



Analisi per sottogruppi

Paolo Bruzzi, Genova

STATISTICAL

PREDETERMINED ANALYSES

METHODOLOGICAL

INTENTION TO TREAT

Analysis of RCTs: Reference criteria

STATISTICAL

All statistical analyses must be explicitly predetermined (endpoint, transformations, test, timing, subgroups)

MULTIPLICITY

Multiplicity

With an increasing number of analyses, the probability of finding, **BY CHANCE**, some noteworthy difference increases



- Five consecutive reds at the roulette wheel
- Two cases with the same inherited mutation.
- Three long-term survivors with advanced NSCLC

Multiplicity

If I look for any possible treatment effect, **BY CHANCE**, I will always find some difference:

- Overall mortality, cause-specific mortality (50 causes)
- QoL (six different domains)
- Incidence of AEs (50 possible) and favorable events (50 possible)

If, afterward, I focus on the one(s) showing a difference, I can always demonstrate that a treatment is effective (less toxic, etc.)

Possible sources of multiplicity in a clinical trial

Multiple endpoints

Transformation of endpoints

(e.g. cumulative incidence at 3, 5, 10 years vs survival curve) Statistical test

Summary effect measure

Missing data

Interim analyses

Subgroup analyses

Critical distinction: Planned multiple tests vs data-derived tests (post hoc analyses)

PLANNED MULTIPLE TESTS

- Predetermined (study protocol)
- Finite number
- Statistical correction possible

POST HOC ANALYSES

- Number potentially infinite
 - The observation of an association induces a test of significance
 - Intensive crosstabulations in search of associations
- Lack of any statistical rationale/validity

Multiplicity: General rules

- Before the start of the study, the number, time, and types of analyses are declared
 - E.g. two analyses in subjects <50 or >50 years old, or three analyses after 100, 200, and 300 events (final)
- A set of rules is established to decide if the study has led to a positive result (or to stop the study)
- These rules are built in such a way that the **overall** probability of an α error is the desired one (e.g. 5%)

Analysis of results: Strategy

- Primary analysis (p value)
- Interim analyses
- Secondary analyses
 - Other endpoints
 - Subpopulations
 - Subgroups (interactions)
 - Multivariate analyses

P R O T O C O L

For each analysis, the statistical plan establishes the statistical method to be used, when and how it will be conducted, and the decision rules (stop/go, positive/negative, p levels, etc.)

Methods to control multiplicity

Adjusted p values

Hierarchical test procedures

Closed test procedures

NOTE

Adjusted p values

Alpha spending function

Alpha split

Others

SAME MEANING

In each analysis, a p value <0.05 is used to ensure that the overall probability of a "significant" result (rejecting the null hypothesis when true) is <5%

MODIFIED SIGNIFICANCE LEVELS (all $<\alpha$) that make the overall probability of a false-positive result equal to the desired α level (usually 5%)

EXAMPLE

If:

- Desired significance level = 0.05 (5%)
- Four analyses are planned
- Study is positive if in at least one of these four analyses p <0.05/4 = 0.0125 (Bonferroni)

NB the correct formula is:

Significance level for n analyses at the overall significance level p: required p at each analysis = $1-(1-0.05)^{1/n}$

Therefore, for n=4 and p=0.05, the required p is 0.127

SINGLE-STEP METHODS

Examples: Bonferroni, Simes, Dunnett...

STEPWISE METHODS

The rejection or non-rejection of a particular hypothesis may depend on the decision made on other hypotheses

Examples: Holm, Hochberg, step-down Dunnett...

ALL METHODS INVOLVE A LOSS OF STATISTICAL POWER

- Limit the number of analyses
- Increase the sample size
- Use hierarchical procedures
- The Bonferroni method is the most conservative (least powerful)
 - Other more complex methods generally used

Methods to control multiplicity

Adjusted p values

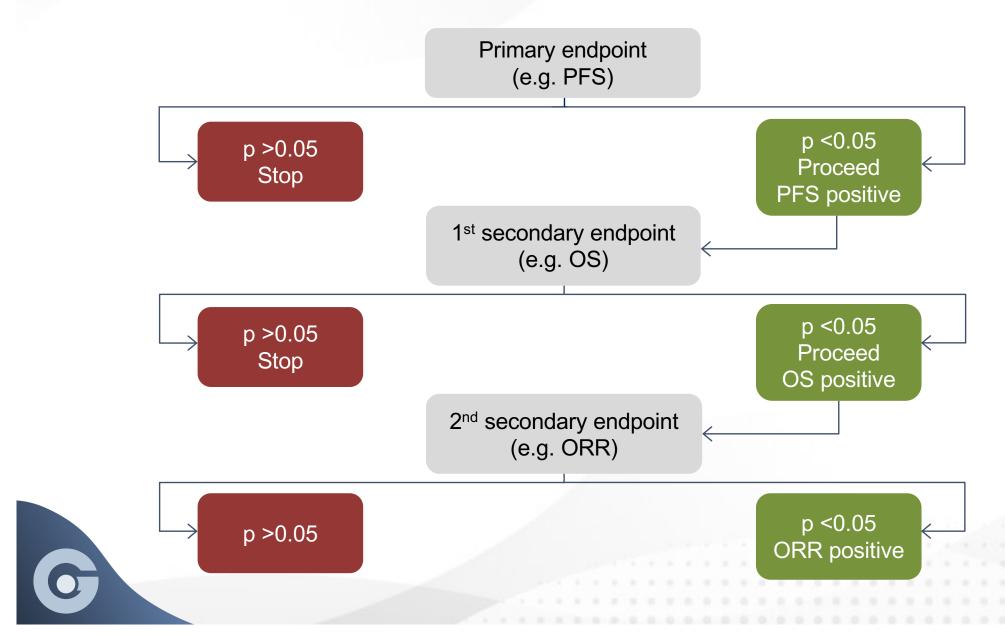
Hierarchical test procedures

Closed test procedures

Hierarchical test procedures

- Hypotheses are ordered in sequence and tested at level α until the first non-rejection
- In practice, first test:
 - If p >0.05 → stop (negative study)
 - If positive → second test
- NO CORRECTION OF THE SIGNIFICANCE LEVEL IS REQUIRED
- Sequence based on relevance, power, plausibility, etc.

Hierarchical test procedures (example)



Hierarchical test procedures

- No loss of power for the first analysis
- Frequently used for multiple endpoints
- Risk of missing relevant treatment effects if the order of tests is incorrect (e.g. OS then PFS)

Methods to control multiplicity

Adjusted p values

Hierarchical test procedures

Closed test procedures

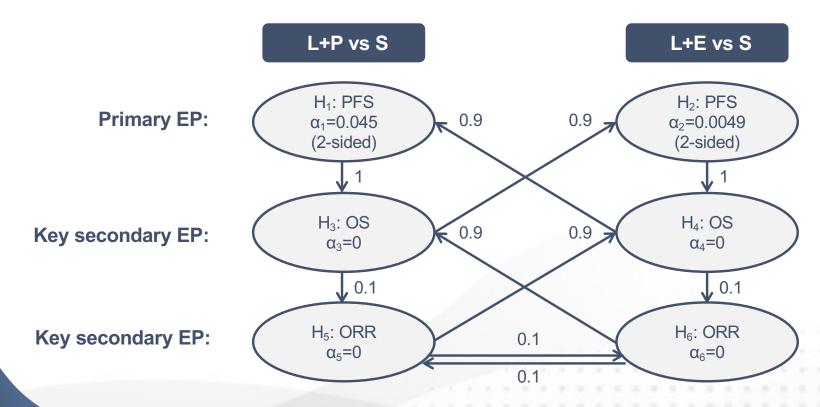
Closed test procedures

- General principle to build procedures for multiple tests
- Used to protect the α error while maintaining efficiency (reduced loss of power)
- Many of the previously mentioned procedures (e.g. Holm, hierarchical) are based on this principle
- Widely used: Maurer & Bretz graphical method

Statistical Plan AMENDMENT 07 (2020)

MULTIPLICITY

To adjust for multiplicity and control the overall FWER, the graphical approach of Maurer and Bretz (Maurer et al., 2013) will be used in the primary endpoint of PFS and the key secondary efficacy endpoints (OS and ORR). No multiplicity adjustment will be made for other secondary endpoint analyses.



Analysis of results: Strategy

- Primary analysis (p value)
- **►** Interim analyses
- Secondary analyses
 - Other endpoints
 - Subpopulations
 - Subgroups (interactions)
 - Multivariate analyses
- Unplanned analyses: merely exploratory aims
 - Used to plan other studies

P R O T O C O I

Possible sources of multiplicity in a clinical trial

Multiple endpoints

Transformation of endpoints

(e.g. cumulative incidence at 3, 5, 10 years vs survival curve) Statistical test

Summary effect measure

Missing data

Interim analyses

Subgroup analyses

Subgroup analyses

The aim of these is to provide information on the opportunity to treat different groups of patients differently, informing the development of:

PERSONALIZED THERAPIES

Leading toward the era of



PRECISION MEDICINE



Prognostic and predictive factors

Subgroup analyses can inform potential prognostic and predictive factors

PROGNOSTIC FACTORS

- Predict outcome (with the same treatment)
- Do not require a randomized trial to identify
- Used in clinical decision-making (informing risk/benefit and cost/benefit)

FOR EXAMPLE:

Nodal status in early-stage breast cancer

- Strong prognostic effect (HR: 2)
- All adjuvant therapies have the same effect, regardless of nodal status

PREDICTIVE FACTORS

- Predict the efficacy of the treatment in different patients
- Identified solely by subgroup analyses from randomized trials

FOR EXAMPLE:

- Hormone receptors: efficacy of hormonal therapy in breast cancer
- PD-L1 expression: efficacy of immunotherapy in solid tumors
- Tumor grade (differentiation): efficacy of chemotherapy in NHL



Issues with subgroup analyses: Methodological

There are several methodological factors that can affect the validity of a subgroup analysis:

Retrospective vs prospective

Not very important

Planned vs unplanned

- Must always be included in the statistical plan
- If not planned = scientific exercise

Bias

- Blinded classification/analyses
- Selection of compliers, responders, and treated patients (?)
- Always compare subgroup HRs with those from the overall population

STUDY PROTOCOL!

Issues with subgroup analyses: Statistical

The common statistical issues with subgroup analyses are:

IMPROPER SIGNIFICANCE TESTING

MULTIPLICITY

Widely used method (wrong)

The test for significance is repeated in each subgroup at the conventional significance level

Hypothetical trial with 120 patients/arm

72 responses (60%) exp. therapy vs 48 (40%) control **p <0.002**

3 SUBGROUPS:

<30 years	p=0.01
30–50 years	p=0.2 n.s.
>50 years	p=0.2 n.s.

Conclusions?

3 SUBGROUPS:

<30 years	48/80 (60%) vs 32/80 (40%) p=0.01
30–50 years	12/20 (60%) vs 8/20 (40%) p=0.2
>50 years	12/20 (60%) vs 8/20 (40%) p=0.2

3 SUBGROUPS:

<30 years	p=0.07 n.s.
30–50 years	p=0.07 n.s.
>50 years	p=0.07 n.s.

Conclusions?

3 SUBGROUPS:

<30 years	24/40 (60%) vs 16/40 (40%) p=0.07 n.s.
30–50 years	24/40 (60%) vs 16/40 (40%) p=0.07 n.s.
>50 years	24/40 (60%) vs 16/40 (40%) p=0.07 n.s.

Subgroup analyses

SUBGROUP SPECIFIC P: MEANINGLESS E MISLEADING

- Correct analysis: Test of interaction (Heterogeneity of the effect)
 - New null hypothesis: The effect is the same in all subgroups
 - The observed variation in the effect is compared to that expected by chance alone (small subgroups: large variations)

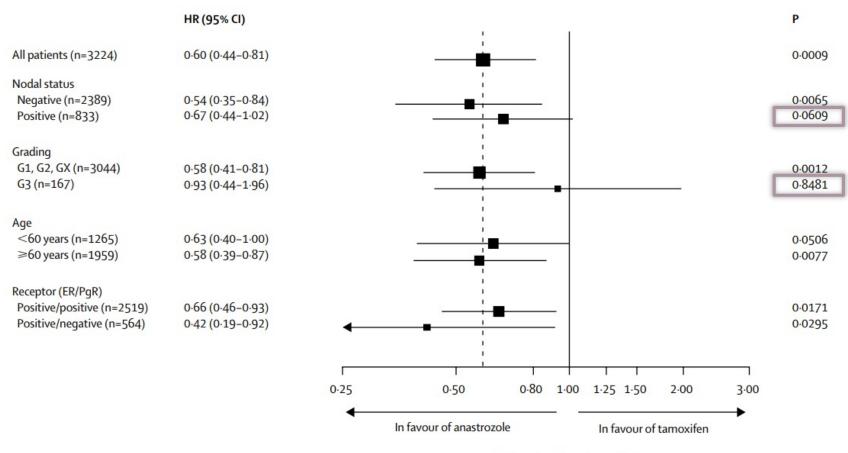
<30 years	60% vs 40%	OR=2.25
30–50 years	60% vs 40%	OR=2.25
>50 years	60% vs 40%	OR=2.25

No evidence of interaction p=1

<30 years	50% vs 50%	OR=1
30–50 years	60% vs 40%	OR=2.25
>50 years	70% vs 30%	OR=5.3

Evidence of interaction p < 0.05

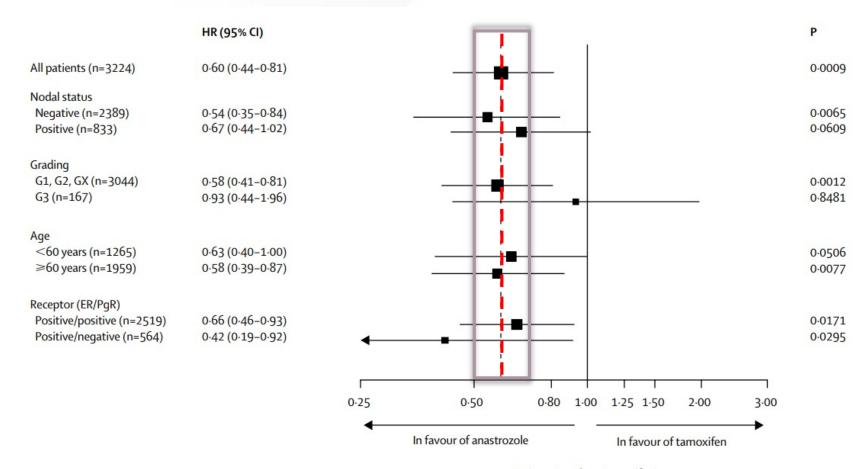
Improper significance testing



HR (anastrozole vs tamoxifen)

In this forest plot, the subgroup-specific p values are incorrect

Proper subgroup analysis



HR (anastrozole vs tamoxifen)

Including the overall treatment effect level makes it easier to see if the effect in a subgroup differed significantly from the overall treatment effect

Improper significance testing

A **TEST OF INTERACTION** is required to correctly analyze subgroup analyses; this assesses the heterogeneity of the treatment effect

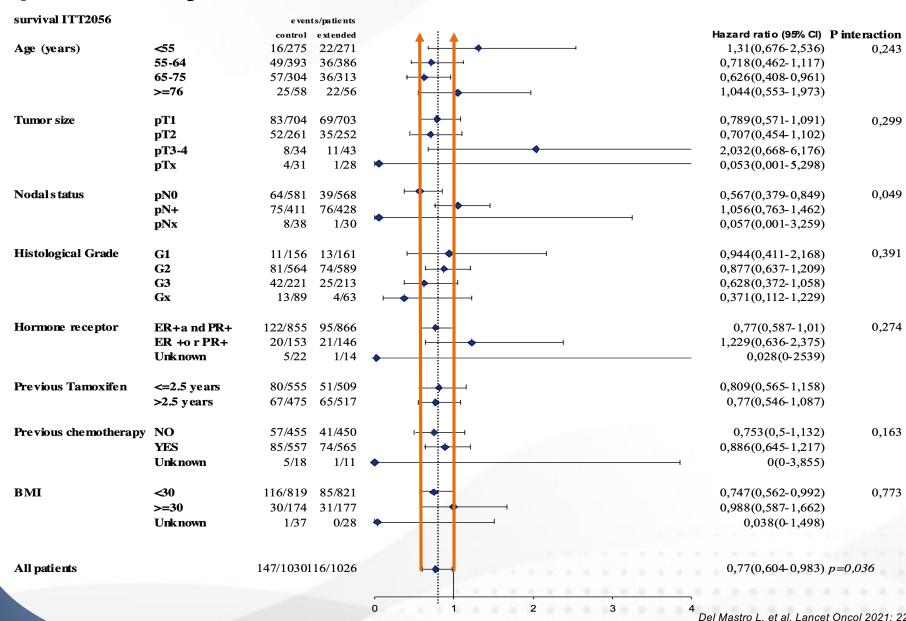
TEST OF INTERACTION

New null hypothesis: The treatment effect is the same across all subgroups

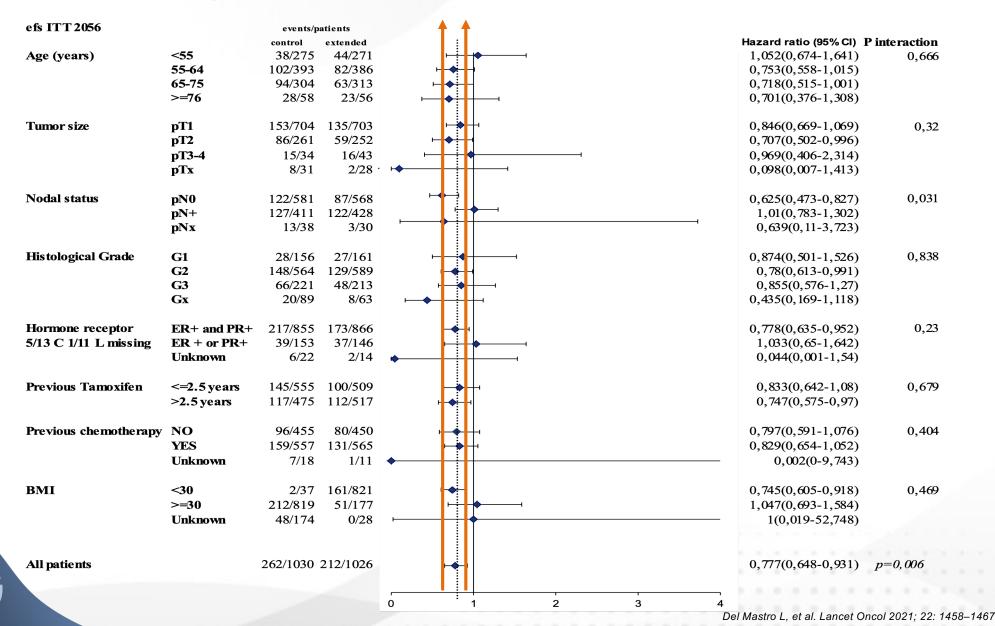
- The observed variation in the treatment effect is compared with that expected by chance alone
- Small subgroups, large variations

Subgroup-specific p values can be misleading and meaningless

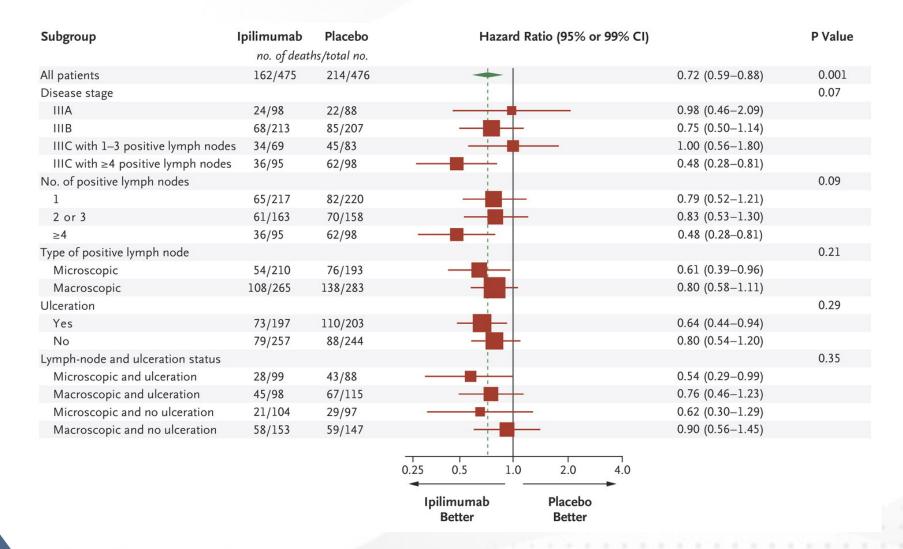
Appropriate analysis – Overall survival



Appropriate analysis – Event-free survival



Appropriate significance testing



Subgroup analysis methodology

In modern trial design, subgroup analyses are carefully designed to ensure validity of the results

SUBGROUP ANALYSIS METHODOLOGY

- Careful planning to prevent selection and assessment biases
- Test for interaction: H0=the (lack of) effect is the same in all subgroups
 - No subgroup-specific p values should be calculated
- Controlling for multiplicity:
 - Planned vs post hoc analyses
 - Exploratory vs confirmatory analyses
 - p value corrections

Larger data sets allow more POWERFUL subgroup analyses

Strategies for addressing multiplicity

POSITIVE PRIMARY RESULTS

If all preceding statistical analyses are positive, planned subgroup analyses remain valid

NEGATIVE PRIMARY RESULTS

If any preceding statistical analyses were negative, **planned subgroup analyses will be invalid** if not corrected for multiplicity

α split required

e.g. 2% × primary analysis 1% each × 3 interaction tests

STATISTICAL ANALYSIS PLAN!

Conclusions

- Most modern trials are adequately designed to avoid problems of multiplicity in the statistical analyses
- However, these analyses are focused on statistical testing rather than estimation. Modern statistical and medical sciences are moving away from p values and are more interested in estimating the effects of the treatments

Multiplicity-corrected, statistically significant treatment effect estimates are biased, as they represent overestimations of the true treatment effects