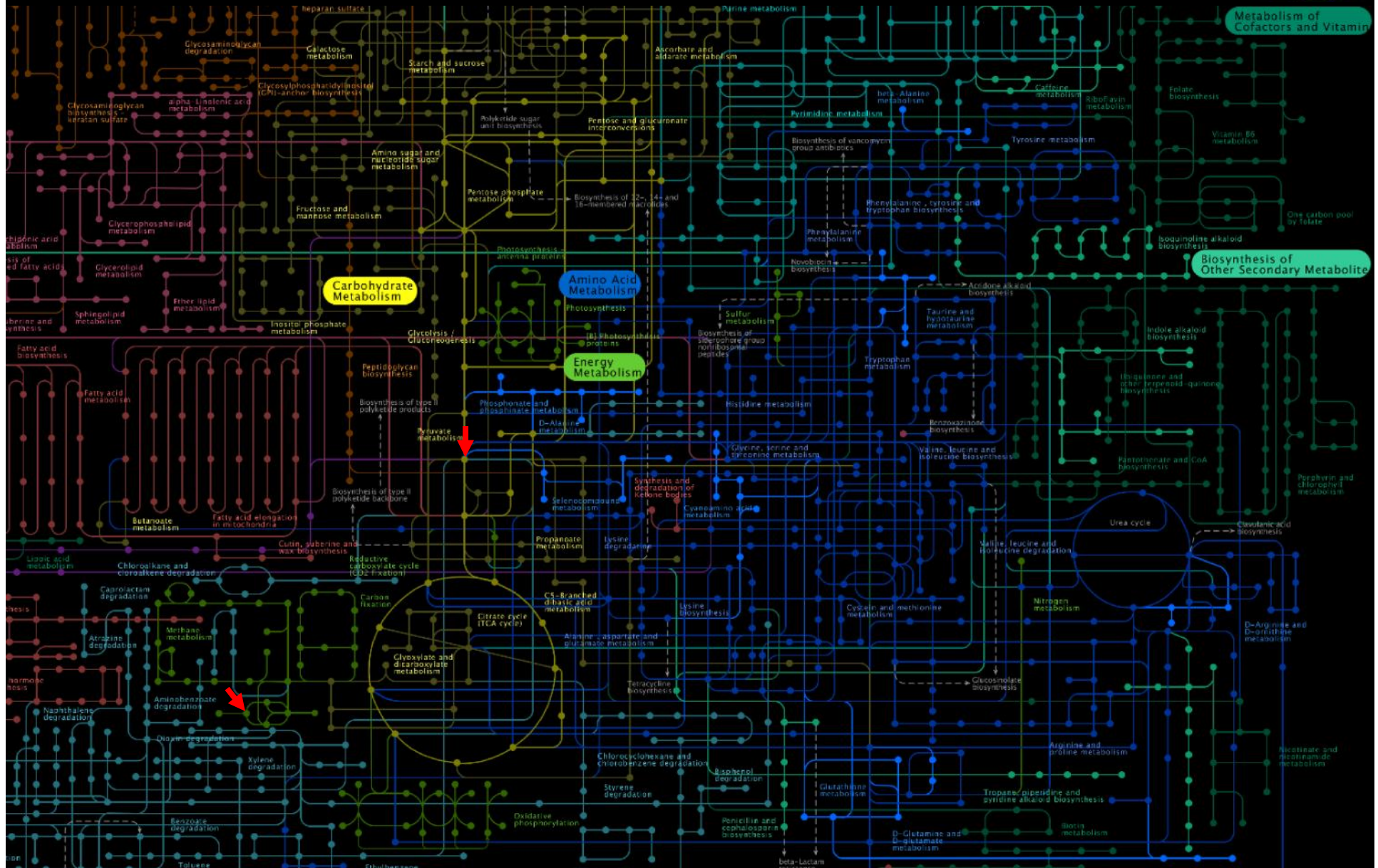
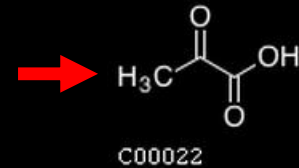
Goals?



PYRUVATE METABOLISM



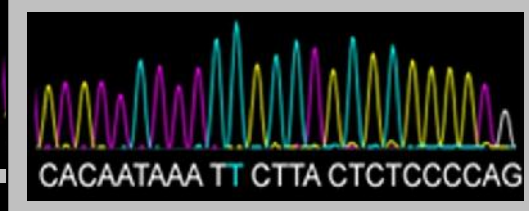
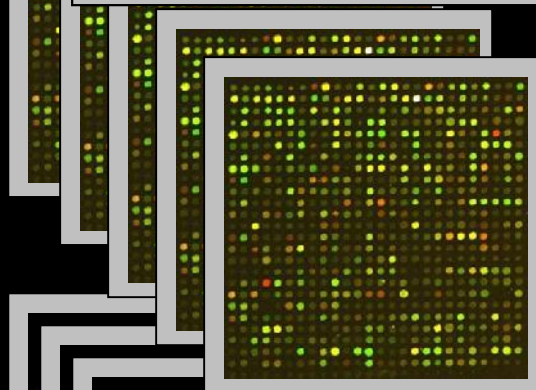
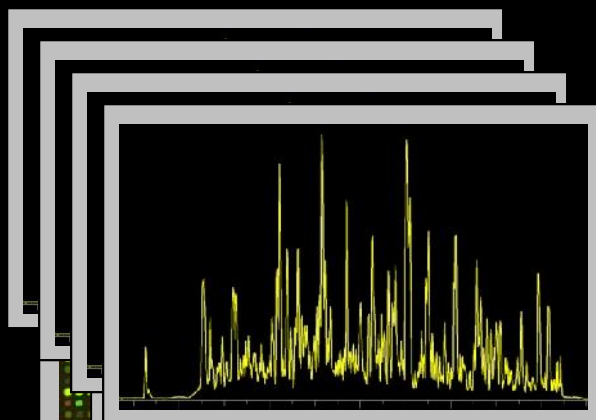
Metabolomics



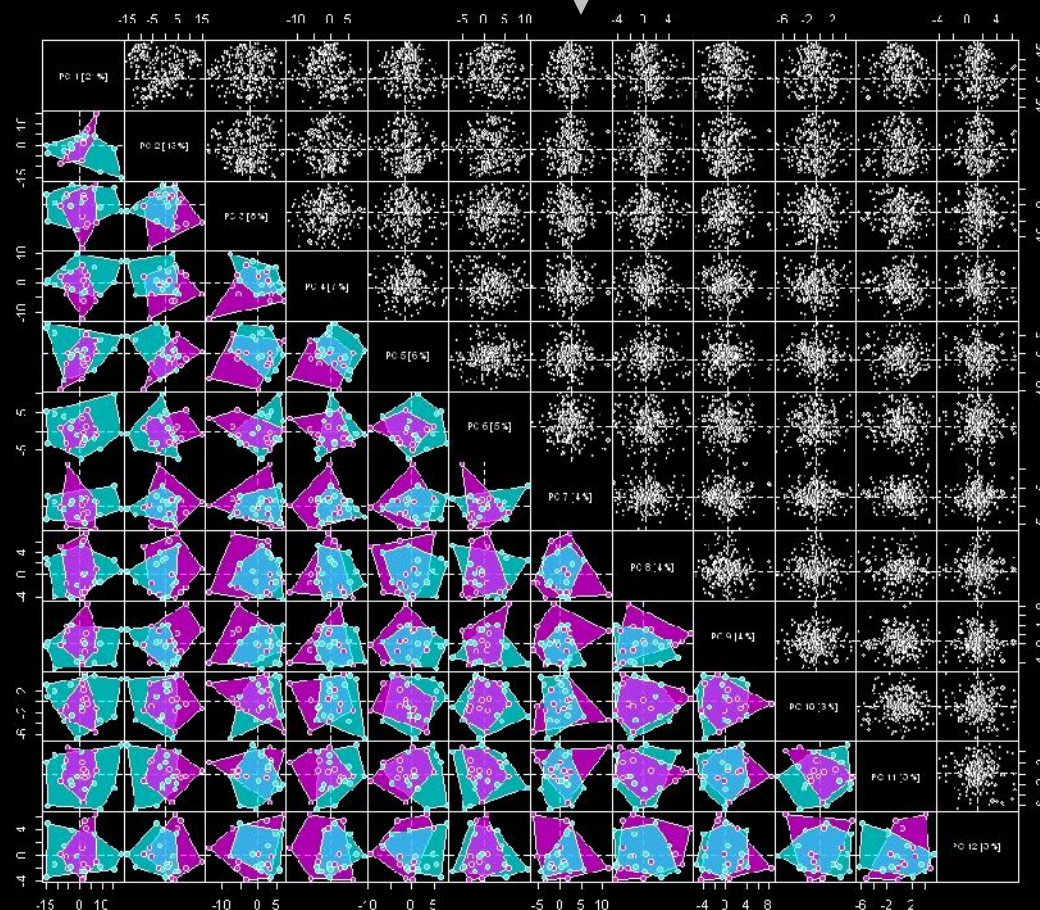
Analytical Dimensions



variables



Samples



Analyzing Metabolomic Data



- Pre-analysis
- Data properties
- Statistical approaches
- Multivariate approaches
- Systems approaches

code
grape variety
organ
vineyard
harvest time
sample prep
File
SetupX Class ID
SetupX Sample ID

user provided data

retention index

unique bin id

quantification ion

compound name

bin mass spec

PubChem CID

peak height

replaced values

VOC	BinBase	Name	RI	Quant	Ion	VOC	BB	ID	mass spec	PubChem ID
hexanol (2-)	333677	57	46468	30:10818.0	CID 13577	620412	672256	183937	730148	273004
45973	338054	36	45973	30:9405.0	3-	450312	43678	157573	40	38927
heptanol (2-)	348144	45	45431	33:208.0	37	709301	421494	127448	12	47019
45887	348085	95	45887	30:856.0	31-	31934	22043	5131	8924	20906
hexadecanal (2e,4e)	35	35	1:396.0	33	CID 637564	52139	41386	25676	20224	20113
45684	35	35	0:639.0	31-		223971	151719	85093	103616	127220
octen-3-ol (1-)	387675	57	45685	30:46.0	31	908257	745576	333106	347280	511140
hepten-2-ol (5-methyl)	390009	43	47178	36:68.0	37	112115	212384	88624	111432	247464
octanone (2-)	382202	58	45345	36:2.0	37	20142	17032	9030	9892	12978
pentyl	45700	30:250.0	31	CID 19602	23822	17261	14399	220620	16734	21440
linalool	47418	31:440.0	37		43824	38230	25162	3132	3465	43886
45926	397471	51	45926	32:6219.0		310716	205812	217432	245172	275715
hexenyl acetate (3e)	358524	67	45927	30:8661.0		266349	394820	475080	267456	310821
50987			0:508.0	31-		120245	59168	3		3085
hexyl acetate			0:450.0	31	CID 8908	15908	14407	1		3088
heptadecenal (2e,4e)	400037	51	400037	30:125.0	34	79904	57109	36030	37943	34344

Pre-analysis

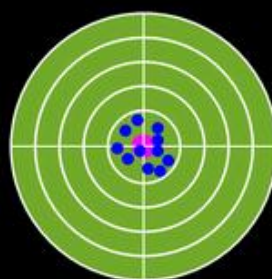
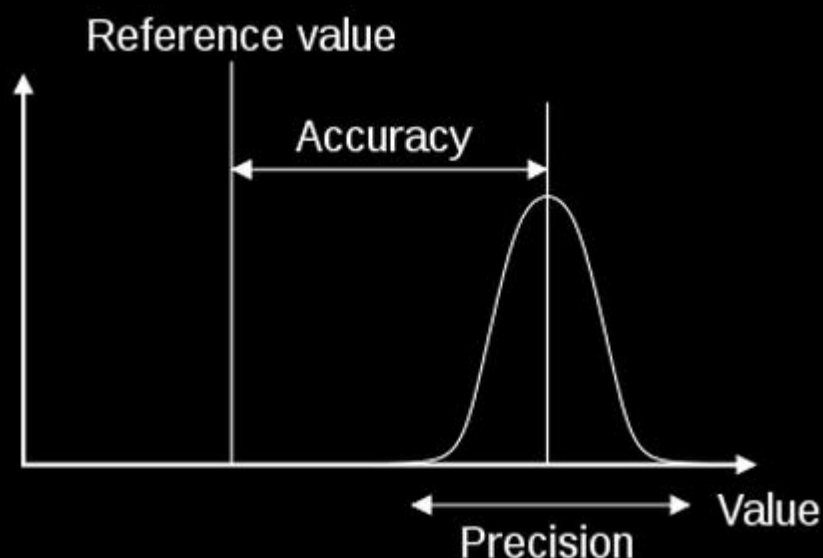


Data quality metrics

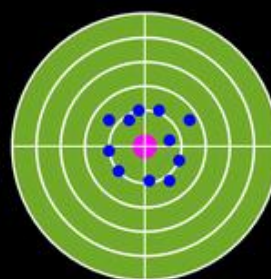
- precision
- accuracy

Remedies

- normalization
- outliers detection
- missing values imputation



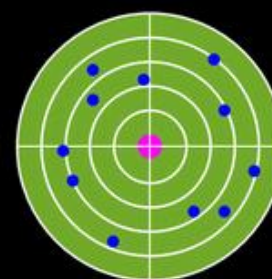
Accurate &
Precise



Accurate &
Imprecise



Inaccurate &
Precise

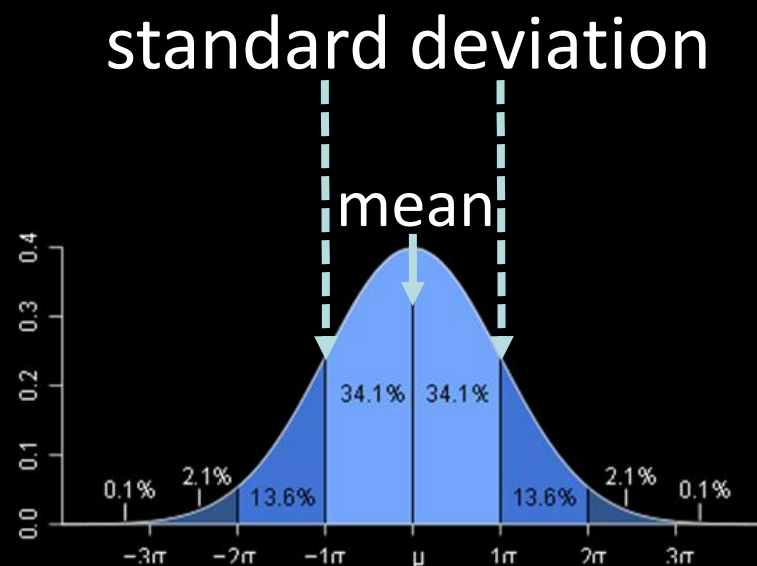
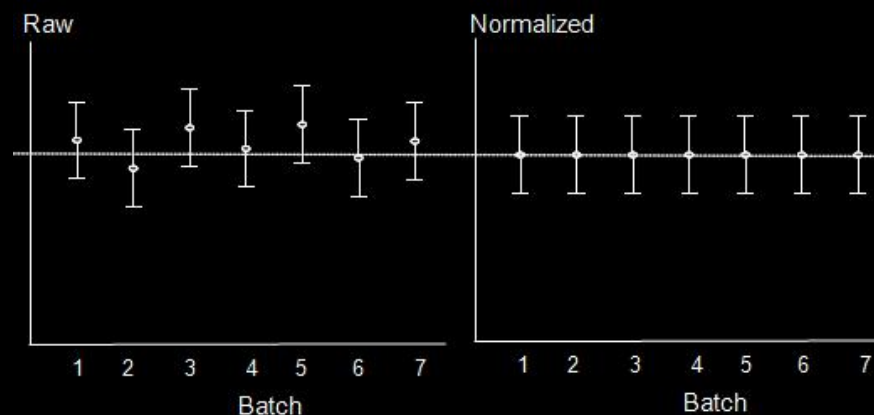


Inaccurate &
Imprecise

Normalization



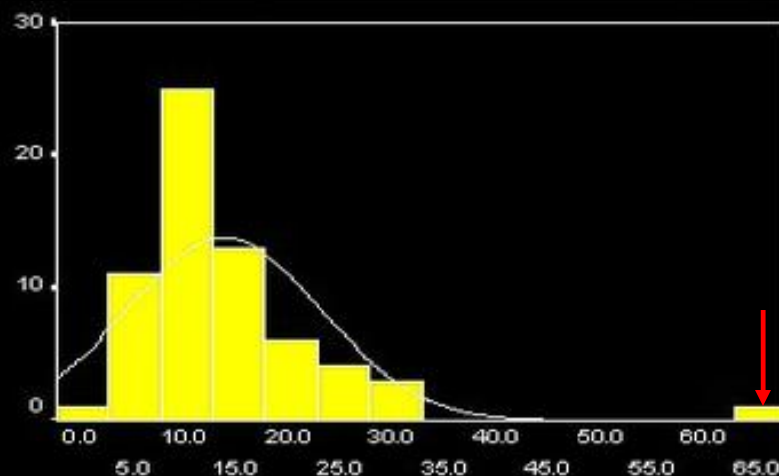
- sample-wise
 - sum, adjusted
- measurement-wise
 - transformation (normality)
 - encoding (trigonometric, etc.)



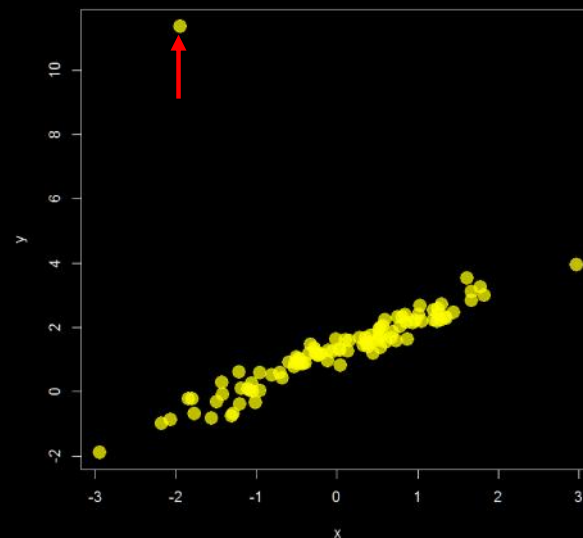
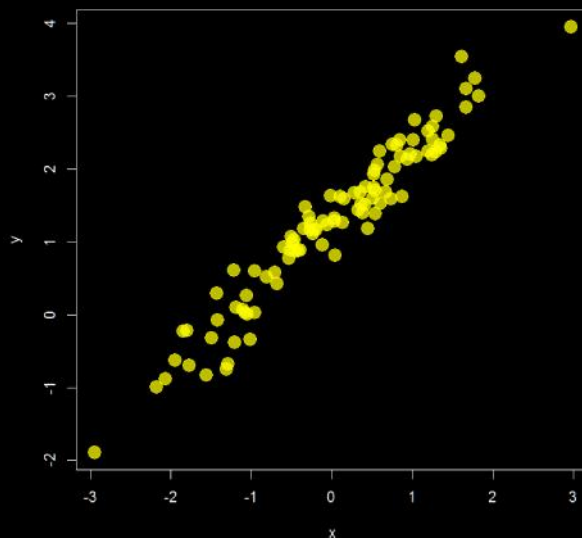
Outliers



- single measurements (univariate)



- two compounds (bivariate)



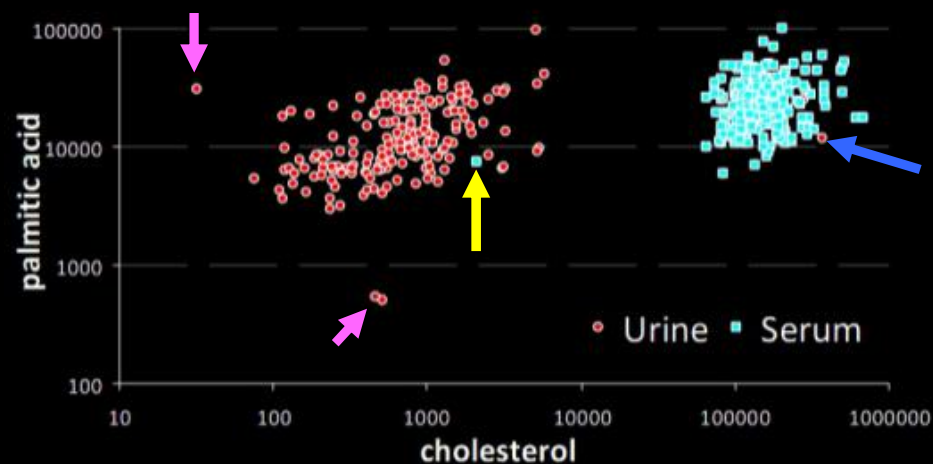
Outliers



univariate/bivariate

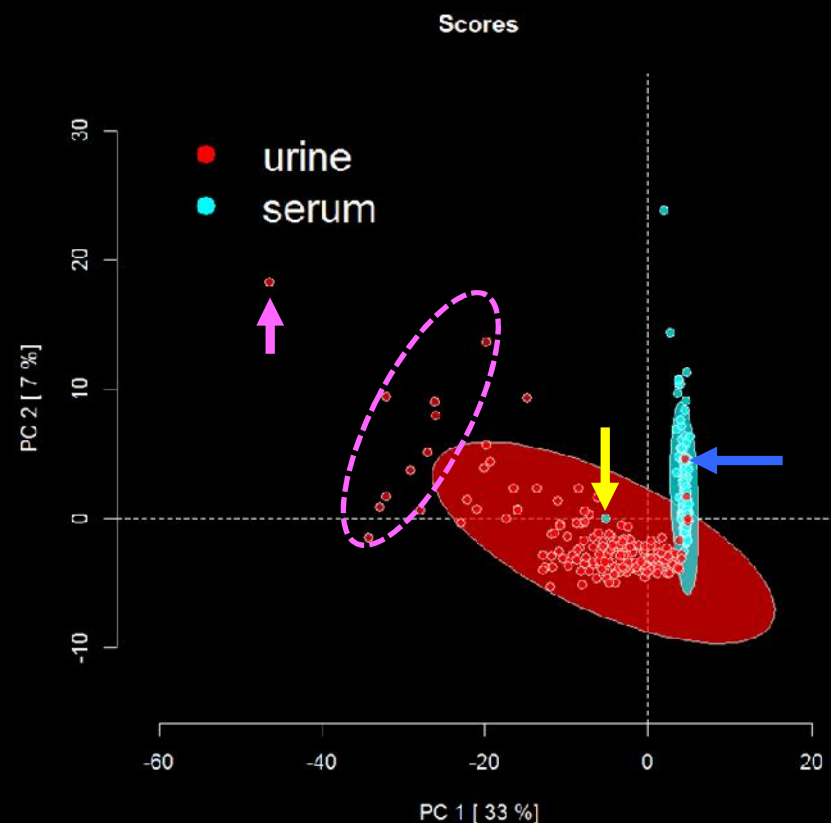
vs.

multivariate



→ outliers?

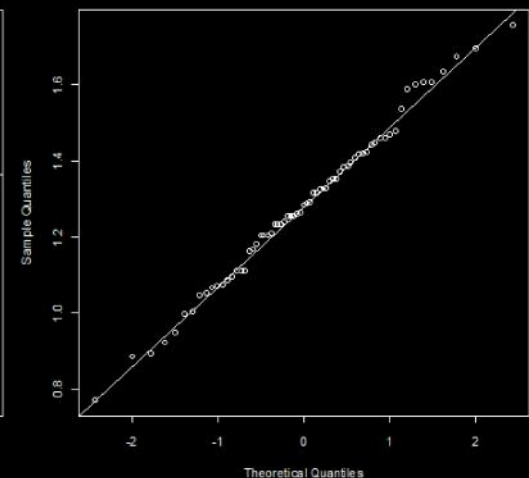
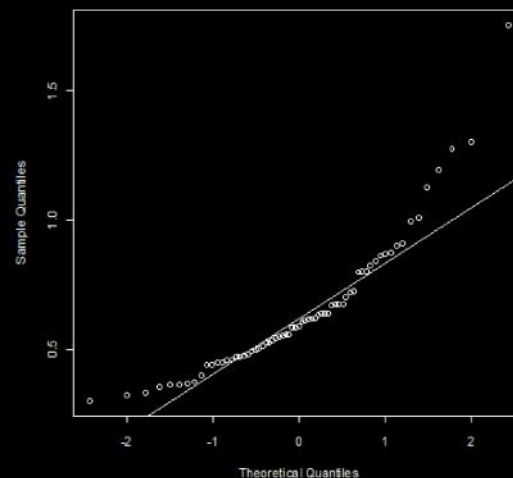
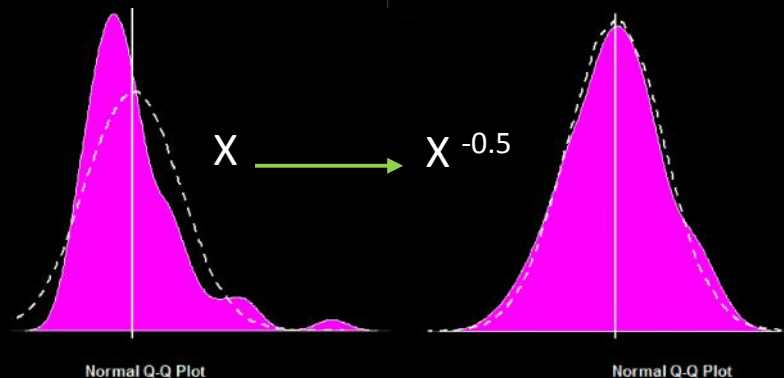
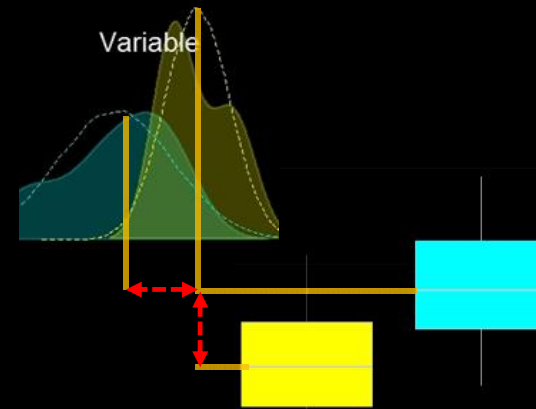
→ mixed up samples



Transformation

- logarithm (shifted)
- power (BOX-COX)
- inverse

Quantile-quantile (Q-Q) plots are useful for visual overview of variable normality



Missing Values Imputation

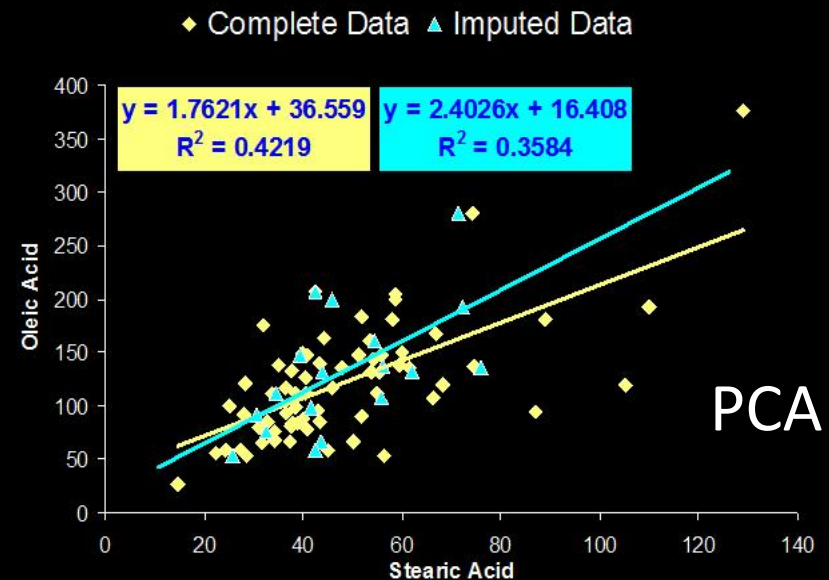
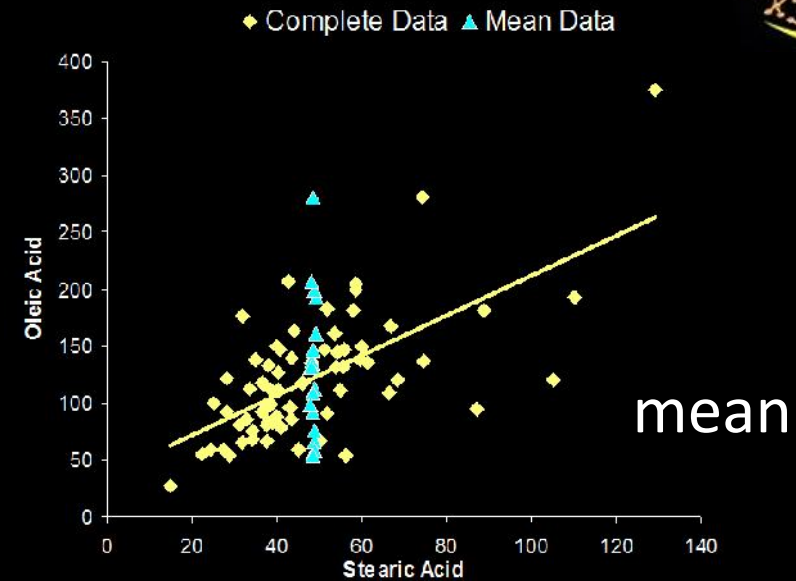


Why is it missing?

- random
- systematic
 - analytical
 - biological

Imputation methods

- single value (mean, min, etc.)
- multiple
- multivariate



Goals for Data Analysis

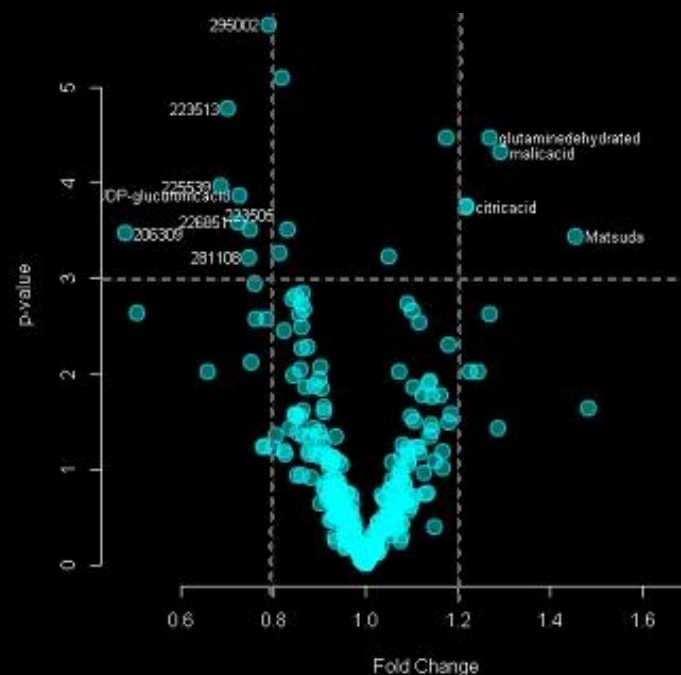


Exploration

Classification

Prediction

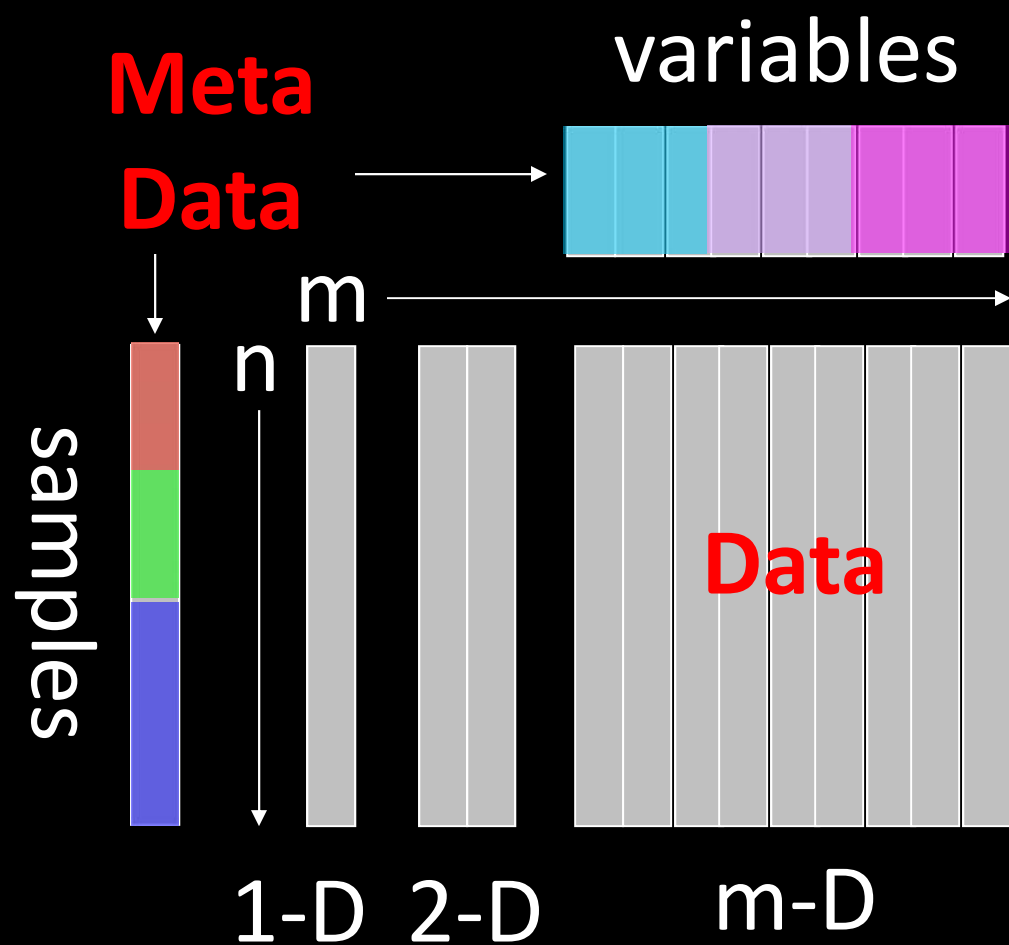
- Are there any trends in my data?
 - analytical sources
 - meta data/covariates
- Useful Methods
 - matrix decomposition (PCA, ICA, NMF)
 - cluster analysis
- Differences/similarities between groups?
 - discrimination, classification, significant changes
- Useful Methods
 - analysis of variance (ANOVA)
 - partial least squares discriminant analysis (PLS-DA)
 - Others: random forest, CART, SVM, ANN
- What is related or predictive of my variable(s) of interest?
 - regression
- Useful Methods
 - correlation



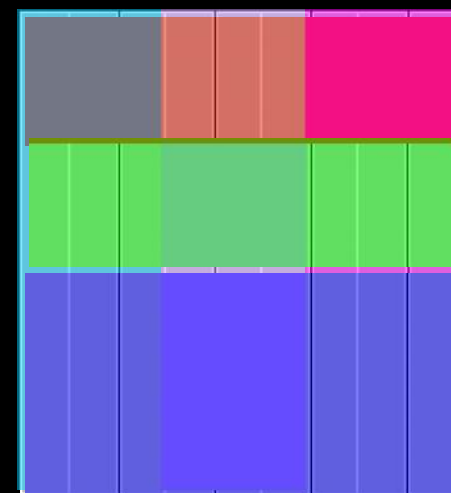
- ## • Data Types

- [illegible]

Data Complexity



**Experimental
Design =
complexity**



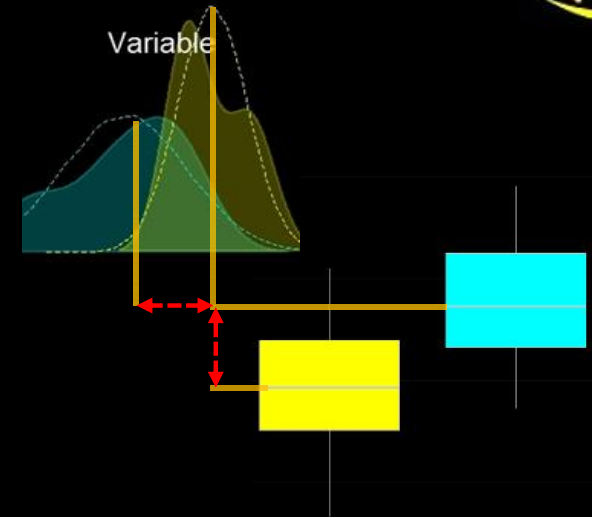
Variable # = dimensionality

Univariate Analyses

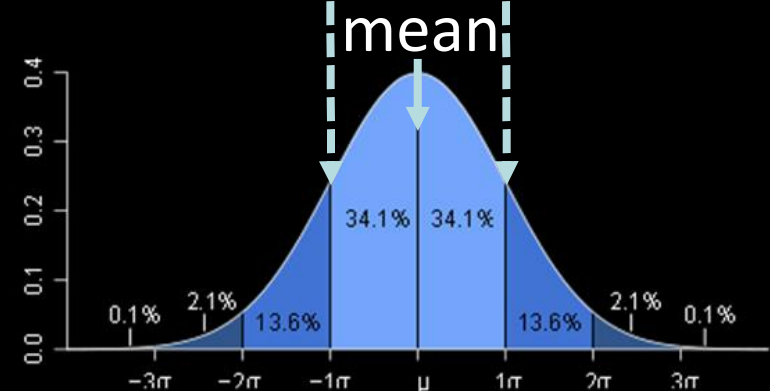


univariate properties

- length
- center (mean, median, geometric mean)
- dispersion (variance, standard deviation)
- Range (min / max)

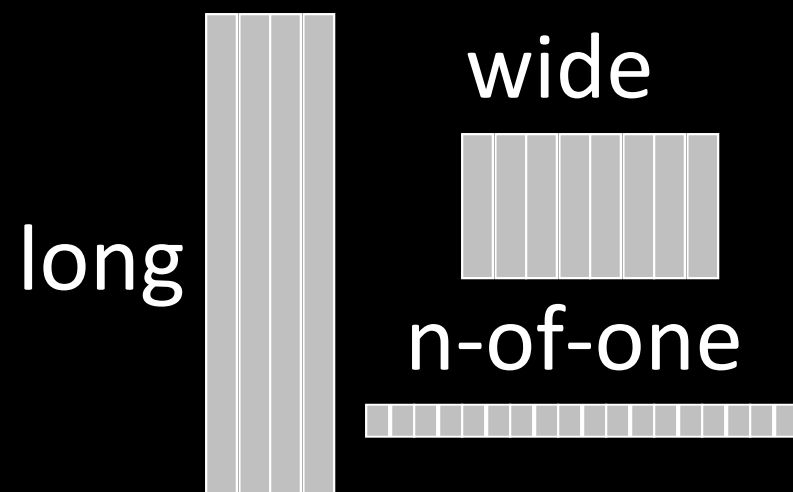
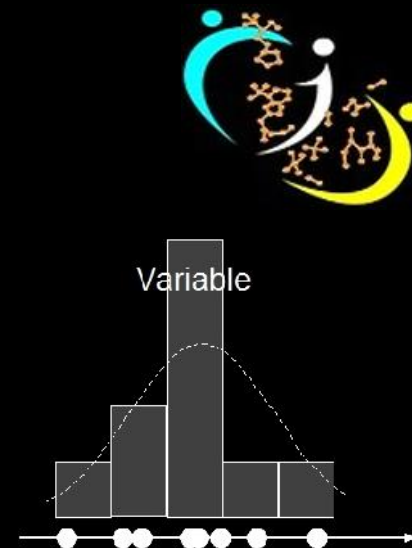


standard deviation



Univariate Analyses

- sensitive to distribution shape
 - parametric = assumes normality
- error in Y, not in X ($Y = mX + \text{error}$)
- optimal for long data
- assumed independence
- false discovery rate

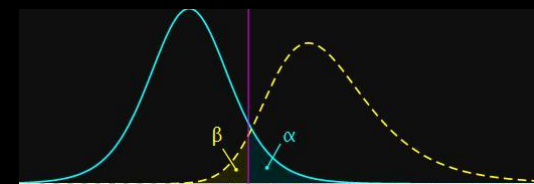


False Discovery Rate (FDR)

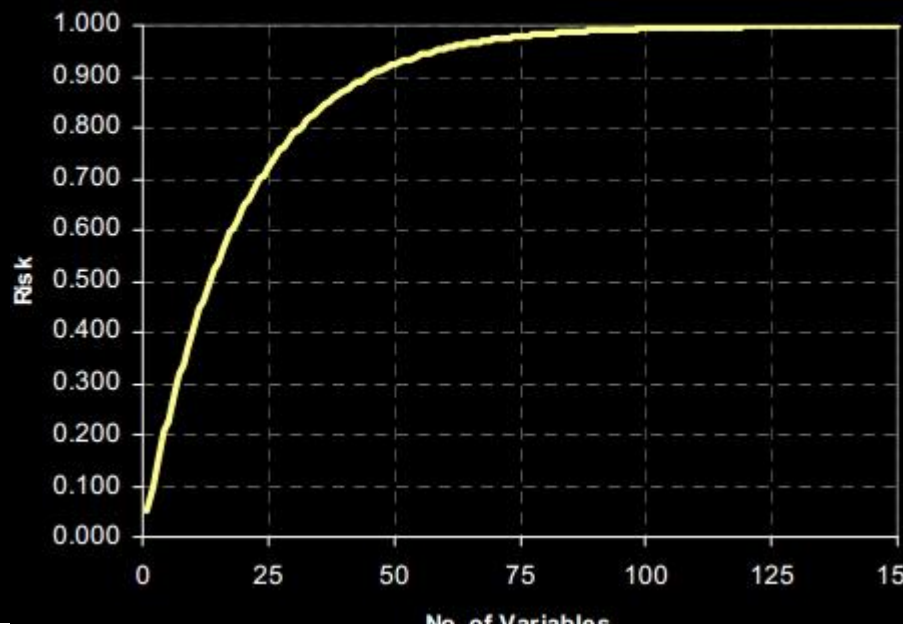


univariate approaches do not scale well

- Type I Error: False Positives
 - Type II Error: False Negatives
 - Type I risk =
 - $1-(1-p.value)^m$
- m = number of variables tested



Risk of Spurious Result



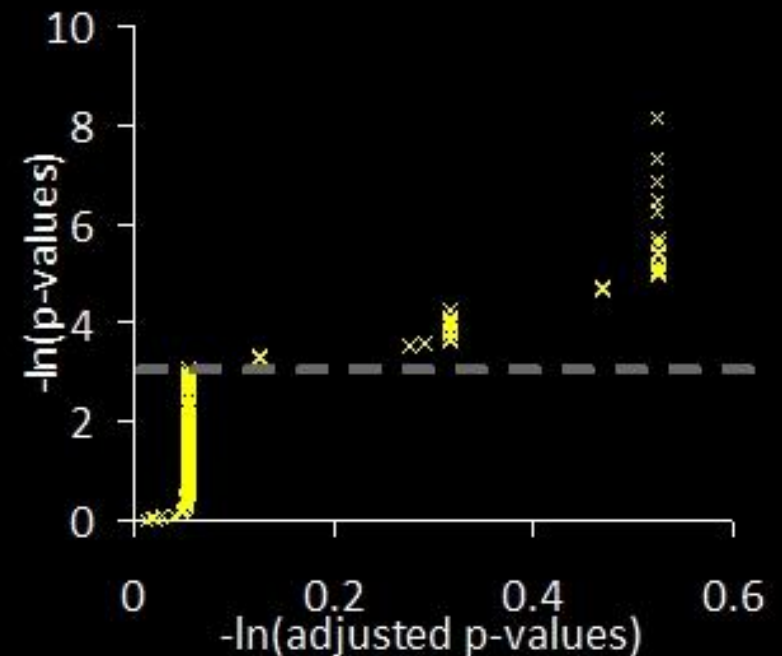
FDR correction

Example:

Design: 30 sample, 300 variables

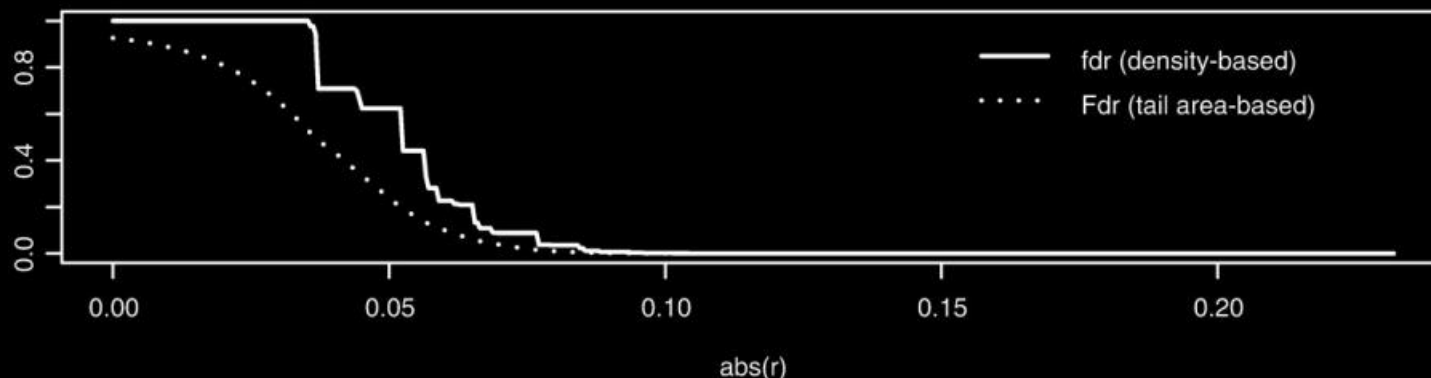
Test: t-test

FDR method: Benjamini and Hochberg (fdr) correction at $q=0.05$



Results

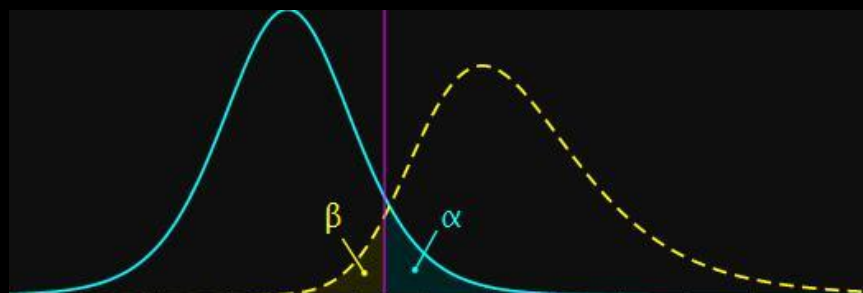
FDR adjusted p-values (fdr) or estimate of FDR (Fdr, q-value)



Achieving “significance” is a function of:

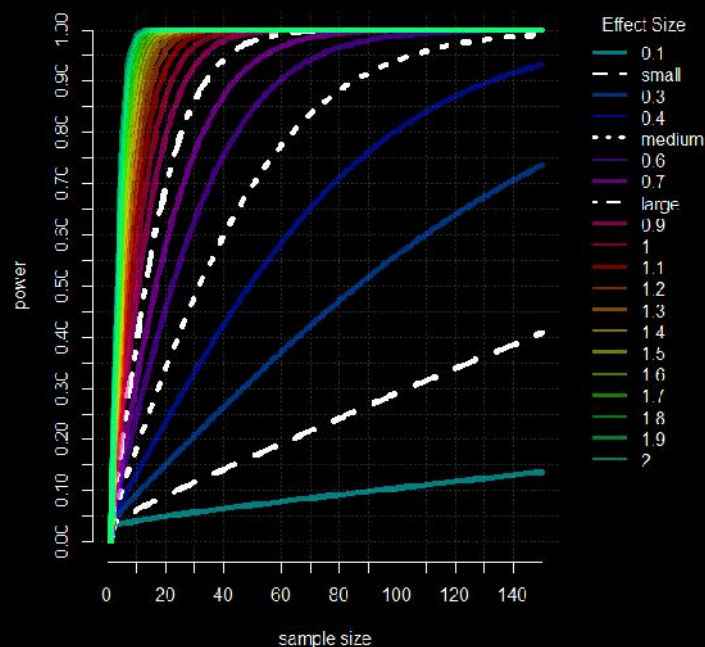


significance level (α) and power ($1-\beta$)



effect size (standardized difference in means)

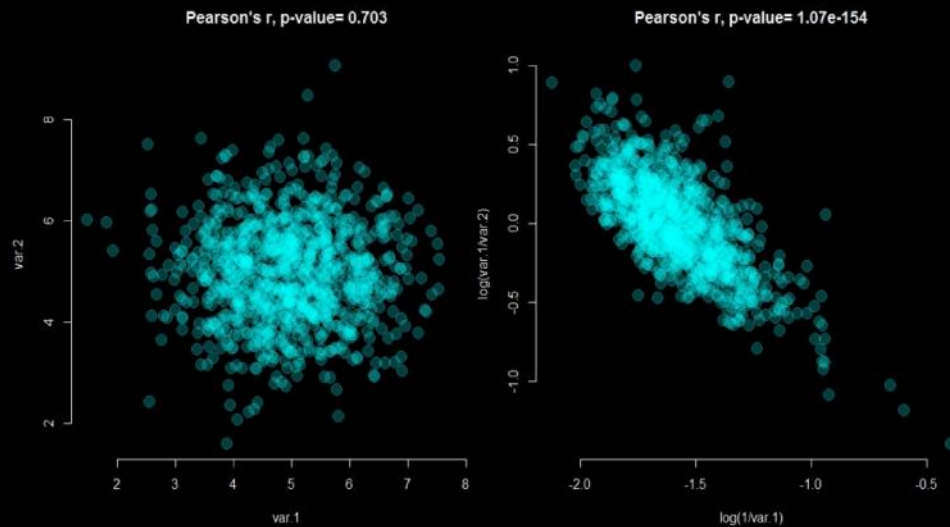
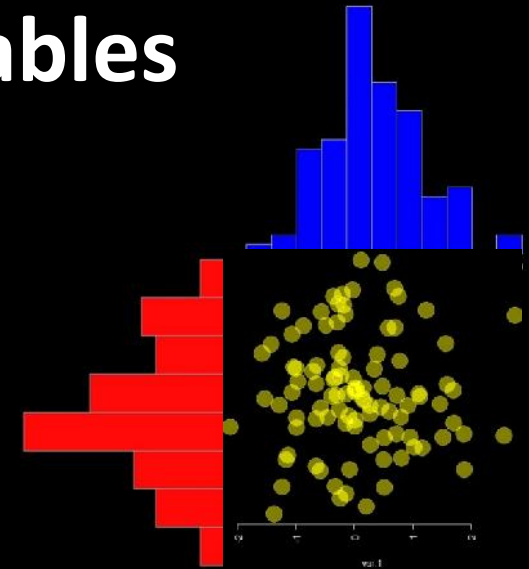
sample size (n)



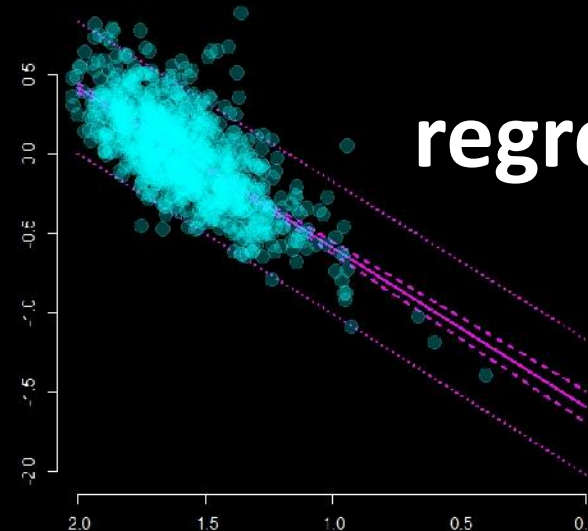
Bivariate Data

relationship between two variables

- correlation (strength)
- regression (predictive)



correlation

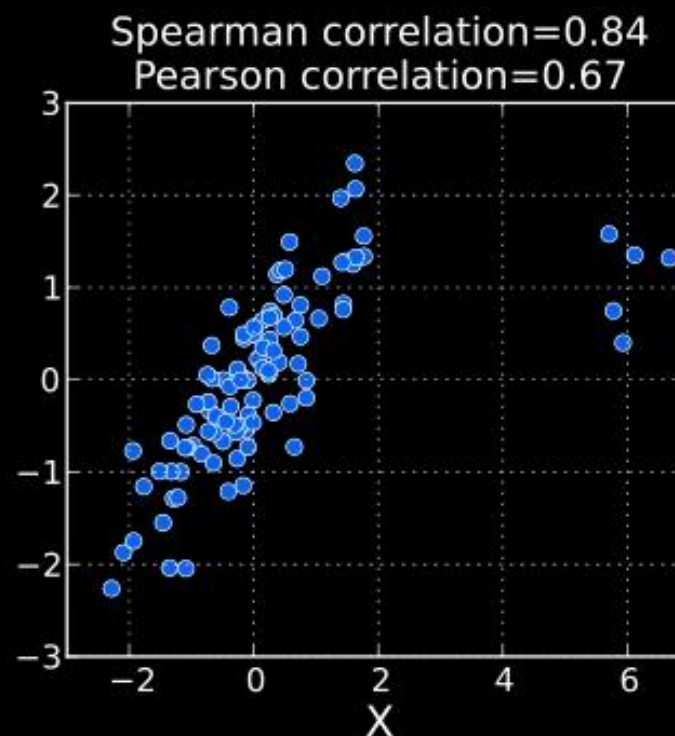
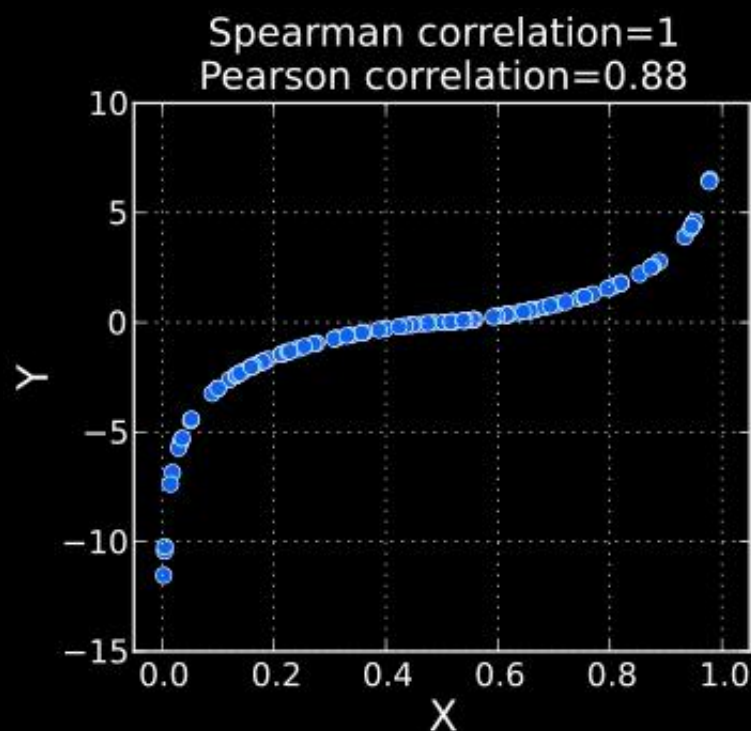


regression

Correlation

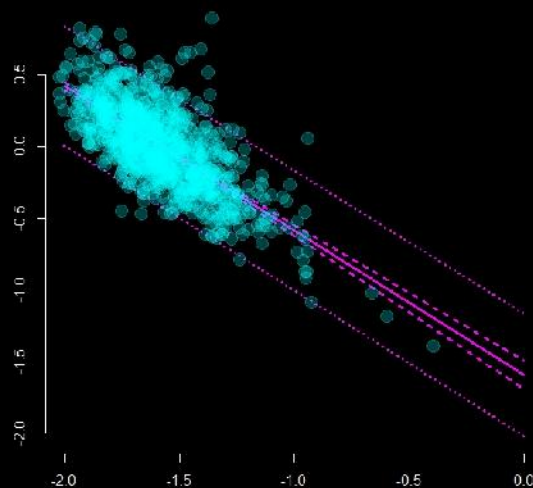
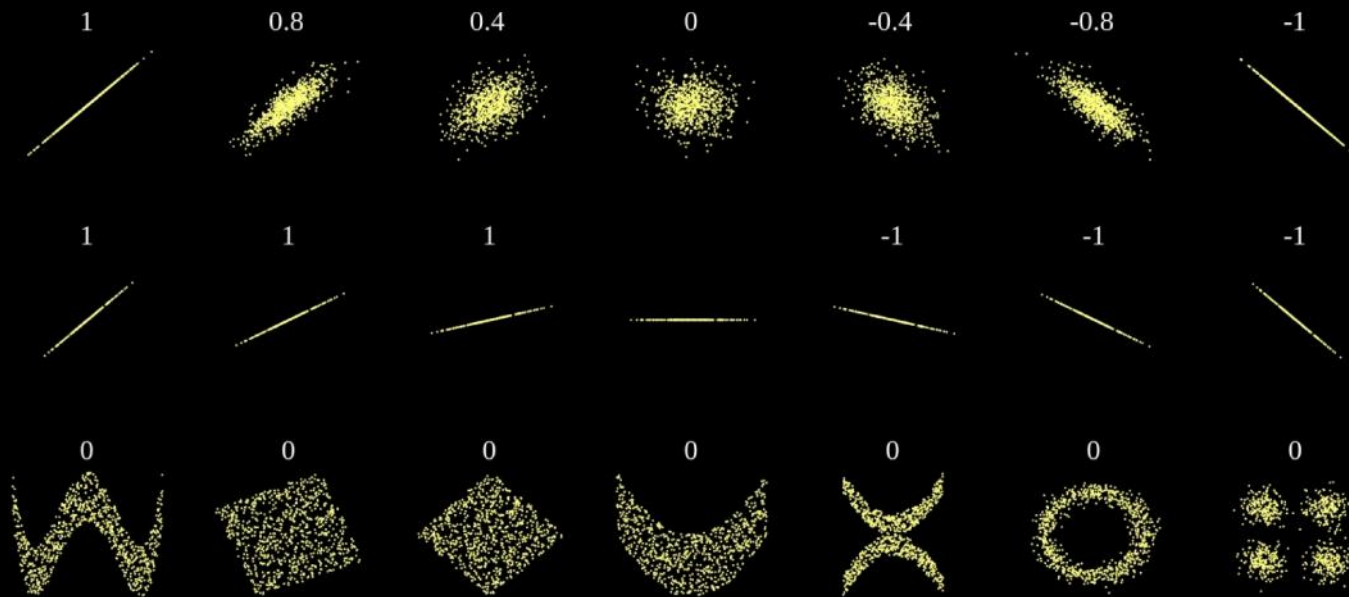


- Parametric (Pearson) or rank-order (Spearman, Kendall)



- correlation is covariance scaled between -1 and 1

Correlation vs. Regression



Regression describes the least squares or best-fit-line for the relationship ($Y = m \cdot X + b$)

Bivariate Example

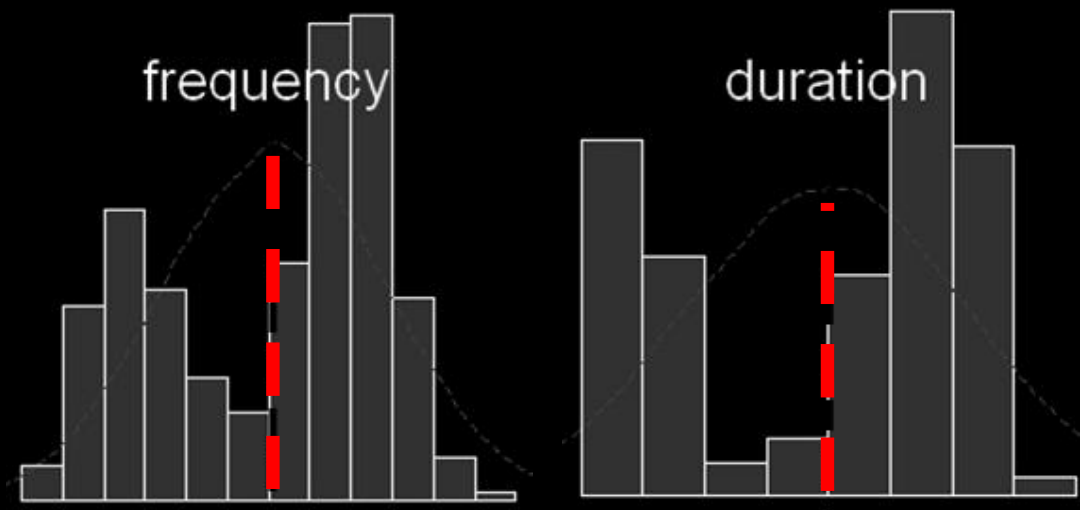
Goal: Don't miss eruption!

Data

- time between eruptions
 - 70 ± 14 min
- duration of eruption
 - 3.5 ± 1 min



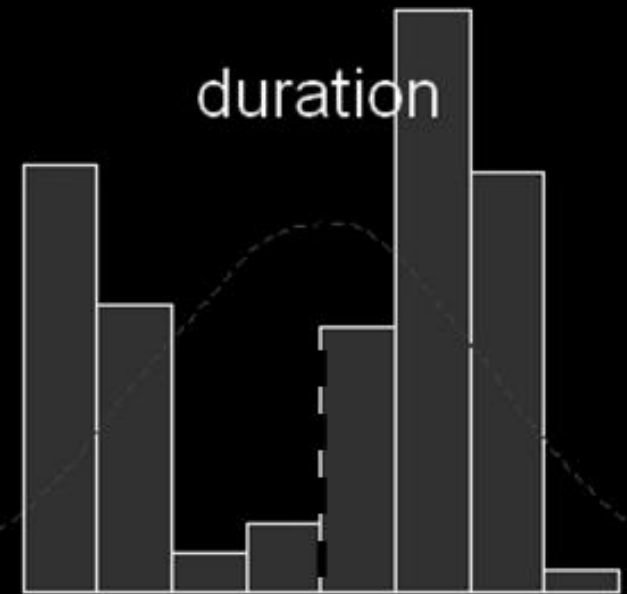
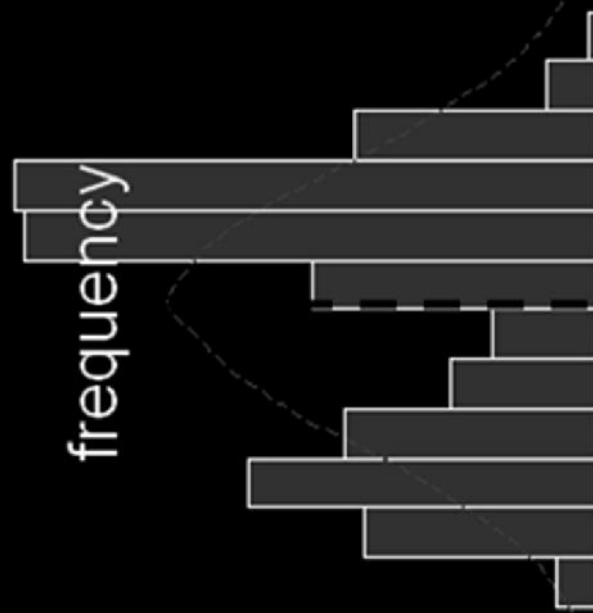
Old Faithful, Yellowstone, WY



Azzalini, A. and Bowman, A. W. (1990). A look at some data on the Old Faithful geyser. *Applied Statistics* **39**, 357–365

Bivariate Example

Two cluster pattern for both duration and frequency

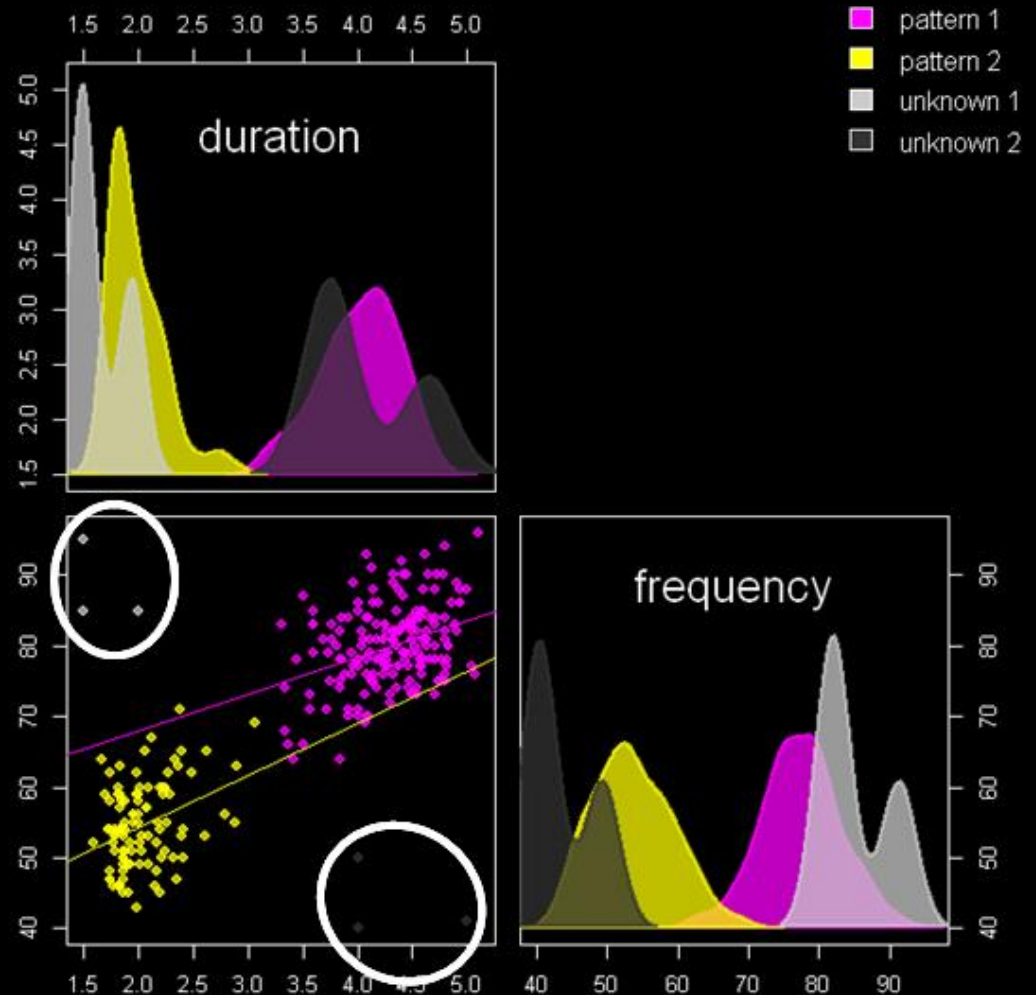


Bivariate Example

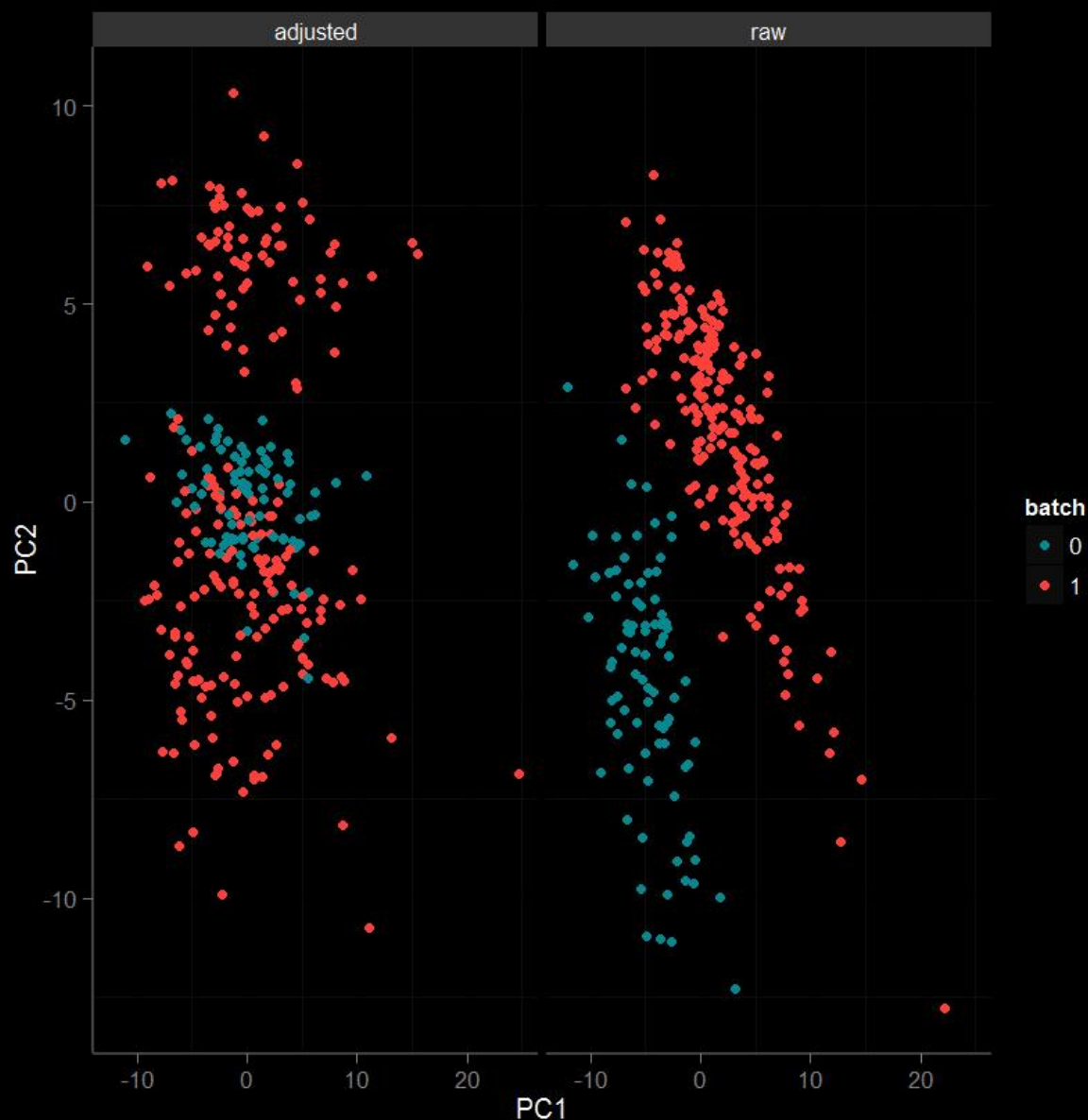


Noted deviations from
two cluster pattern

- Outliers?
- Covariates?



Covariates



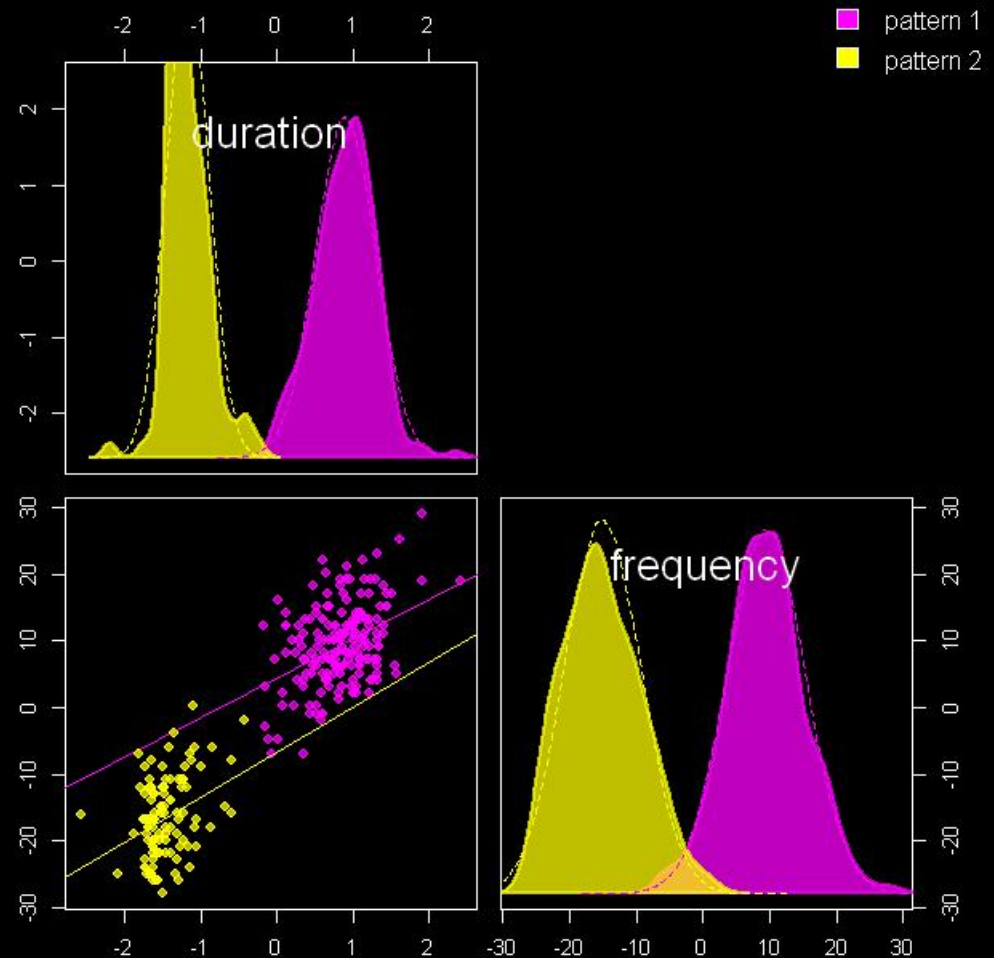
Trends in data which mask primary goals can be accounted for using covariate adjustment and appropriate modeling strategies

Bivariate Example

Noted deviations from
two cluster pattern
can be explained by
covariate:

Hydrofraking ☹️

Covariate adjustment
is an integral aspect of
statistical analyses
(e.g. ANCOVA)



Summary



Data exploration and pre-analysis:

- increase robustness of results
- guards against spurious findings
- Can greatly improve primary analyses

Univariate Statistics:

- are useful for identification of statically significant changes or relationships
- sub-optimal for wide data
- best when combined with advanced multivariate techniques

Resources



Web-based data analysis platforms

- MetaboAnalyst(<http://www.metaboanalyst.ca/MetaboAnalyst/faces/Home.jsp>)
- MeltDB(<https://meltdb.cebitec.uni-bielefeld.de/cgi-bin/login.cgi>)

Programming tools

- The **R** Project for Statistical Computing(<http://www.r-project.org/>)
- Bioconductor(<http://www.bioconductor.org/>)

GUI tools

- imDEV(<http://sourceforge.net/projects/imdev/?source=directory>)