# **Advanced Statistical Topics**

## Day 1 - Multiple testing, Randomization approaches, dimension reduction

Claus Thorn Ekstrøm UCPH Biostatistics ekstrom@sund.ku.dk

November 6th, 2023



# Practicalities

- Every day from **8.15 to 15** in room **7.0.08/7.0.18/35.3.13**
- Internet access: eduroam or KU-guest
- R
- Breaks as we go along
- Please read before classes

## Course overview

- Multiple testing, dimension reduction, randomization tests
   High-dimensional data, penalized regression, bootstrap, cross-validation
- 3. Classification/regression tress, random forests
- 4. Multiple imputation techniques



#### web.stanford.edu/~hastie/CASI/

#### Data sizes. The $N \ll P$ problem



#### The "Big Data" revolution

- 1. "Big *P* small *N*" problem with many modern large-scale-datasets: registers, images, \*-omics, ...
- 2. Need to reduce the dimension in some way
- 3. How do we evaluate significance when we have used the data for feature selection?
- 4. Multiple testing becomes an issue --- not just for high-dimensional data

#### Manhattan plot



#### Multiple comparison problems

Errors committed when testing a single null hypotheses,  $H_0$ 

Analysis result	H0 true	H0 false
Reject	Q	1-β
Don't reject	1-α	β

lpha is the significance level

1-eta is the power

#### Multiple comparison problems

The family-wise error rate (FWER) is the probability of making at least one type I error (false positive).

For *m* tests we have

 $FWER = P(\cup (p_i \leq lpha))) = 1 - P( ext{no false positives}) = 1 - (1 - lpha)^m \leq mlpha$ 

where the third equality only holds under independence, but the inequality holds due to Boole's inequality.

#### Multiple comparison problems

Number of errors committed when testing m null hypotheses.

Analysis result	H_0 true	H_0 false	Total
Reject	V	S	R
Don't reject	U	Т	m-R
Total	$m_0$	$m-m_0$	m

Here R, the number of rejected hypotheses/discoveries. V, S, U and T are unobserved. The FWER is

$$FWER = P(V > 0) = 1 - P(V = 0)$$

#### **Bonferroni correction**

The most conservative method but is free of dependence and distributional assumptions.

$$FWER = 1 - P(V=0) = 1 - (1-lpha)^m \leq mlpha$$

So set the significance level for each individual test at  $\alpha/m$ .

In other words we reject the *i*th hypothesis if

$$mp_i \leq lpha \Leftrightarrow p_i \leq rac{lpha}{m}$$

#### Sidak correction

$$1-(1-lpha)^m=lpha^* \Leftrightarrow lpha=1-\sqrt[m]{1-lpha^*}$$

Slightly less conservative than Bonferroni (but not much). Requires independence!

#### Holm correction

1. Compute and order the individual p-values:  $p_{(1)} \leq p_{(2)} \leq \cdots \leq p_{(m)}$ . 2. Find  $\hat{k} = \min\{k : p_{(k)} > \frac{\alpha}{m+1-k}\}$ 3. If  $\hat{k}$  exists then reject hypotheses corresponding to  $p_{(1)} \leq p_{(2)} \leq \cdots \leq p_{(\hat{k}-1)}$ 

#### Holm correction

Controls the FWER: Assume the (ordered) k is the first wrongly rejected true hypothesis. Then  $k \leq m - (m_0 - 1)$ .

Hypothesis k was rejected so

$$p_{(k)} \leq rac{lpha}{m+1-k} \leq rac{lpha}{m+1-(m-(m_0-1))} \leq rac{lpha}{m_0}$$

Since there are  $m_0$  true hypotheses then (Bonferroni argument) the probability that one of them is significant is at most  $\alpha$  so FWER is controlled.

#### Practical problems

• While guarantee of FWER-control is appealing, the resulting thresholds often suffer from low power.

In practice, this tends to wipe out evidence of the most interesting effects

• FDR control offers a way to increase power while maintaining some principled bound on error

#### False discovery rate

Number of errors committed when testing m null hypotheses.

Analysis result	H_0 true	H_0 false	Total
Reject	V	S	R
Don't reject	U	Т	m-R
Total	$m_0$	\$m-m_0\$	m

Proportion of false discoveries is  $Q = \frac{V}{R}$ . [Set to 0 for R = 0]

The false discovery rate is  $FDR = E(Q) = E(\frac{V}{R})$ 

### **Estimating FDR**



#### **Estimating FDR**



#### **Estimating FDR**



#### Estimating FDR — BH step-up

Benjamini-Hochberg step-up procedure to control the FDR at  $\alpha$ .

1. Compute and order the individual p-values:  $p_{(1)} \leq p_{(2)} \leq \cdots \leq p_{(m)}$ . 2. Find  $\hat{k} = \max\{k : \frac{m}{k} \cdot p_{(k)} \leq \alpha\}$ 3. If  $\hat{k}$  exists then reject hypotheses corresponding to  $p_{(1)} \leq p_{(2)} \leq \cdots \leq p_{(\hat{k})}$ 

#### Estimating FDR — BH step-up

*p*-values

$$egin{array}{rcl} ilde{p}_{(1)}&=&\min\{ ilde{p}_{(2)},mp_{(1)}\}\ dots&&dots\ dots\ dots\$$

Note that each  $p_i$  is smaller or equal to the criterium in Holm's method so controls the FWER.

#### Estimating FDR — BH step-up

If iid of the  $m_0$  tests (and all tests independent) and ordered so the  $m_0$  true tests comes first. Control FDR at level q:

$$egin{aligned} E(V/R) &= \sum_{r=1}^m E[rac{V}{r} \mathbbm{1}_{R=r}] = \sum_{r=1}^m rac{1}{r} E[V \mathbbm{1}_{R=r}] \ &= \sum_{r=1}^m rac{1}{r} E[\sum_{i=1}^m \mathbbm{1}_{p_i \leq rac{qr}{m}} \mathbbm{1}_{R=r}] = \sum_{r=1}^m rac{m_0}{r} \, [\mathbbm{1}_{p_1 \leq rac{qr}{m}} \mathbbm{1}_{R=r}] = \cdots \ &= \sum_{r=1}^m rac{m_0}{r} \, [\sum_{i=1}^{m_0} \mathbbm{1}_{p_1 \leq rac{qr}{m}} \mathbbm{1}_{R=r}] \ &= q rac{m_0}{m} \leq q \end{aligned}$$

q values

The q-value is defined to be the FDR analogue of the p-value.

$$q ext{ value}(p_i) = \min_{t \geq p_i} \widehat{ ext{FDR}}(t)$$

The *q*-value of an individual hypothesis test is the minimum FDR at which the test may be called significant.

## q values

- When all *m* null hypotheses are true then FDR control is equivalent to FWER control.
- FDR approach generally gives more power than FWER control and fewer Type I errors than uncorrected testing.
- The FDR bound holds for certain classes of dependent tests. In practice, it is quite hard to "break"

## **Exercises**

# **Randomzation tests**

- Bootstrap
- Permutation tests

# Evaluating complex methods and data

When we have complex data (or perhaps just big data combined with simple methods) and non-parametric methods then we still with to evaluate them.

How stable are the results?

## The bootstrap/jackknife procedures

Whenever we provide an estimate (mean, proportion, ...) we *also* want to infer its precision!

We may or may not be able to formulate a full parametric (or semiparametric model).

The bootstrap procedure allows us to estimate the standard error even in complicated situations or for non-standard statistics.



#### Statistics 101: multiple samples

Different samples will result in different outcomes.

If we had the means to produce several samples we would know the sampling distribution.









#### Resampling

Have observations  $x_i \sim F$  and an estimate

$$\hat{ heta}=s(x)$$

for some estimation algorithm.

Want the SE of  $\hat{\theta}$ .

#### The jackknife estimate

Estimate  $\hat{\theta} N$  times - once with each observation removed.

$$\hat{ heta}_{(-i)}=s(x_{(-i)})$$
 .

Then the jackknife estimator is

$$SE(\hat{ heta}) = \sqrt{rac{N-1}{N} \sum_{i=1}^{N} (\hat{ heta}_{(-i)} - \hat{ar{ heta}}_{(-)})^2}$$

Reduces to the standard error if  $\hat{\theta}$  is a sample average.

#### The jackknife estimate

- Non-parametric, no assumptions on F, samples of size N-1
- In general: the jackknife standard error is upwardly biased.
- Only to be used with smooth, differentiable statistic

```
x <-c(8.26, 6.33, 10.4, 5.27, 5.35, 5.61, 6.12, 6.19, 5.2,
7.01, 8.74, 7.78, 7.02, 6, 6.5, 5.8, 5.12, 7.41, 6.52, 6.21,
12.28, 5.6, 5.38, 6.6, 8.74)
CV <- function(x) sqrt(var(x))/mean(x)
CV(x)
```

[1] 0.2524712

library("bootstrap") ; res <- jackknife(x, CV)</pre>

#### The jackknife estimate

res

\$jack.se
[1] 0.05389943

\$jack.bias
[1] -0.009266436

\$jack.values

[1] 0.2563873 0.2565586 0.2384298 0.2507329 0.2513200
[6] 0.2530603 0.2557374 0.2560293 0.2501992 0.2580969
[11] 0.2541045 0.2577524 0.2581067 0.2551946 0.2571038
[16] 0.2541711 0.2495662 0.2581975 0.2571609 0.2561093
[21] 0.2020978 0.2529980 0.2515338 0.2573745 0.2541045

\$call

jackknife(x = x, theta = CV)

#### Estimating bias

When  $\hat{ heta}$  is unbiased then

$$E(\hat{ar{ heta}}) = rac{1}{N}\sum_{i=1}^N E({\hat{ heta}}_{(-i)}) = heta_{i=1}$$

But if the procedure has bias

$$E(\hat{\theta}) = \theta + \frac{a}{N} + \frac{b}{N^2} + \text{rest}$$

then we can estimate the size of the bias from jackknife results.

#### Estimating bias

$$E(\hat{ar{ heta}}-\hat{ heta})=rac{a}{N(N-1)}+ ext{rest}$$

Thus,

$$ext{bias}_{ ext{jack}} = (N-1)(\hat{ar{ heta}} - \hat{ heta}) = rac{a}{N}$$

Furthermore, the bias-corrected jackknife estimate,

$${\hat heta}_{
m jack} = {\hat heta} - {
m bias}_{
m jack}$$

is an unbiased estimate of  $\theta$  up to second order.

#### Improvement on the jackknife

Instead of removing 1 observation at a time, remove d. Then there are  $\binom{N}{d}$  sets

$$SE = \sqrt{rac{N-d}{dinom{N}{d}}\sum(({\hat{ heta}}_{(Z)}-{\hat{ar{ heta}}}_{(-)})^2}$$

... or use the bootstrap!

#### Nonparametric bootstrap

If we could draw extra samples from the population it would be easy!

Use the sample as the population and generate "fake samples"

 $Population \rightarrow Sample \rightarrow "Fake sample"$ 

#### Nonparametric bootstrap

Get a random bootstrap sample from the sample *with replacement* 

$$x^* = (x_1^*, x_2^*, \dots, x_N^*)$$

Then we can get

$${\hat heta}^* = s(x^*)$$
 .

Do that B times and get information about the full distribution.

#### Jackknife vs bootstrap

- Jackknife provides stable results (will always get the same result) whereas bootstrap varies.
- Jackknife only estimates the variance of the point estimator whereas the bootstrap provides information on the distribution.

$$SE=\sqrt{rac{\sum({\hat{ heta}}^{*b}-{\hat{ heta}}^{*-})^2}{B-1}}$$

#### Nonparametric bootstrap in R

results <- bootstrap(x, 200, CV)</pre>



#### Parametric bootstrap

The nonparametric bootstrap made no assumptions about the distribution. Use distribution information if known.

- Fit model to data
- Draw *B* samples of random numbers from the fitted model
- Use those for bootstrap

Useful for small sample sizes (assuming the model holds), difficult evaluations, ... Sampling from the "wrong" distribution and forgetting the uncertainty. Retains the information in the explanatory variables but needs the error distribution.

#### Parametric bootstrap - resample residuals

- Fit model to data, keep predictions  $\hat{y}_i$  and compute a vector of residuals,  $\hat{\epsilon}_i = y_i \hat{y}_i$ .
- Create new sets of observations  $y^* = \hat{y}_i + \hat{\epsilon}_j$  using a random residual.
- Refit the model using the new set of response variables, and compute the statistic
- Do that *B* times

Retains the information in the explanatory variables. What to resample?

#### Rough R code

x <- c(5, 9, 8, 4, 7, 4, 2)
# Non-parametric bootstrap
x.star <- sample(x, replace = TRUE)</pre>

# Parametric bootstrap for assumed Gaussianity
x.star <- rnorm(length(x), mean = mean(x), sd = sd(x))</pre>

# Mean approximates the mean for Gaussian distribution for residuals
resids <- x - mean(x)
x.star <- mean(x) + sample(resids, replace=TRUE)</pre>

#### What to do?

Depends on the situation.

- The structure of the data might make some options easier.
- Belief about the parametric model would improve efficiency.
- Belief about the bias of the estimate would influence the choice.

#### Bootstrap confidence intervals

Standard 95% confidence intervals

 $\hat{ heta} \pm 1.96SE$ 

#### Could get that directly from the bootstrap results.

mean(results\$thetastar) + c(-1.96, 1.96)\*sd(results\$thetastar)

[1] 0.1497948 0.3262128

Only really works if the distribution is symmetric

#### Bootstrap percentile confidence intervals

Generate the "full" distribution. Cut off 2.5% at each end. Use an improvement that depends on the precision of the percentiles.

bcanon(x, 2000, CV)

\$confpoints

alpha bca point [1,] 0.025 0.3035158 [2,] 0.050 0.3307189 [3,] 0.100 0.3595159 [4,] 0.160 0.3875672 [5,] 0.840 0.6210273 [6,] 0.900 0.6629274 [7,] 0.950 0.7128380 [8,] 0.975 0.7218803

#### Bootstrap percentile confidence intervals

#### Can also use the t distribution

boott(x, CV)

\$confpoints

	0.001	0.01	0.025	0.05	0.1
[1,]	0.2793437	0.3234873	0.3348129	0.3488559	0.3805315
	0.5	0.9	0.95	0.975	0.99
[1,]	0.47362 0.	.6681084 0	.7423936 0	.7955465 1	.383483
[1,]	1.596856				
\$thet NULL	za –				
\$g NULL					

#### Permutation / randomization tests

#### Exchangeability

Exchangeability corresponds to the situation where the labels or order identifying the individual observations are uninformative.

Thus, the simultaneous distribution will not change after permutation

$$P(x_1,\ldots,x_N)=P(x_{\pi(1)},\ldots,x_{\pi(N)})$$

Required for the bootstrap and permutation tests to work. Fulfilled for iid.

In general, exchangeability fails when your sample can be stratified into sub-groups.

## Permutation tests

Permute the data in order to *remove* association between the variables of interest.

Can then get an idea of the null distribution.

What is the *actual* null hypothesis?

How can we adapt to that using permutation tests?

#### Principal component analysis

Principal component analysis (PCA) extracts a low dimensional set of features from a high dimensional data set with a motive to capture as much information as possible.

It is always performed on a symmetric correlation or covariance matrix of correlations among variables. This means the matrix should be numeric and have standardized data.

Dimension reduction of the *covariates* only - outcome not used! Unsupervised

#### Principal component analysis

Typically based on singular value decomposition (SVD):

 $X = U\Sigma V^t$ 

- U matrix of Eigenvectors of  $XX^t$  ( n imes n )
- $\Sigma$  diagonal matrix of Eigenvalues ( n imes p )
- $V^t$  matrix of with Eigenvectors of  $X^t X$  ( p imes p )

PCA reduces p dimensions of X to k principal components.

- $U\Sigma$  gives principal components
- $V^t$  gives loading factors (weights in rows)

#### Important features of PCA

- PCs are ordered by the decreasing amount of variance explained
- PCs are orthogonal i.e. uncorrelated to each other
- The columns of X should be mean-centered, because then the covariance matrix is  $\approx X^t X$ .

#### Principal component analysis

Warning: Using `size` aesthetic for lines was deprecated in ggplot2 3.4.0. i Please use `linewidth` instead. This warning is displayed once every 8 hours. Call `lifecycle::last\_lifecycle\_warnings()` to see where this warning was generated.



#### Principal component analysis

Compute the covariance matrix of the predictor data set x.
 Calculate the eigenvalues and corresponding eigenvectors of this covariance matrix

3. Eigenvectors correspond to orthogonal directions, sort by eigenvalue.

Reduce dimensionality so pick a unit vector u, and replace each data point with its projection  $u^t x$ . Normalize first.

These new data points have variance  $u^t \Sigma u$  if  $var(x) = \Sigma$ . Find u s.t.  $u^t \Sigma u$  is maximized which is the largest eigenvector.

#### Example: Chemicals in Italian wine

#### 3 types of wine (V1).

wine <- read.table(</pre> "http://archive.ics.uci.edu/ml/machine-learning-databases/wine/wine.data", sep=",")

	Cvs	Alcohol	Malic a	acid	Ash	Alcalinit	cy of a	ash
1	1	14.23		1.71	2.43		15	5.6
2	1	13.20	-	1.78	2.14		11	L.2
3	1	13.16	, i i i i i i i i i i i i i i i i i i i	2.36	2.67		18	3.6
4	1	14.37	-	1.95	2.50		16	5.8
5	1	13.24		2.59	2.87		21	L.0
6	1	14.20		1.76	2.45		15	5.2
	Magr	nesium To	otal phe	enols	s Flav	/anoids		
1		127		2.80	9	3.06		
2		100		2.65	5	2.76		
3		101		2.80	9	3.24		
4		113		3.85	5	3.49		
5		118		2.80	9	2.69		
6		112		3.2	7	3.39		
	Nont	flavanoid	d pheno <sup>-</sup>	ls Pi	roantł	nocyanins	Color	intensity
1			0.2	28		2.29		5.64
2			0.2	26		1.28		4.38
3			0.3	30		2.81		5.68
4			0.2	24		2.18		7.80
5			0.3	39		1.82		4.32

#### Example: Chemicals in Italian wine

sdc <- as.data.frame(scale(wine[,2:14])) # standardise</pre> wine.pca <- prcomp(sdc)</pre> # run PCA summary(wine.pca)

Importance of components:

	PC1	PC2	PC3	PC4
Standard deviation	2.169	1.5802	1.2025	0.95863
Proportion of Variance	0.362 0	9.1921 (	9.1112	0.07069
Cumulative Proportion	0.362 0	0.5541 (	9.6653	0.73599
	PCS	5 P(	C6	PC7 PC8
Standard deviation	0.92370	0.801	93 0.74	231 0.59034
Proportion of Variance	0.06563	3 0.0493	36 0.04	239 0.02681
Cumulative Proportion	0.80162	2 0.8509	98 0.89	337 0.92018
	PCS	9 PC10	9 PC	11 PC12
Standard deviation	0.53748	3 0.5009	9 0.475	17 0.41082
Proportion of Variance	0.02222	2 0.0193	3 0.017	37 0.01298
Cumulative Proportion	0.94240	0.961	7 0.979	07 0.99205
	PC13	3		
Standard deviation	0.32152	2		
Proportion of Variance	0.00795	5		
Cumulative Proportion	1.00000	9		

#### Scree plot



#### wine.pca

#### Loadings (weights) for pca

#### wine.pca\$rotation

	PC1	PC2
Alcohol	-0.144329395	-0.483651548
Malic acid	0.245187580	-0.224930935
Ash	0.002051061	-0.316068814
Alcalinity of ash	0.239320405	0.010590502
Magnesium	-0.141992042	-0.299634003
Total phenols	-0.394660845	-0.065039512
Flavanoids	-0.422934297	0.003359812
Nonflavanoid phenols	0.298533103	-0.028779488
Proanthocyanins	-0.313429488	-0.039301722
Color intensity	0.088616705	-0.529995672
Hue	-0.296714564	0.279235148
OD280/OD315 of diluted wines	-0.376167411	0.164496193
Proline	-0.286752227	-0.364902832
	PC3	PC4
Alcohol	-0.20738262 -	-0.01785630
Malic acid	0.08901289	0.53689028
Ash	0.62622390 -	-0.21417556
Alcalinity of ash	0.612080 <u>35</u>	0.06085941
Magnesium	0 13075693 -	-0 35179658

#### Biplot

biplot(wine.pca, col=c("black","red"), lwd=3)



#### PCA

- Each PC is a linear combination of the existing (original) variables.
- Each pair of PCs are orthogonal.
- PCA on unnormalized variables will lead to large loadings for variables with high variance.
- Interpretation

Now that we have the PCs - then what?

Use the new components as replacement predictors in a regression model.

#### Principal component regression

pc <- predict(wine.pca)[,1:4] # Select first 4 pcs
model <- lm(wine[,1] ~ pc) # Not super optimal model ...
tidy(model)</pre>

# A tibble:  $5 \times 5$ 

	term	estimate	std.error	statistic	p.value
	<chr></chr>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>	<dbl></dbl>
1	(Intercept)	1.94	0.0260	74.4	2.66e-133
2	pcPC1	0.319	0.0120	26.5	8.07e- 63
3	pcPC2	-0.00559	0.0165	-0.338	7.36e- 1
4	рсРСЗ	0.000991	0.0217	0.0456	9.64e- 1
5	pcPC4	0.0591	0.0272	2.17	3.15e- 2

- Virtually no limit to the number of predictors
- Correlated/noisy predictors do not undermind regression fit
- PCs carry maximum amount of variance possible