

臨床試験における PRO/QOLデータの評価方法

山口 拓洋

東北大学大学院医学系研究科

東北大学病院臨床試験データセンター



Division of Biostatistics, Tohoku
University Graduate School of
Medicine

2010-2019
9th Anniversary



1

内容

- これまでのご講演をベースに
- SISAQOL-IMI CONSORTIUM RECOMMENDATION を参照しつつ、PRO研究(臨床試験の目的の一つとしてPRO評価を行う)の統計学的事項について導入的な説明をおこなう

International standards for the analysis of quality-of-life and patient-reported outcome endpoints in cancer randomised controlled trials: recommendations of the SISAQOL Consortium



Corneel Coens, Madeline Pe*, Amylou C Dueck, Jeff Sloan, Ethan Basch, Melanie Calvert, Alicyn Campbell, Charles Cleeland, Kim Cocks, Laurence Collette, Nancy Devlin, Lien Dorme, Hans-Henning Flechtner, Carolyn Gotay, Ingolf Griebisch, Mogens Groenvold, Madeleine King, Paul G Kluetz, Michael Koller, Daniel C Malone, Francesca Martinelli, Sandra A Mitchell, Jammbe Z Musoro, Daniel O'Connor, Kathy Oliver, Elisabeth Piault-Louis, Martine Piccart, Chantal Quinten, Jaap C Reijneveld, Christoph Schürmann, Ashley Wilder Smith, Katherine M Soltys, Martin J B Taphoorn, Galina Velikova, Andrew Bottomley; for the Setting International Standards in Analyzing Patient-Reported Outcomes and Quality of Life Endpoints Data Consortium*

Patient-reported outcomes (PROs), such as symptoms, function, and other health-related quality-of-life aspects, are increasingly evaluated in cancer randomised controlled trials (RCTs) to provide information about treatment risks, benefits, and tolerability. However, expert opinion and critical review of the literature showed no consensus on optimal methods of PRO analysis in cancer RCTs, hindering interpretation of results. The Setting International Standards in Analyzing Patient-Reported Outcomes and Quality of Life Endpoints Data Consortium was formed to establish PRO analysis recommendations. Four issues were prioritised: developing a taxonomy of research objectives that can be matched with appropriate statistical methods, identifying appropriate statistical methods for PRO analysis,

Lancet Oncol 2020; 21: e83-96

*Joint first authors

European Organisation for Research and Treatment of Cancer, Brussels, Belgium
(C Coens MSc, M Pe PhD, L Collette PhD, L Dorme MSc, F Martinelli MSc, J Z Musoro PhD,

推奨ポイント

- 適切な統計解析方法を考慮したうえでの研究目的の分類
- 特定のPRO研究目的を達成するのに適した統計解析方法
- 欠測データに関する統計用語の標準化
- 欠測データの取り扱い

適切な統計手法を考慮したうえでの研究目的の分類

- 検証か探索か？
- 優越性か同等性/非劣性か？
- 個人内かグループ内か、どちらの変化に興味があるか？
 - 個人内: 各個人が臨床的に意味のある変化(事前に設定)があったかどうか？
 - グループ内: 各グループが平均的に意味のある変化(事前に設定)があったかどうか？

研究目的と評価項目

改善

- 改善するまでの時間
- 特定の時点における改善の程度、スコア
- 特定の時点における改善者の割合

悪化

- 悪化するまでの時間
- 特定の時点における悪化の程度、スコア
- 特定の時点における悪化者の割合

安定状態

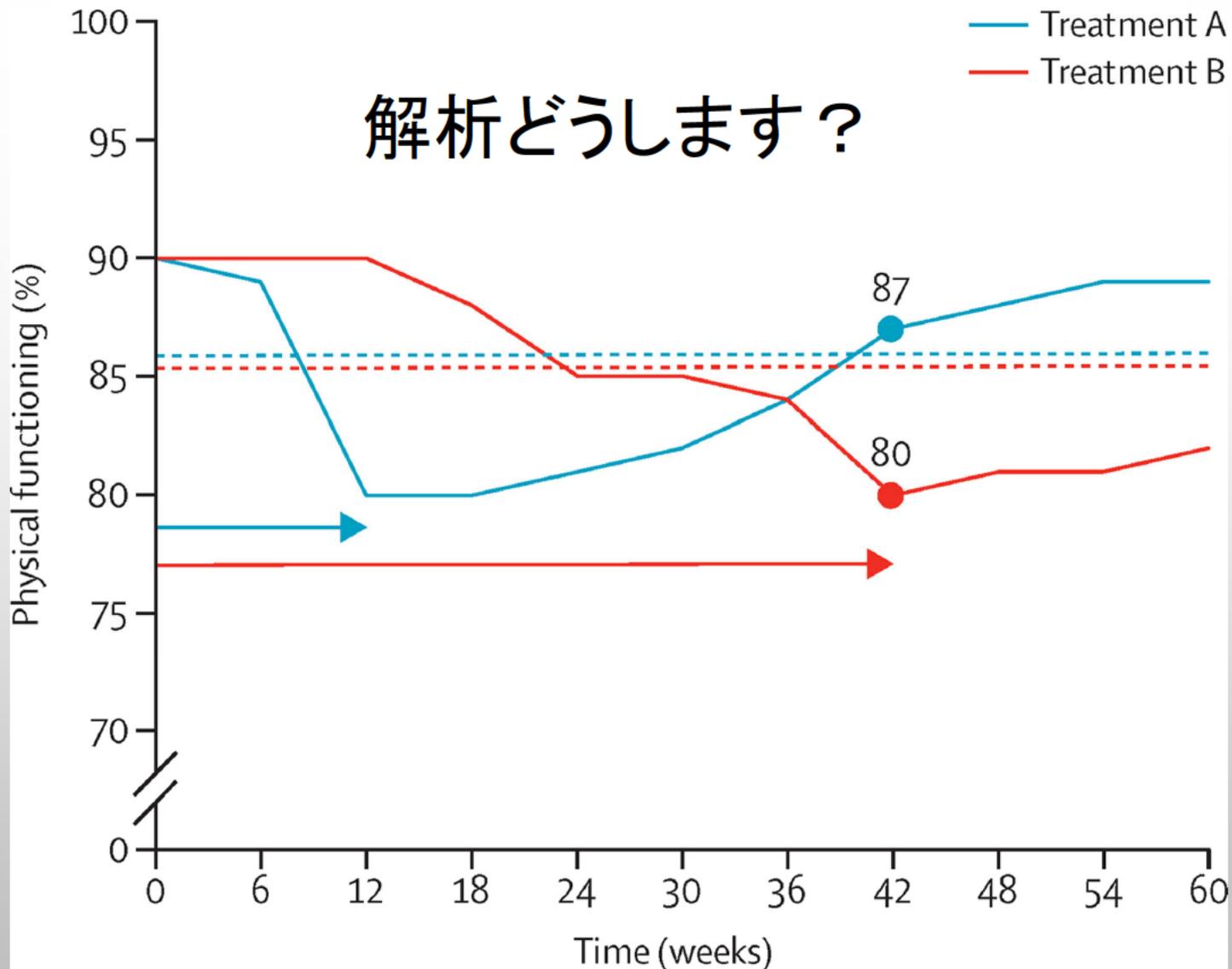
- 安定状態になるまでの時間
- 特定の時点における安定状態者の割合

全体の効果

- 要約指標
- 経時的な推移（反応パターン、プロファイル）

特定のPRO研究目的を達成するのに適した 統計解析方法

- イベント発生までの時間
 - データの記述、ログランク検定、COX回帰モデルなど
- 特定の時点におけるスコア
 - データの記述、 t 検定、重回帰モデルなど(欠測なければ)、線型混合モデルなど
- 特定の時点における割合
 - データの記述、カイ2乗検定、M-H検定、ロジスティック回帰モデルなど(欠測なければ)、ロジスティック混合効果モデルなど(見解がまとまっていない)
- 要約指標
 - 見解がまとまっていない
- 反応パターン、プロファイル
 - データの記述、線型混合モデルなど



The Lancet Oncology 2016 17, e510-e514 DOI: (10.1016/S1470-2045(16)30510-1)

考えられる解析方法

- 特定の時点におけるスコアの平均値のグループ間比較、有効/無効(カットオフ値の定義が重要)の割合のグループ間比較
- ベースラインと特定の時点におけるスコアの個人内変動(変化量、変化率)の平均値のグループ間比較
- 多時点におけるスコアの推移のグループ間比較
- 多時点におけるスコアから算出される曲線下面積(AUC)のグループ間比較
- 特定のイベントが発生(例:スコア2ポイント減少)するまでの時間のグループ間比較
- ...

まずはグラフ表示

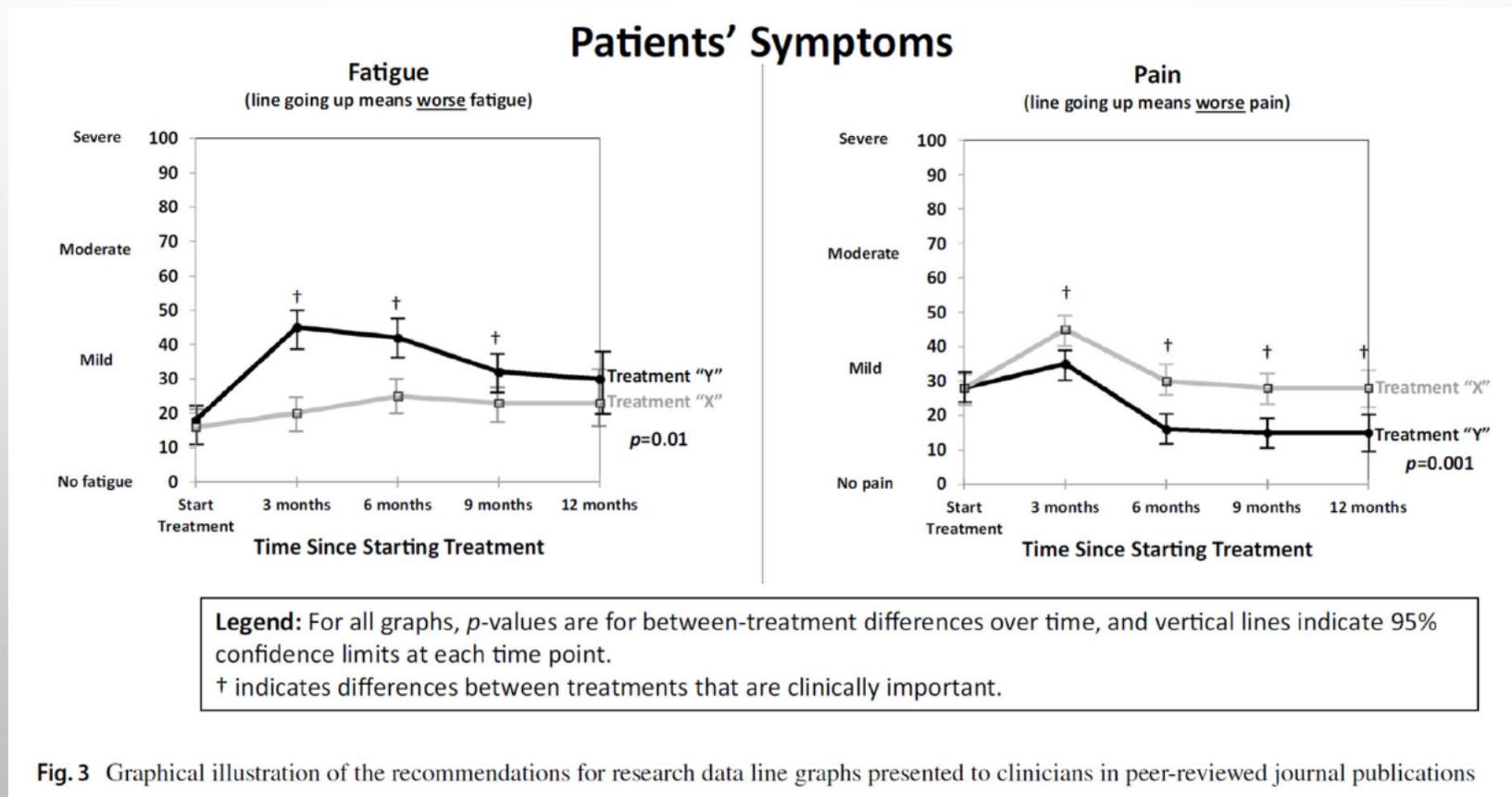
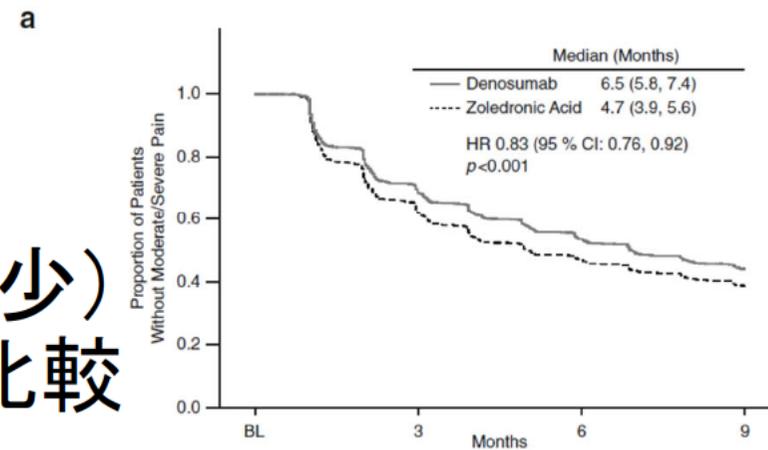


Fig. 3 Graphical illustration of the recommendations for research data line graphs presented to clinicians in peer-reviewed journal publications

特定のイベント(例:スコア2ポイント減少)が発生するまでの時間のグループ間比較

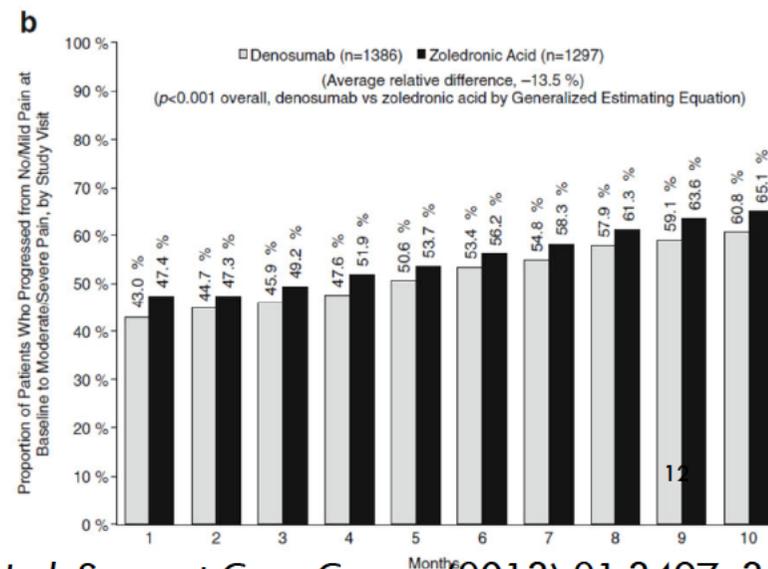
- イベントの定義(カットポイント等、臨床的な妥当性も含む)を明確にする必要
- 各群におけるイベント発生率、発生割合の推定とグループ間比較(ハザード比の推定、ログランク検定)

Fig. 2 Onset of moderate to severe pain. a Time to first report of moderate or severe pain among patients with no or mild pain at baseline. b Proportion of patients at each study visit who progressed from a baseline score of no or mild pain to moderate or severe pain



Risk Set:

Zoledronic Acid 4 mg Q4W	1297	712	506	386
Denosumab 120 mg Q4W	1386	859	621	480



Von Moos et al. Support Care Cancer (2013) 21:3497–3507.

特定の時点におけるスコアの平均値、有効/無効の割合のグループ間比較

- 主要な評価時点を事前に設定
- 各群における平均値の推定とグループ間比較(平均値の差の推定、 t 検定)
 - スコアなどの連続データ
- 各群における割合の推定とグループ間比較(割合の差/比の推定、カイ2乗検定)
 - 有効/無効、改善/非改善など
 - 定義を事前に規定
- 主要評価項目を「NRSの平均値」とした場合に、副次評価項目に「NRSが2点以上下がった患者の割合」「NRSが33%低下した患者の割合」を加えるなど
 - 臨床的に意味のある差(MINIMUM IMPORTANT DIFFERENCE; MID)について考察する必要
 - カットポイントの選択と感度分析

ベースラインと特定の時点におけるスコアの 個人内変化(変化量、変化率)の平均値の グループ間比較

- (主に連続データでの)ベースライン値の取り扱い
- 無視(ランダム化試験など)
 - 前述
- 個人ごとにベースラインとの差(変化量)、比(変化率)の算出、
各群における平均値の推定とグループ間比較(平均値の差の推定、 t 検定)
- ベースライン値を共変量(調整変数)として回帰分析

多時点におけるスコアから算出される 曲線下面積 (AUC) のグループ間比較

- 個人ごとに
AUC、あるいは、
単位時間あたりの
平均スコアの算出、
各群における
平均値の推定と
グループ間比較
(平均値の差の推定、
*t*検定)
- どの時点までを
用いるか？

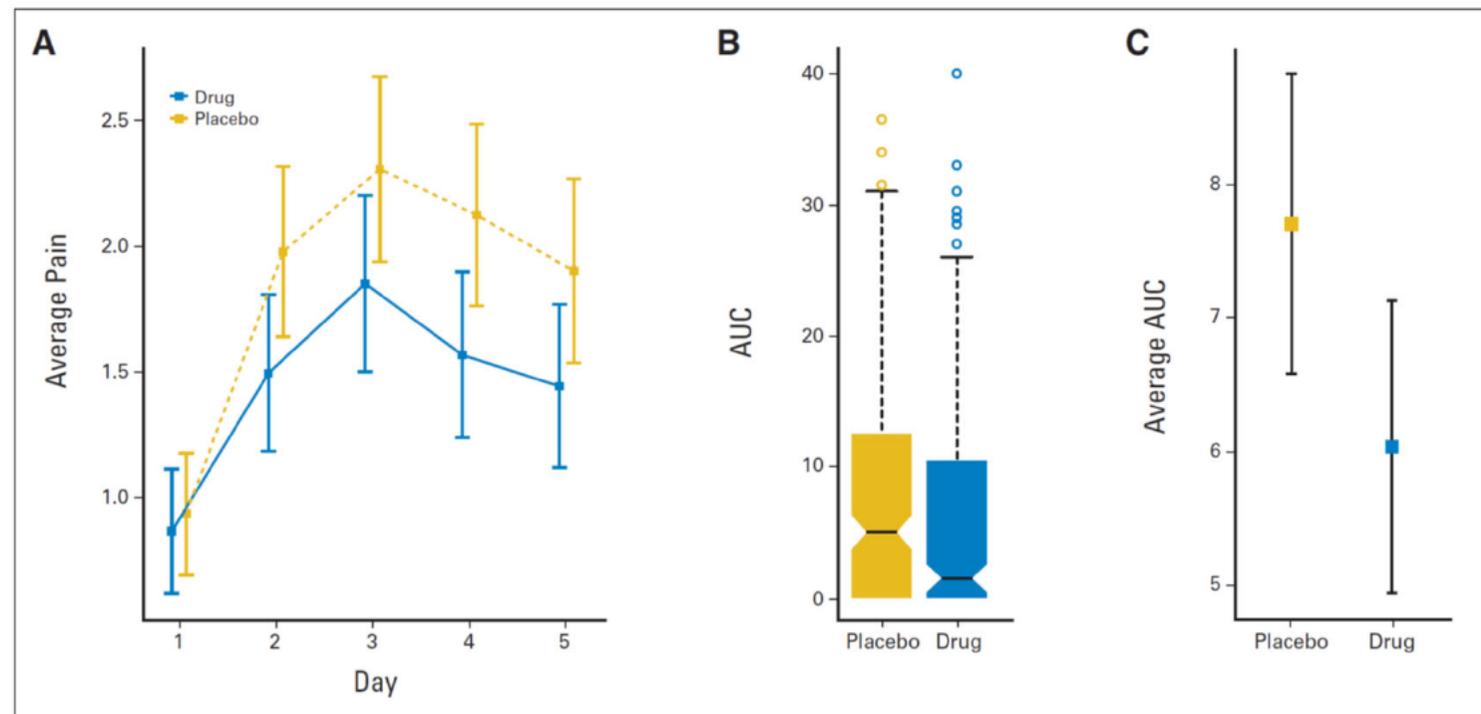


Fig 3. (A) Mean bone pain by day. Mean across patients and 95% CIs for daily bone pain, down by day and arm, on a scale of 0 to 10. Day 1 denotes the day of pegfilgrastim injection. (B) Area under the curve (AUC) distributions. Notched box plots of AUC for bone pain. Nonoverlapping notches indicate a statistically significant difference in the medians. (C) Mean AUC bone pain. Mean and 95% CIs for the average AUC. The possible range of AUC in (B) and (C) is 0 to 40.

Kirshner *et al.* *J Clin Oncol* (2012) 30:1974-1979.

多時点におけるスコアの推移のグループ間比較

- 平均値や割合の経時的な推移(トレンド)の推定とグループ間比較
- 経時データ解析
 - 混合効果モデル
MIXED EFFECT MODEL
- 個人内で繰り返し測定されたデータの相関を考慮

乳房再建術後のQOLを調査する多施設共同研究

研究の背景
乳房再建は患者さんのQOL(Quality of life, 生活の質)や満足度を向上させるための術式である。欧米ではPRO(Patient-reported outcome, 患者報告アウトカム)を用いたエビデンス作りが進められているが、日本ではまだ調査報告が少ない。

日本語版BREAST-Qの開発 (乳房再建用のPRO) 欧米と日本の違い (文化・医療・女性の体型)
↓
日本独自のデータ収集が必要

研究内容

SAQLA study
Satisfaction And Quality of Life After breast reconstruction
(一次乳房再建術後の満足度とQOLを調査する多施設前向きコホート研究)

- ✓ 対象: 7施設で乳がんに対する乳房全摘と一次乳房再建を行う患者さん
- ✓ 電子アンケートを用いてPROを経時的に評価する。

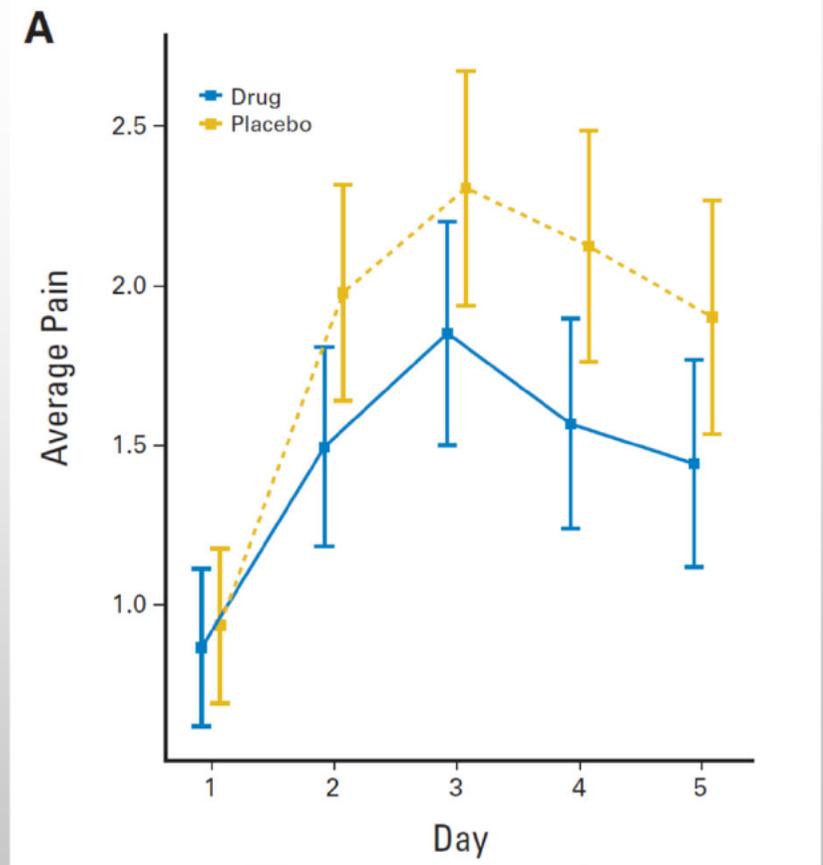
BREAST-Q (乳房の満足・心理社会的健康感・胸部と上肢の健康感・性的健康感・ドナーサイト(背骨や腰部)の健康感・放射線治療の影響)、SF-8 (身体的健康・精神的健康)、乳房再建の負担感など

- ✓ 主に術式による違いについて検討する。

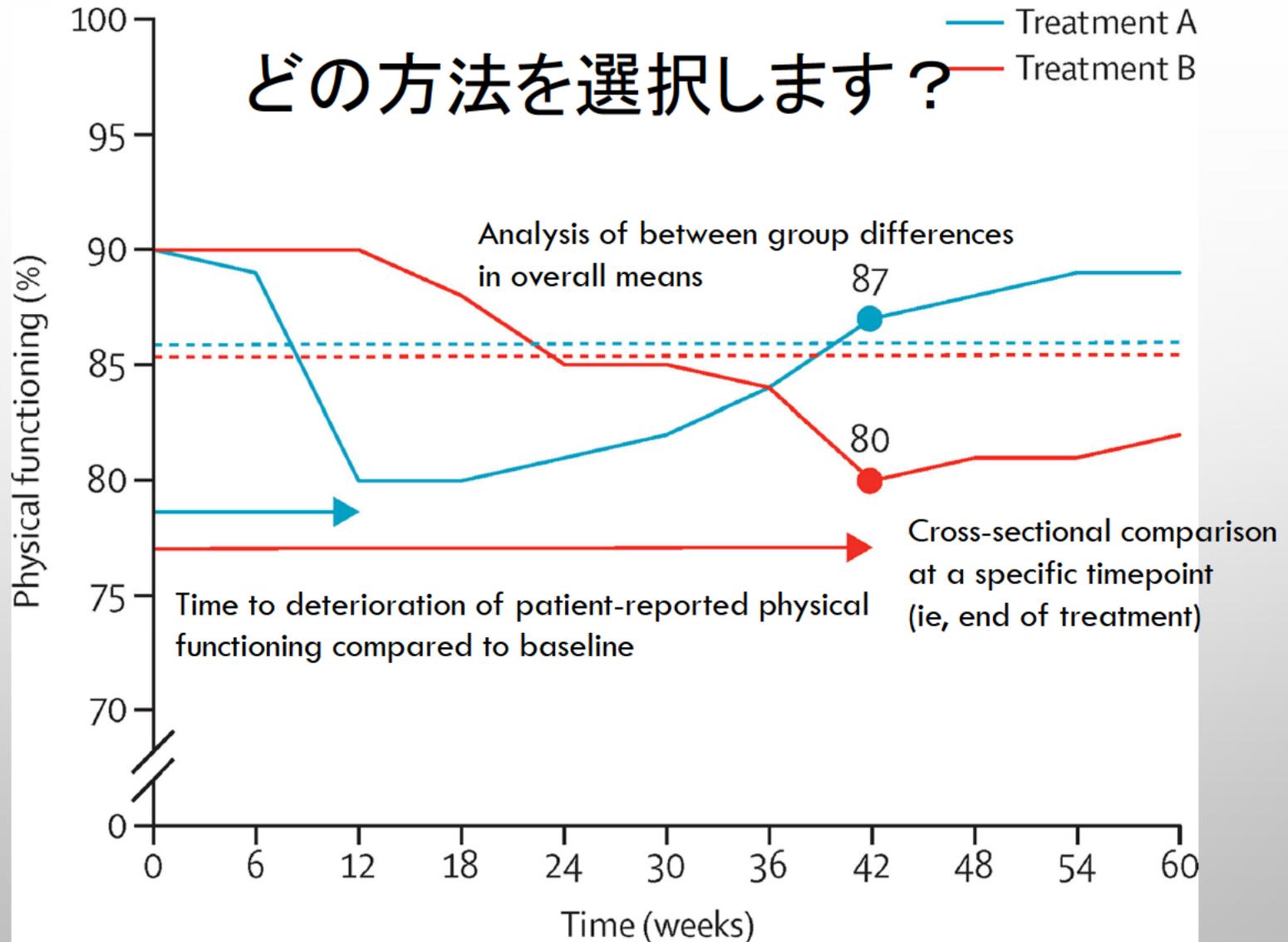
広背筋皮弁 腹部皮弁 人工乳房

研究代表者
岡山大学 木股 敬裕

発表者
岡山大学 雑賀 美帆



どの方法を選択します？



The Lancet Oncology 2016 17, e510-e514 DOI: (10.1016/S1470-2045(16)30510-1)

PROデータの特徴

- 多くの質問票は多次元であるため、複数のスコア(各ドメインスコア、総合スコアなど)などを生じうる
- 複数の時点で通常評価される

解析ストラテジー

- (研究目的、用いる評価尺度、データ収集方法、、、)
- 臨床的な意味づけ、どのような治療効果を推定したいか？
- PRO研究目的と評価項目
- MID
- 目的を達成するのに適した統計解析方法、統計学的な性能(統計家に相談)

Japan Clinical Oncology Group

ポリシー No. 30

タイトル：PRO/QOL 研究

適用範囲：

研究グループ、PRO/QOL 研究事務局、プロトコール審査委員会、データセンター/運営事務局

PRO/QOL 研究 PRO/QOL Research

6.5. 統計学的事項

6.5.1. PRO/QOL 研究における統計解析

QOL 尺度などの多くの質問票は多次元であるため、複数のスコア（各ドメインスコア、総合スコアなど）などを生じうる。さらに、PRO/QOL 評価は複数の時点で通常行われる。スコアの経時変化などを確認するためのグラフ表示が重要であるとともに、評価項目、解析方法などを事前に規定しておく必要がある。評価項目には、特定の時点におけるスコアや有効/無効（定義が重要、例えばスコアの 33% の低下）、特定のイベント（例えば、スコアの 2 点の低下）が発生するまでの時間、観察期間全体を通してのスコアの推移や曲線下面積（Area Under the Curve : AUC）などが考えられる。さらには、臨床的に意味のある差（MID（3. 用語説明 参照））について考察する必要がある。状況に応じて、解析の多重性についても考慮する。

研究のアウトプットを見据える

- 評価項目（観察測定スケジュール含む）と対応させて
- **TLF (TABLE, LIST, FIGURE)** を考える
 - 論文、報告書等にそのまま載せる
 - 解析計画書に付ける場合が多い
 - 記載例（モックアップ）とも呼ばれる

PROのデータ、例えば毎週データを取りたくなると思いますが
解析に使用しないデータは取る必要なし！
アウトプットを考えてデータを取る必要がある

観察測定スケジュール

6.1 Treatment Phase for the MK-3475 Arm

Treatment Cycle / Scheduled Time Scheduling Window (Days): ²	Screening (Visit 1)	Treatment Cycles ¹														End of Treatment
	-42 to -1	1	2	3	4	5	6	7	8	9	10	11	12	13	14 through 35	Discontinuation Visit/Safety follow up Visit
			±3	±3	±3	±3	±3	±3	±3	±3	±3	±3	±3	±3	±3	
Administrative Procedures																
Informed Consent	X															
Informed Consent for Future Biomedical Research (optional)	X															
Inclusion/Exclusion Criteria	X															
Subject Identification Card	X															
Demographics and Medical History	X															
Prior and Concomitant Medications	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
NSCLC Disease Details and Prior Treatment	X															
Clinical Procedures / Assessments																
Review Adverse Events		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Full Physical Examination	X															X
Directed Physical Examination		X	X	X	X	X	X	X	X	X	X	X	X	X	X	
Vital Signs and Weight ¹⁶	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
12-Lead ECG	X															
ECOG Performance Status	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Patient Reported Outcomes (PRO)																
EuroQol EQ-5D ¹⁵		X	X	X			X			X			X		X ¹⁵	X
EORTC QLQ-C30 ¹⁵		X	X	X			X			X			X		X ¹⁵	X
EORTC QLQ LC-13 ¹⁵		X	X	X			X			X			X		X ¹⁵	X

PEMBROLIZUMAB VERSUS CHEMOTHERAPY FOR PD-L1-POSITIVE NON-SMALL-CELL LUNG CANCER (KEYNOTE-024)

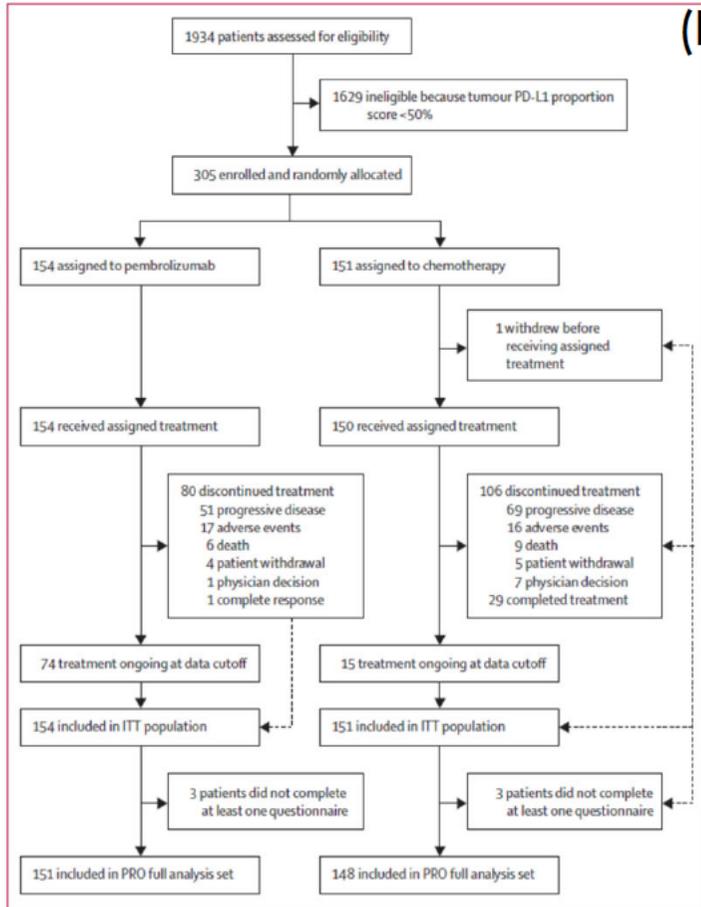
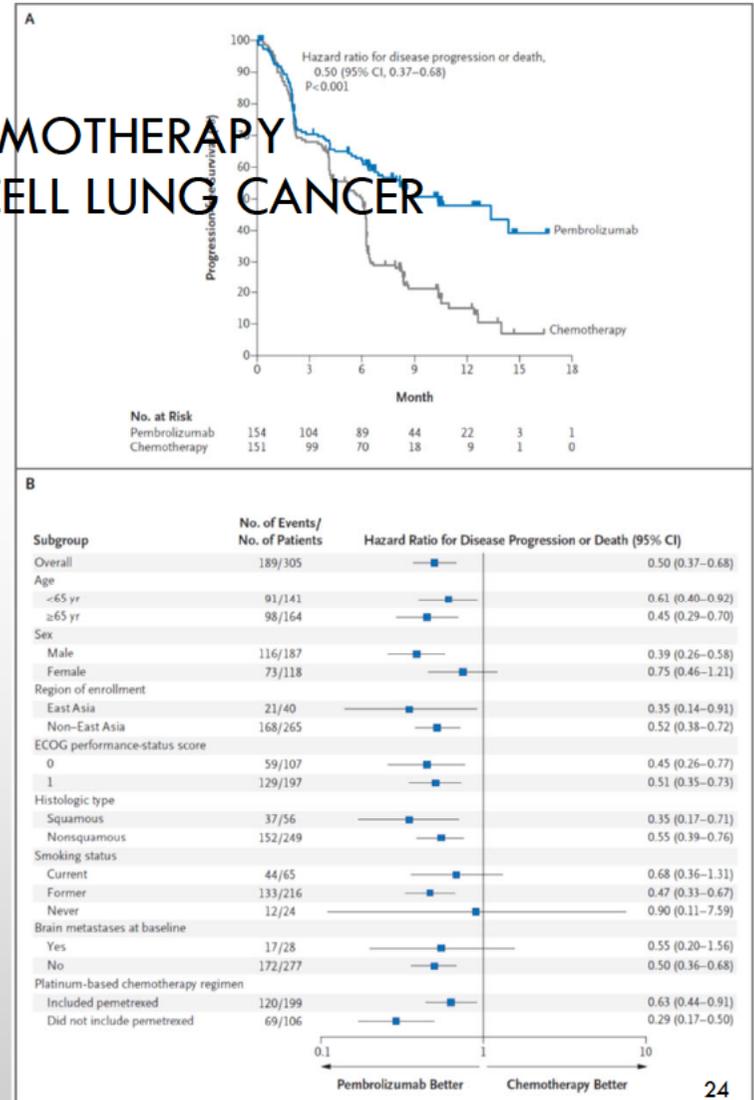


Figure 1: Trial profile
Data cutoff date: May 9, 2016. PD-L1=programmed cell death-ligand 1. ITT=intention-to-treat. PRO=patient-reported outcome. Reproduced from Reck and colleagues,¹ by permission of Massachusetts Medical Society.



Brahmer JR et al. *Lancet Oncol* 2017; 18: 1600-9.
Reck M et al. *N Engl J Med* 2016; 375: 1823-33.

- PATIENTS WITH ADVANCED NSCLC HAVE A HIGH BURDEN OF SYMPTOMS, SUCH AS FATIGUE, COUGH, DYSPNOEA, ANOREXIA, WEIGHT LOSS, AND PAIN THAT CAN HAVE A SUBSTANTIAL NEGATIVE EFFECT ON HEALTH-RELATED QUALITY OF LIFE (HRQOL) AND FUNCTIONING. THEREFORE, THE EFFECT OF NOVEL TREATMENTS ON SYMPTOM CONTROL AND HRQOL NEEDS TO BE CONSIDERED ALONGSIDE SURVIVAL.

- WE EVALUATED PROS AS PRESPECIFIED EXPLORATORY ENDPOINTS. THE TWO KEY PROS WERE CHANGE FROM BASELINE TO WEEK 15 IN QLQ-C30 GHS/QOL SCORE AND TIME TO DETERIORATION OF THE COMPOSITE OF COUGH, CHEST PAIN, AND DYSPNOEA IN THE QLQ-LC13 (DEFINED AS TIME TO FIRST ONSET OF A 10-POINT OR GREATER DECREASE FROM BASELINE IN COUGH, CHEST PAIN, OR DYSPNOEA, CONFIRMED BY A SECOND ADJACENT 10-POINT OR GREATER DECREASE FROM BASELINE IN ANY OF THESE THREE SYMPTOMS); PATIENTS WITH NO CONFIRMED DECREASE FROM BASELINE WERE CENSORED AT THE DATE OF THEIR LAST OBSERVATION.

- WE CHOSE THE WEEK 15 TIMEPOINT FOR COMPARISONS WITH BASELINE WHEN PATIENTS IN BOTH GROUPS WERE STILL ON TREATMENT TO MINIMISE LOSS OF DATA DUE TO DEATH OR DISEASE PROGRESSION.

	Pembrolizumab group (n=151)	Chemotherapy group (n=148)
QLQ-C30		
Baseline	145 (96%)	137 (93%)
Week 3	127 (84%)	122 (82%)
Compliance	127/144 (88%)	122/138 (88%)
Week 6	120 (79%)	110 (74%)
Compliance	120/138 (87%)	110/131 (84%)
Week 15	109 (72%)	92 (62%)
Compliance	109/129 (84%)	92/117 (79%)
Week 24	98 (65%)	75 (51%)
Compliance	98/111 (88%)	75/92 (82%)
QLQ-LC13		
Baseline	145 (96%)	136 (92%)
Week 3	127 (84%)	122 (82%)
Compliance	127/144 (88%)	122/138 (88%)
Week 6	120 (79%)	110 (74%)
Compliance	120/138 (87%)	110/131 (84%)
Week 15	109 (72%)	91 (61%)
Compliance	109/129 (84%)	91/117 (78%)
Week 24	98 (65%)	75 (51%)
Compliance	98/111 (88%)	75/92 (82%)
EQ-5D-3L		
Baseline	144 (95%)	137 (93%)
Week 3	127 (84%)	122 (82%)
Compliance	127/144 (88%)	122/138 (88%)
Week 6	120 (79%)	110 (74%)
Compliance	120/138 (87%)	110/131 (84%)
Week 15	108 (72%)	92 (62%)
Compliance	108/129 (84%)	92/117 (79%)
Week 24	98 (65%)	75 (51%)
Compliance	98/111 (88%)	75/92 (82%)

Compliance was defined as the proportion of patients who completed the indicated questionnaire among those who were expected to complete it at each timepoint (excluding those patients missing by design). QLQ=European Organisation for the Research and Treatment of Cancer Quality of Life Questionnaire. C30=Core 30 items. LC13=Lung Cancer 13 items. EQ-5D-3L=European Quality of Life 5 Dimensions-3 Level. PRO=patient-reported outcome.

Table 1: Completion and compliance for the PRO instruments

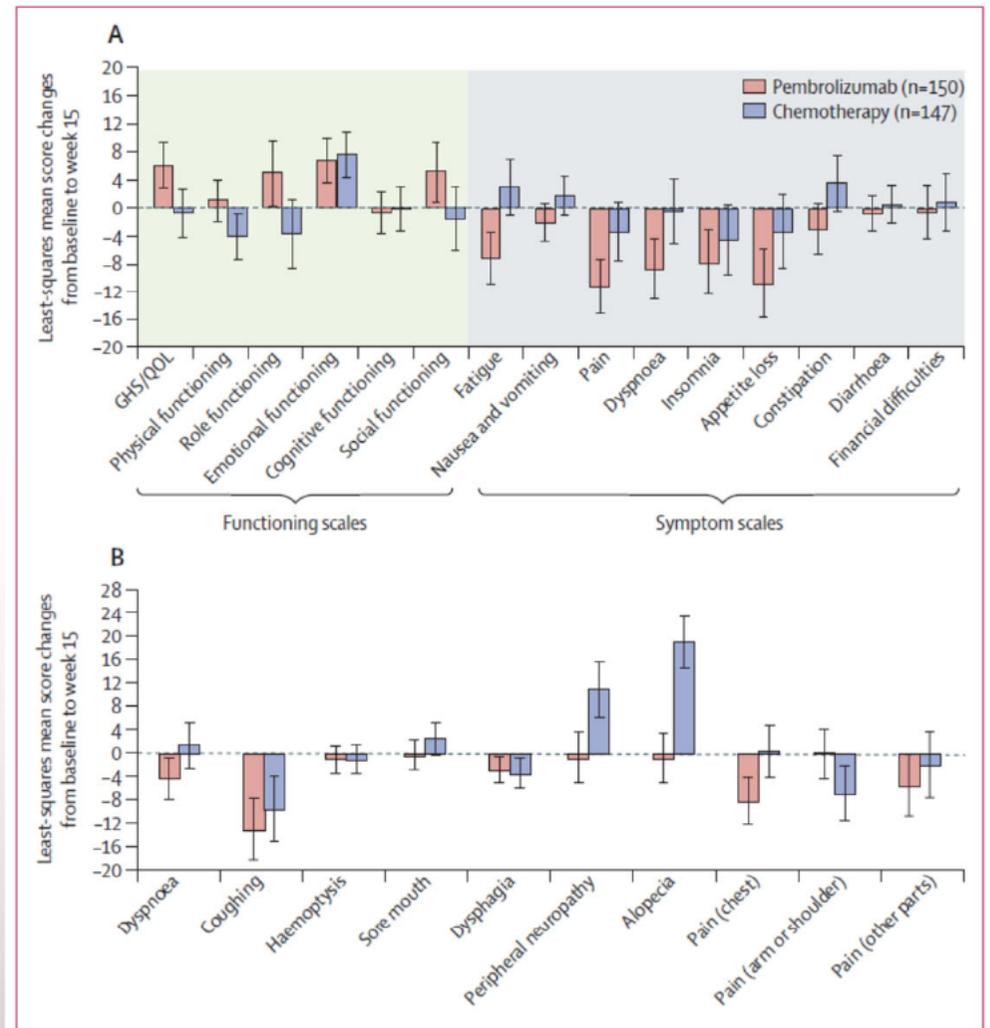


Figure 2: Change from baseline to week 15 in (A) QLQ-C30 functioning and symptom scales and (B) QLQ-LC13 symptoms

For GH5 and functioning scales, higher scores denote improved functioning; for symptom scales, higher scores denote worse symptoms. Mean score changes are based on a constrained longitudinal data analysis model. Error bars represent SEs. QLQ-C30=Quality of Life Questionnaire Core 30 items. QLQ-LC13=Quality of Life Questionnaire Lung Cancer 13 items. Brahmer JR et al. *Lancet Oncol* 2017; 18: 1600-9.

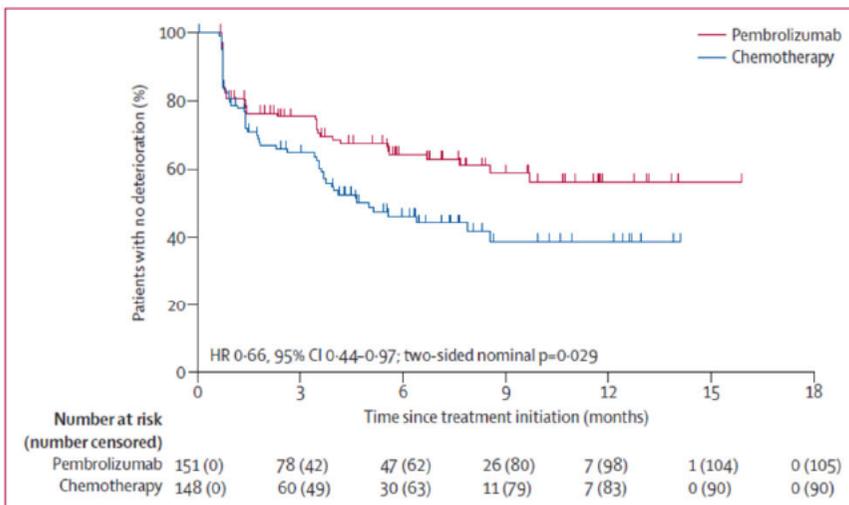


Figure 3: Time to deterioration of the composite of cough, chest pain, and dyspnoea in the QLQ-LC13
 Deterioration defined as time to the first onset of a 10-point or greater decrease from baseline in cough, chest pain or dyspnoea, confirmed by a second adjacent 10-point or greater decrease from baseline in any of these three symptoms. QLQ-LC13=Quality of Life Questionnaire Lung Cancer 13 items. PRO=patient-reported outcome. HR=hazard ratio.

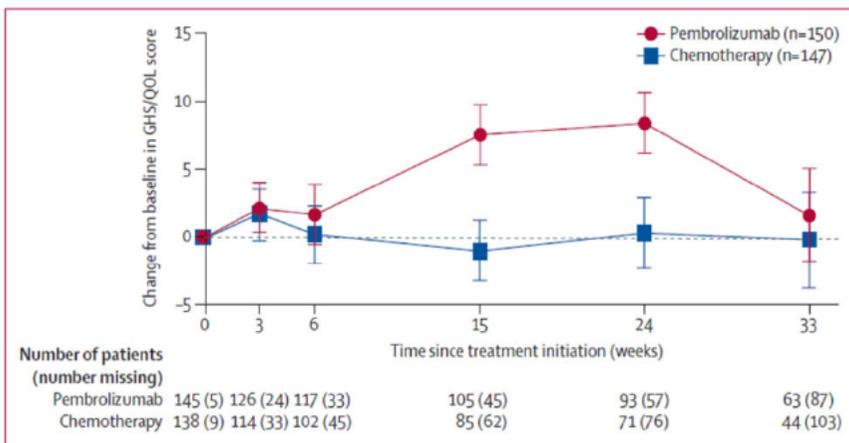


Figure 4: Change from baseline in GHS/QOL score at each study timepoint, without imputation for missing data
 Datapoints represent means and error bars denote SEs. For GHS, positive change from baseline denotes improved functioning. GHS=global health status. QOL=quality-of-life.

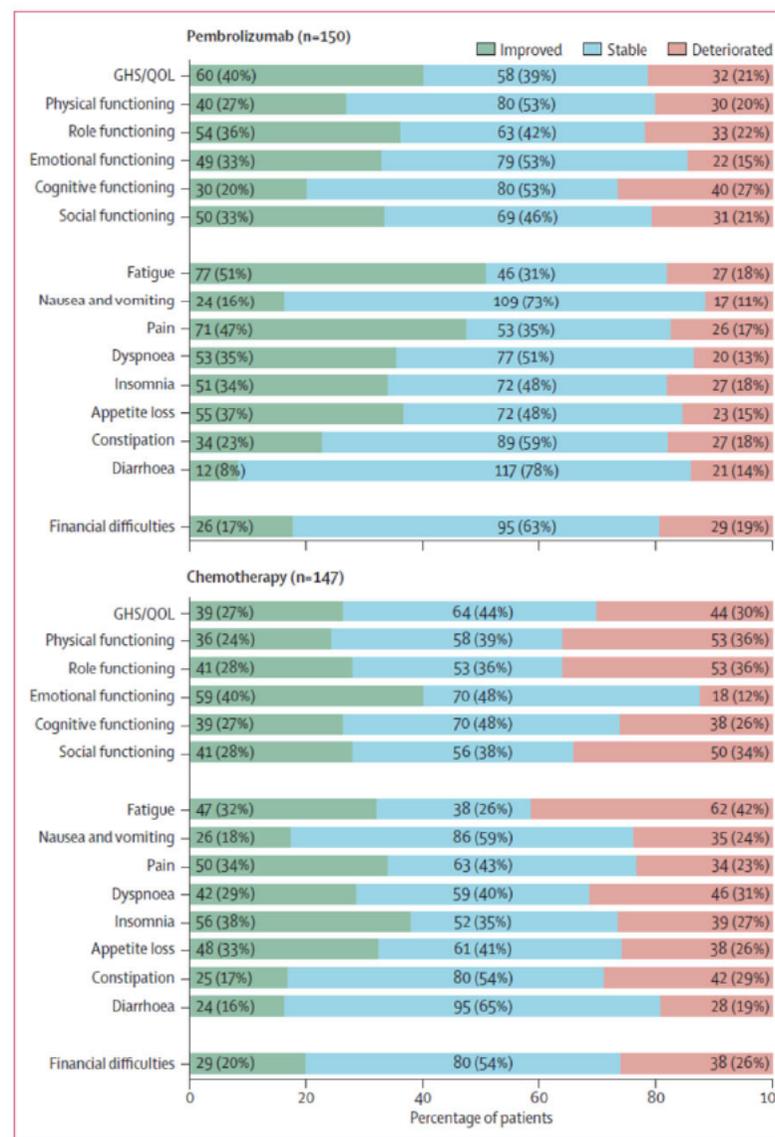


Figure 5: Proportion of patients with improved, stable, and deteriorated QLQ-C30 scores
 GHS=global health status. QOL=quality of life. QLQ-C30=European Organisation for the Research and Treatment of Cancer Quality of Life Questionnaire Core 30 items.

Brahmer JR et al. Lancet Oncol 2017; 18: 1600-9.

データの測定と欠測

- 解析に用いる意味があるが収集されなかったデータ(単に収集されなかったデータのことではない)
- PRO調査にデータの欠測はつきものであり、データの欠測が生じないような調査計画を第一に考えるべき
- データの欠測に結論が影響を受けにくい解析方法、あるいは、データの欠測理由を十分考慮した解析方法の適用が望まれる
- そのためには、データの欠測理由がわかるような調査計画を考慮すべき
- 特定の時点におけるPROデータの欠測
 - 尺度中の全項目ではなく一部の項目が欠測
 - いくつかの尺度のスコアリングマニュアルなどには、欠測の項目があった場合の対処方法(例えば、総合スコアの計算方法)が示されているが、適用が適切かどうか十分に確認する必要
 - PRO評価全体が行われていない
 - 解析においては、データが欠測した理由についての仮定(データの欠測メカニズムの仮定)が必要
 - COMPLETE CASE ANALYSIS、いくつかの補完法、モデルにもとづく方法など、さまざまな統計学的アプローチがあり、時点ごとの欠測状況を集計し、欠測理由に応じて適切な方法を用いなければならない
 - プロトコルにはこれら欠測の対処法について十分記載する必要
 - ただし、欠測値に対処する方法で、普遍的に適用可能と薦められる方法はない
欠測値に対処する方法により解析結果がどの程度変わり易いかを、欠測の数が多い場合には特に検討すべきである(感度解析)

6.5.2. データの欠測 (**missing data**) の取り扱い

PRO/QOL 評価には、データの欠測がつきものである。データの欠測が生じないような調査計画を第一に考えるべきである。また、データの欠測に結論が影響を受けにくい解析方法、あるいは、データの欠測理由を十分考慮した解析方法の適用が望まれる。そのためには、データの欠測理由がわかるような調査計画を立案すべきである。

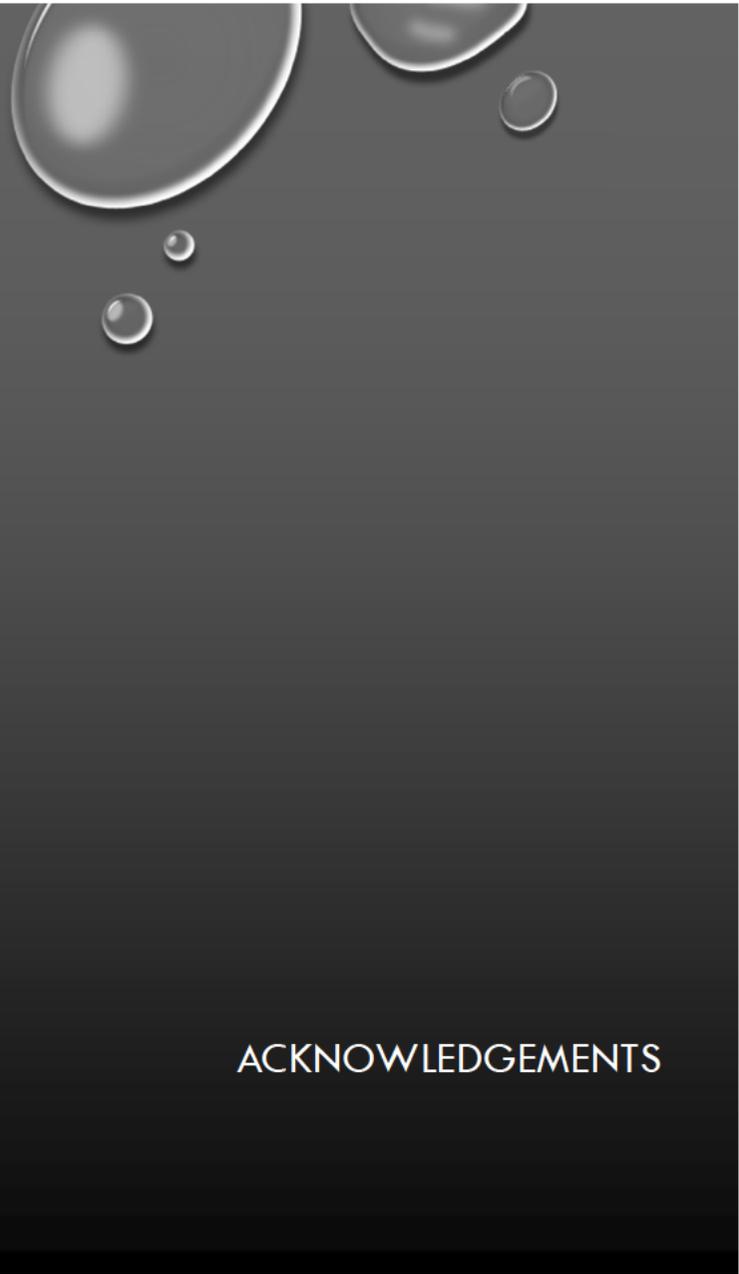
特定の時点における PRO/QOL データの欠測には、2つのレベルがある。(1) 尺度中の全項目ではなく一部の項目が欠測、(2) PRO/QOL 評価全体が行われていない、である。

(1) については、いくつかの尺度のスコアリングマニュアルなどには、欠測の項目があった場合の対処方法（例えば、総合スコアの計算方法）が示されているが、適用が適切かどうか十分に確認する必要がある。

(2) については、解析においてデータが欠測した理由についての仮定（データの欠測メカニズムの仮定）が必要となる。Complete case analysis、いくつかの補完法、モデルにもとづく方法など、さまざまな統計学的アプローチがあり、時点ごとの欠測状況を集計し、欠測理由に応じて適切な方法を用いなければならない。プロトコールにはこれら欠測の対処法について十分記載する必要がある。ただし、欠測値に対処する方法で、普遍的に適用可能と薦められる方法はない。欠測値に対処する方法により結果が異なる場合があるため、解析結果がどの程度変わり易いかを、欠測の数が多い場合には特に検討すべきである。

まとめ

- 原点回帰
- なぜPROデータを収集するのか、真面目に考えてみる



ACKNOWLEDGEMENTS

- TAKASHI KAWAGUCHI (TOKYO UNIV. OF PHARMACY AND LIFE SCIENCES)
- TEMPEI MIYAJI (UNIV. OF TOKYO)
- YASUO OHASHI