



OPEN

Using genomic scars to select immunotherapy beneficiaries in advanced non-small cell lung cancer

H. C. Donker¹, B. van Es^{2,3✉}, M. Tamminga^{1,4}, G. A. Lunter⁵, L. C. L. T. van Kempen⁶, E. Schuurin⁷, T. J. N. Hiltermann¹ & H. J. M. Groen¹

In advanced non-small cell lung cancer (NSCLC), response to immunotherapy is difficult to predict from pre-treatment information. Given the toxicity of immunotherapy and its financial burden on the healthcare system, we set out to identify patients for whom treatment is effective. To this end, we used mutational signatures from DNA mutations in pre-treatment tissue. Single base substitutions, doublet base substitutions, indels, and copy number alteration signatures were analysed in $m = 101$ patients (the discovery set). We found that tobacco smoking signature (SBS4) and thiopurine chemotherapy exposure-associated signature (SBS87) were linked to durable benefit. Combining both signatures in a machine learning model separated patients with a progression-free survival hazard ratio of $0.40^{+0.28}_{-0.17}$ on the cross-validated discovery set and $0.24^{+0.31}_{-0.14}$ on an independent external validation set ($m = 56$). This paper demonstrates that the fingerprints of mutagenesis, codified through mutational signatures, select advanced NSCLC patients who may benefit from immunotherapy, thus potentially reducing unnecessary patient burden.

In non-small cell lung cancer (NSCLC), response to immunotherapy is low. Radiology-assessed response is typically around 20–25%¹, while the percentage of patients achieving durable benefit (DB), defined as progression-free survival (PFS) $\geq \frac{1}{2}$ year, is only slightly higher. Predictors can help to narrow down specific subpopulations for which treatment is particularly effective. Programmed death ligand 1 (PD-L1) protein expression in tumor tissue, and tumor mutational burden (TMB), defined as the number of acquired amino acid sequence-changing mutations², is currently used to predict the efficacy of immunotherapy in NSCLC³. The explanation behind the predictive value of TMB is that the accumulation of mutations in coding DNA results in a high diversity of neoantigens. In turn, these neoantigens induce a broad anti-tumor adaptive immune response. *In silico* analysis of neoantigens found a near-perfect relation between TMB and inferred number of neoantigens⁴. Once the cancer cell expresses neoantigens, the cancer cell can be eliminated through immune recognition and cell killing^{1,2}.

While several studies show a clear association of TMB with response to treatment, not all do⁵, motivating the search for an improved proxy for treatment efficacy. The chief advantage of TMB is that it is easy to compute by pooling all amino acid sequence-changing mutations from whole-exome sequencing data², irrespective of their genomic context. However, this ignores the fact that somatic mutations are generated by a range of processes. Nucleotide context is a clue to the genesis of mutations⁶. Analysis of mutation spectra, which partition mutations by alteration and DNA context, have revealed that both exogenous and endogenous DNA mutational processes can be linked to specific signatures^{7–9}. Over the last decade, increasingly large datasets, such as the COSMIC database, have enabled the systematic identification of mutational signatures^{10,11}. And in many cases, also the elucidation of their aetiology, particularly for single base substitutions⁷. For instance, specific DNA mutational

¹Department of Pulmonary Diseases, University of Groningen, University Medical Centre Groningen, Hanzeplein 1, P.O. Box 30.001, 9700 RB Groningen, The Netherlands. ²Central Diagnostic Laboratory, University Medical Centre Utrecht, Utrecht University, Heidelberglaan 100, 3508 GA Utrecht, The Netherlands. ³MedxAI, Theophile de Bockstraat 77-1, 1058 VA Amsterdam, The Netherlands. ⁴Department of Internal Medicine, Twente Hospital, Enschede, The Netherlands. ⁵Department of Epidemiology, University of Groningen, University Medical Centre Groningen, 9713 GZ Groningen, The Netherlands. ⁶Department Of Pathology, University of Antwerp, University Hospital Antwerp, 2650 Edegem, Belgium. ⁷Department of Pathology and Medical Biology, University of Groningen, University Medical Centre Groningen, Hanzeplein 1, P.O. Box 30.001, 9700 RB Groningen, The Netherlands. ✉email: bes3@umcutrecht.nl

signatures or “genomic scars” left by e.g. polycyclic aromatic hydrocarbons¹², ultraviolet light exposure¹³, and platinum chemotherapy¹⁴ have all been validated in experiments.

Evidence is accumulating that mutational signatures in cancer are therapeutically relevant⁹. Given the partial success of TMB for predicting immunotherapy efficacy in NSCLC², a logical next step is to try and improve TMB by sieving out irrelevant mutations. Mutational signatures, which disentangle mutations by their proposed root cause, may help to pinpoint relevant alterations. In fact, specific mutational signatures (e.g., APOBEC A3A) have been hypothesised to be promising candidates for immune stimulation treatments¹⁵. Along this line, previous studies performed de novo identification of single base substitution signatures in NSCLC and linked these to immunotherapy efficacy^{16,17}. Given the relatively small datasets used in these studies, only three signatures were identified in both cases^{16,17}. Specifically, Wang et al.¹⁶ found that patients with durable clinical benefit were enriched in signatures with similarities to COSMIC signatures SBS2 and SBS13 that are associated with damage from APOBEC, a family of enzymes that are part of the innate anti-retroviral defense that operates by generating mutations in single-stranded DNA^{9,18}. Signatures similar to clock-like singlet SBS1, capturing mutations that steadily accrue with age, were linked to non-response to immunotherapy by Chong et al.¹⁷.

Here, instead of de novo analysis we directly use previously catalogued mutational signatures like in earlier work^{4,19}. We expand previous efforts by interrogating an order of magnitude more signatures, including the recently developed copy number signatures^{20,21}, for their role in eliciting an immune response. In addition, we validate the presence of single base substitution signatures at the RNA level. The primary goal of this work is to develop and validate an immunotherapy efficacy model for advanced NSCLC patients to help reduce ineffective treatment. We hypothesise that the genomic scars in pre-treatment tumor tissue, decomposed into mutational signatures, can help identify durable immunotherapy beneficiaries (Fig. 1a).

Result and discussion

In total, $M = 157$ patients with advanced NSCLC were analysed, $m = 101$ in the discovery and $m = 56$ in a separate validation set. The population had a mean age of 63.2, consisted of an approximately equal split between males and females (53.5%). The majority of the patients (85/147, 57.8%) did not achieve durable benefit from immunotherapy (Table 1). In the Discovery dataset, 21 patients (20.8 %) received combination therapy.

Mutational signatures in the discovery set. Whole genome sequencing revealed a mean of 7.80 Mb^{-1} [median 5.43 Mb^{-1} ; inter quantile range (IQR) $3.90\text{--}9.85 \text{ Mb}^{-1}$] non-synonymous variants in the discovery set, consistent with previous reports on whole exome sequencing^{22,23}.

Mutational signature data of $p = 47$ single base substitutions (SBS), $p = 11$ doublet base substitution (DBS), $p = 17$ short insertion and deletion (indel), and $p = 20$ copy number signatures were determined (see Fig. 1b for a schematic overview) and analysed, after discarding sequencing artefact attributed signatures (Supplementary Table 1). Ranking mutational signature attribution W (see the Methods for details) by the median (across samples) shows that nine SBS signatures are present in the majority of samples (Fig. S1a, Supplementary Material), with signature SBS4, a lung cancer-specific signature²⁴ that is linked to tobacco smoking^{6,7,12}, having the highest median value (0.91 Mb^{-1}), although exhibiting considerable variability (IQR $0.00\text{--}3.14 \text{ Mb}^{-1}$). Signature attributions of indel and doublet base substitution were an order of magnitude lower (Fig. S1b, Supplementary Material) with the two highest median signatures, ID3^{7,10} (median: 0.11 Mb^{-1} , mean: 0.17 Mb^{-1}) and DBS2^{7,25} (median: 0.06 Mb^{-1} , mean: 0.10 Mb^{-1}), both attributed to tobacco smoking, like SBS4. Whole genome copy number deconvolution revealed that homologous recombination deficiency²¹ signature CN17 had the highest median attribution (median: 37.2, mean: 46.5, Fig. S1c, Supplementary Material).

Mutational signatures linked to durable benefit. Next, we looked for pre-immunotherapy mutational signatures that were determinants of therapy efficacy. Univariate analysis ($m = 93$ patients) singled out two single base substitution signatures (Fig. 1c). Tobacco smoking signature SBS4 was significantly different ($q = 0.014$, B-HK-S test) in patients who derive DB from immunotherapy. This finding underpins earlier work that found that smoking-attributed transversion-high tumors²⁶, or enrichment in smoking signature¹⁹, had improved outcome in ICI-treated NSCLC. Similarly, SBS87—whose mutations coincide with thiopurine chemotherapy exposure^{9,27}—was also found to differ between both groups of patients ($q = 0.017$, B-HK-S test). Note that, according to the clinical records available to us, none of the patients has been treated with thiopurine-related compounds.

We subsequently considered the mutual dependence between SBS4 and SBS87. Part of this correlation is mediated through the outcome (DB versus non-DB) and the total number of amino-acid sequence-changing mutations (i.e., TMB). We, therefore, normalised both signatures by TMB and correlated separately for patients with DB and non-DB. In both outcome groups, signatures SBS4 and SBS87 were unrelated (Kendall $\tau = -3.6 \cdot 10^{-2}$, $p = 0.76$ and $\tau = -8.2 \cdot 10^{-3}$, $p = 0.95$ for DB and non-DB, respectively).

Earlier work linked clock-like mutational signature SBS1—capturing substitutions that steadily accrue with age—with non-response and worse survival¹⁷; we could not replicate the association with DB, even without multiple testing correction ($p = 0.36$, K-S test). Another study linked mutational signatures associated with APOBEC—a family of enzymes that are part of the innate anti-retroviral defense that operates by generating mutations in single-stranded DNA^{9,18}—with improved immunotherapy outcome¹⁶. Validation of APOBEC mutational signatures SBS2 and SBS13¹⁶ showed a trend for SBS13 ($p = 0.04$ K-S test, $q = 0.76$ B-HK-S test) but not for SBS2 ($p = 0.66$ K-S test, $q = 1.0$ B-HK-S test). Given that APOBEC mutagenesis is highly transient, with episodic bursts of mutations²⁸, our samples, which represent a snapshot in time, are perhaps less suited to fully interrogate the relevance of this mutational signature on treatment outcome. None of the doublet, indel, and genome-wide copy number alteration signatures was significantly associated with DB.

| | Discovery (Hartwig DR #094) | Validation [Miao et al. ⁴] | Overall |
|---------------------------------------------------|-----------------------------|----------------------------------------|------------------|
| <i>m</i> | 101 | 56 | 157 |
| Age, μ (σ) | 63.84 (8.72) | 61.49 (8.75)* | 63.16 (8.77) |
| Gender, <i>m</i> (%) | | | |
| Female | 52 (51.49) | 32 (57.14) | 84 (53.50) |
| Male | 49 (48.51) | 24 (42.86) | 73 (46.50) |
| Smoker, <i>m</i> (%) | | | |
| Current | | 14 (25.00) | 14 (8.92) |
| Former | | 29 (51.79) | 29 (18.47) |
| Never | | 13 (23.21) | 13 (8.28) |
| Unknown | 101 (100.00) | | 101 (64.33) |
| Prior therapy, <i>m</i> (%) | | | |
| Chemotherapy | 33 (32.67) | | 33 (21.02) |
| Naive | 18 (17.82) | | 18 (11.46) |
| Radiotherapy | 7 (6.93) | | 7 (4.46) |
| Radiotherapy + chemotherapy | 43 (42.57) | | 43 (27.39) |
| Unknown | | 56 (100.00) | 56 (35.67) |
| Treatment, <i>m</i> (%) | | | |
| α CTLA-4 + α PD-1 | 1 (0.99) | | 1 (0.64) |
| α CTLA-4 + α PD-1 with chemotherapy | 1 (0.99) | | 1 (0.64) |
| α PD-1 | 77 (76.24) | | 77 (49.04) |
| α PD-1 with chemotherapy | 13 (12.87) | | 13 (8.28) |
| α PD-1/ α PD-L1 | | 56 (100.00) | 56 (35.67) |
| α PD-L1 | 3 (2.97) | | 3 (1.91) |
| α PD-L1 with chemotherapy | 2 (1.98) | | 2 (1.27) |
| α PD-L1/cabozantinib | 1 (0.99) | | 1 (0.64) |
| α TNFRSF7 + α PD-1 | 1 (0.99) | | 1 (0.64) |
| α VEGF-A + α PD-1 | 2 (1.98) | | 2 (1.27) |
| Durable benefit, <i>m</i> (%) | | | |
| No | 57 (56.44) | 28 (50.00) | 85 (54.14) |
| Yes | 36 (35.64) | 26 (46.43) | 62 (39.49) |
| Unknown | 8 (7.92) | 2 (3.57) | 10 (6.37) |
| Sequencing, <i>m</i> (%) | | | |
| Whole exome | | 56 (100.00) | 56 (35.67) |
| Whole genome | 101 (100.00) | | 101 (64.33) |
| Tissue, <i>m</i> (%) | | | |
| Formalin-fixed paraffin-embedded | | 56 (100.00) | 56 (35.67) |
| Fresh frozen | 101 (100.00) | | 101 (64.33) |
| Tumor purity, median [Q_1, Q_3] | 0.41 [0.28,0.56] | 0.35 [0.24,0.54] | 0.39 [0.27,0.56] |

Table 1. Patient characteristics and outcome of advanced non-small cell lung cancer. Symbols and abbreviations: α CTLA-4, cytotoxic T lymphocyte-associated antigen-4 inhibitor; α PD-1, programmed death-1 inhibitor; α PD-L1, programmed death-ligand 1 inhibitor; α TNFRSF7, tumor necrosis factor receptor superfamily type 7 inhibitor; α VEGF-A, vascular endothelial growth factor inhibitor; μ , mean; σ , standard deviation; DR, data request; Q_1 , first quartile; Q_3 , third quartile. Age was missing for 15 patients in the validation cohort.

A signature-based classifier predicts immunotherapy benefit. In combining SBS4 and SBS87 signature attributions (representing, per signature, the number of amino-acid sequence changing singlets in DNA) both remained significant in each cross-validated fold, confirming that the aforementioned univariate analysis does not lead to overfitting on the discovery set. The classifier scored an area under the receiver operating characteristic curve (ROC AUC) of $0.74^{+0.11}_{-0.12}$ on the hold-out folds (Fig. S3, Supplementary Material). This was significantly higher than when the classifier was trained on TMB (ROC AUC: $0.65^{+0.12}_{-0.13}$, $p = 0.016$ PPT, Fig. S3a, Supplementary Material). The ROC curves of individual signatures were similar to that of the model (Fig. S4). With an estimated 43.8% patients with DB²⁹ as our classification probability threshold, a sensitivity of $0.56^{+0.16}_{-0.16}$, a specificity of $0.86^{+0.08}_{-0.10}$, and an accuracy of $0.74^{+0.09}_{-0.10}$ was achieved (Table 2) with 8 false positive and 16 false negative classifications (Supplementary Table 3). The classifier was not perfectly calibrated (Fig. S5a), as expected from the conditional independence assumption of the naive Bayes model.

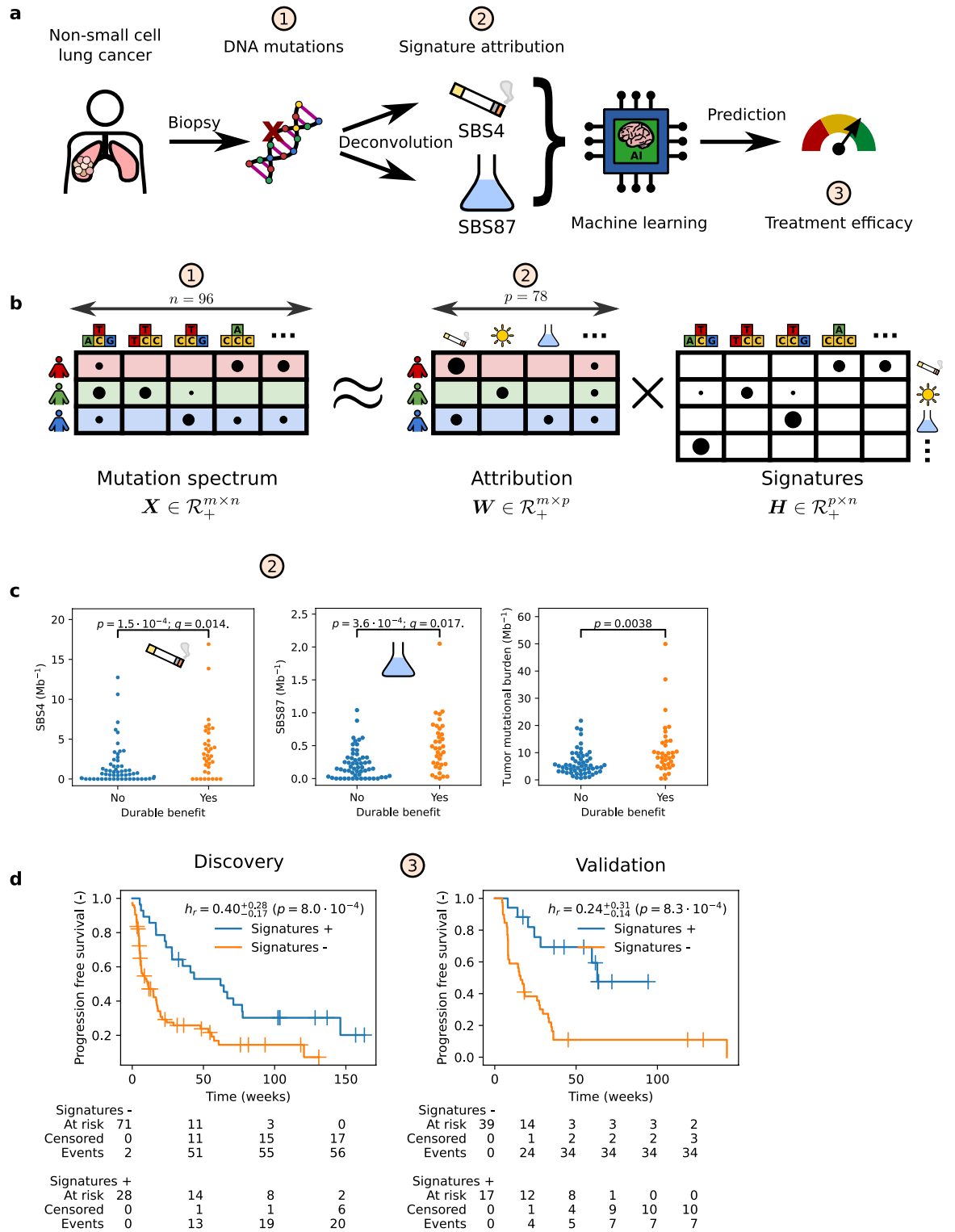


Figure 1. Mutational signatures from pre-treatment tumor tissue predict immunotherapy efficacy in advanced non-small cell lung cancer. (a), Single base substitutions (SBS) are determined from pre-immunotherapy tumor material. After deconvolution into signature attributions, a machine learning classifier uses smoking-associated signature SBS4 and thiopurine chemotherapy-associated signature SBS87 to predict durable benefit (DB) from immunotherapy. (b), Cartoon illustration of SBS signature deconvolution, where we solve for signature attribution W given mutation spectrum X and COSMIC signatures H through $X \approx WH$. Nucleotide pyramids indicate SBS with flanking context; Sun, cigarette, and Erlenmeyer symbols depict example aetiologies; Shading highlights information that pertains to the corresponding patient. For illustration purposes, the size of the dots do not represent actual data. (c) Signatures SBS4 [$q = 0.014$, Benjamini–Hochberg corrected Kolmogorov–Smirnov (B–HK–S) test] and SBS87 ($q = 0.017$, B–HK–S test) are linked to DB (discovery set). For reference, tumor mutational burden is also shown. (d) Patients predicted to have DB (Signatures +, blue line) have superior progression-free survival compared to those predicted to have non-DB (Signatures -, orange line) in the discovery set (left panel). The classifier’s performance replicates in an independent validation set (right). Censored observations are indicated by crosses. Estimates and corresponding 95% confidence intervals are indicated by sub and superscripts.

| | ROC AUC | AP | F_1 | Sensitivity | Specificity | Accuracy | h_r * |
|-----------------------------------------|----------------------------------------|----------------------------------------|----------------------------------------|----------------------------------------|----------------------------------------|----------------------------------------|----------------------------------------|
| Discovery (Hartwig DR #094) | 0.74 ^{+0.11} _{-0.10} | 0.63 ^{+0.18} _{-0.15} | 0.63 ^{+0.14} _{-0.16} | 0.56 ^{+0.16} _{-0.17} | 0.86 ^{+0.09} _{-0.10} | 0.74 ^{+0.09} _{-0.10} | 0.40 ^{+0.28} _{-0.17} |
| Validation (Miao et al. ¹⁸) | 0.69 ^{+0.13} _{-0.15} | 0.71 ^{+0.16} _{-0.19} | 0.57 ^{+0.17} _{-0.19} | 0.46 ^{+0.21} _{-0.21} | 0.86 ^{+0.11} _{-0.14} | 0.67 ^{+0.13} _{-0.13} | 0.24 ^{+0.31} _{-0.14} |
| p | 0.57 | 0.56 | 0.68 | 0.54 | 1.00 | 0.35 | 0.32 |

Table 2. Durable benefit (DB) classifier performance metrics and p -value comparing the performance between the two datasets. ROC AUC area under the receiver operating characteristic curve, AP average precision, h_r hazard ratio of progression-free survival (PFS) comparing patients predicted DB versus predicted non-DB. Estimates and corresponding 95% confidence intervals are indicated by sub and superscripts. *All patients were included in the PFS analysis. For all other classification metrics, only patients where censoring permitted unambiguous DB label assignment were analysed.

To incorporate patients censored prior to the $\frac{1}{2}$ year mark, we compared predicted versus actual outcome using Kaplan–Meier (Fig. 1d) in all $m = 101$ discovery patients. Durable benefit predicted patients had superior median progression-free survival (62 for predicted DB versus 11 weeks for predicted non-DB). Cox regression provided additional confirmation that the classifier significantly predicted outcome based on mutational signatures SBS4 and SBS87 combined (hazard ratio $h_r = 0.40$ ^{+0.28}_{-0.17}, $p = 8.0 \cdot 10^{-4}$).

Independent validation on the data from Miao et al.⁴ reproduces the classifier’s performance. The median survival in the external validation set ($m = 56$ patients) was 63 versus 16 weeks for DB-predicted patients compared to the others (Fig. 1d). All performance metrics were similar (Table 2), highlighting the reproducibility of our approach despite differences in (i) tissue handling (formalin-fixed paraffin-embedded versus fresh frozen), (ii) chemistry (whole exome versus genome capture), and (iii) bioinformatics pipeline. Unlike the discovery set, the ROC AUC of a model trained on TMB was not significantly different from the mutational signature model (ROC AUC 0.78^{+0.12}_{-0.14} versus 0.69^{+0.14}_{-0.14}, respectively, $p = 0.18$ PPT, Fig. S3b, Supplementary Material). This discrepancy was attributed to a difference in TMB distribution between the discovery and validation set ($p = 0.031$, KS test), while both SBS4 and SBS87 signature attributions remained similar in both cohorts ($p = 0.082$ and $p = 0.69$, respectively, KS test). A fair head-to-head comparison between TMB and the mutational signature approach requires a separate, substantially larger, study with more detailed patient characteristics and medical history than presented here (together with a more harmonised tissue handling, sequencing chemistry, and upstream bioinformatics preprocessing).

Error analysis of the discovery set revealed that treatment-naïve patients were overrepresented in the top ten worst predicted false negatives ($p = 0.015$, Fisher exact test). Exclusion of ($m = 18$) treatment-naïve patients slightly improved the model (ROC AUC: 0.79^{+0.11}_{-0.12}) on the discovery set. More strikingly, after training on pre-treated patients only, generalization on the validation patients—of whom prior therapy was unknown—also improved (ROC AUC: 0.71^{+0.14}_{-0.15}). Caution is therefore warranted when applying the classifier to an (in our study, underrepresented) treatment-naïve population.

Subanalysis of patients with smoking status ($m = 54$, validation set only, Table 1) showed that never smokers were more difficult to classify (ROC AUC: 0.40^{+0.10}_{-0.15} versus 0.71^{+0.12}_{-0.13}, $p = 0.048$ UPT), although the numbers were low with only $m = 3$ durable beneficiaries in the never-smokers group.

Note that current/former smokers are known to have a better overall response rate³⁰. Seeing that the performance was lower on the combined set indicates that the mutational signatures provide information that is orthogonal (or, complementary) to smoking status.

Decision curve analysis. Net benefit is a decision-theoretic concept that quantifies the practical utility of a classifier^{31,32}. The integrated net benefit of using the signature-based model is positive (Figs. S6 and S7, Supplementary Material), with a median integrated combined net benefit (see e.g. Talluri et al.³³) of 0.37 (IQR 0.366–0.379) and 0.29 (IQR 0.284–0.296) with respect to net benefits less than zero when treating *all* patients for the discovery and validation datasets, respectively. In general, according to Ref.³² a model can be recommended for clinical use if, across a range of clinically reasonable probability thresholds, it has the highest level of benefit. This range of thresholds can be viewed as the acceptable range of *number-needed-to-treat* (NNT) to have one effective treatment and is clinician/protocol dependent. In Figs. S7 and S8, we see a net benefit for both the treated/untreated with respect to treating all patients for a probability threshold range of roughly 0.3–0.6 and 0.4–0.6 for the discovery and validation datasets respectively. Important to note, although we have fairly poor calibration for the 0.3–0.6 range within which we expect a relative net benefit with respect to the baseline, outside this range the net benefit is roughly equal. That is, within a probability threshold range of 0.3–0.6, so an NNT-range of roughly 1.5–3 patients, we have a net benefit when applying our model. Outside this range, there is no added benefit.

Relation of signature specific mutations to specific genomic loci. To better understand why SBS4 and SBS87 relate to DB, we set out to link their attribution to specific genes (in the Discovery set, $m = 101$ patients). Neither the SBS4 nor the SBS87 signature was correlated to the number of mutations in any of the 23 significantly lung cancer mutated genes in TCGA^{22,23} (Supplementary Table 2). Expanding the search from 23 to the top 2.5%, 5% and 10% highly expressed genes also found no correlation. In contrast, analysis of mutations in 523 genes contained in the clinically relevant TSO500 panel, consisting of genes canonically mutated in cancer, yielded eight genes [*ATM* ($q = 0.029$), *EPHA5* ($q = 0.013$), *LRP1* ($q = 5.3 \times 10^{-4}$), *MTOR* ($q = 0.034$), *NRG1* ($q = 0.036$), *PTPRD* ($q = 1.7 \times 10^{-5}$), *PTPRT* ($q = 6.7 \times 10^{-5}$), *RUNX1T1* ($q = 0.023$)], Kendall τ correlation

test] in which mutation count correlated significantly with smoking signature SBS4 (Fig. 2) after multiple testing correction and exclusion of non-mutated genes. When combined, the aggregated mutation count was also directly linked to DB ($p = 0.0050$, K-S test) in addition to the indirect correlation through SBS4. Four of these genes [namely, *LRP1* ($q = 0.021$), *PTPRD* ($q = 0.021$), *PTPRT* ($q = 0.0091$), and *RUNX1T1* ($q = 0.039$), but no other genes] also correlated with SBS87. Overall, significant genes were large, ranging from 146 Kb (*ATM*) to 2.3 Mb (*PTPRD*)^{34,35}. This was expected since in order to (significantly) correlate, enough mutations must be detected (the correlation with only zeroes is trivially zero). And the larger the gene, the more mutations can accumulate randomly. Functionally, both *ATM* and *EPHA5* interact at the site of DNA repair. Adding *ATM* to a DB logistic regression model with TMB changed the regression coefficient by more than 10 % (0.156 versus 0.117) indicating that *ATM* (but not *EPHA5*) potential confounds TMB. *MTOR* (a paralog of *ATM*) regulates cellular metabolism and the others are tumor suppressor genes. They are involved in cell interactions such as the *PTPRT* and suppression of inflammatory responses such as *STAT3*. *PTPRD* show deleterious mutations in 9% of lung cancers. *NRG1* suppresses the transcription of inflammatory cytokines and was the only gene (out of all eight) that was significant ($p = 0.020$) when added to a logistic regression model with TMB. Compared to TMB, mutations in *NRG1* anti-correlated with DB (as indicated by the negative regression coefficient -2.6). One explanation could be that mutations in *NRG1* disrupt the adaptive immune system's capacity to induce an immune response. SBS87 only correlated with tumor suppressor genes. However, statistically, no enrichment for tumor suppressor or oncogenes was found in the set correlating with SBS4 ($p = 0.25$ and $p = 0.16$, respectively, Fisher exact test) nor with SBS87 ($p = 0.08$ and $p = 0.16$, respectively, Fisher exact test) relative to the TSO500 gene set.

Expression of SBS4/SBS87 mutations is not a sufficient condition for predicting durable benefit. Next, we aimed to explain the predictiveness of SBS4 and SBS87 attributed mutations by looking for differences at the RNA level. Assume that immune recognition, through the presentation of mutated protein fragments on the major histocompatibility complex, is a prerequisite for eliciting an immune response. If transcription is a necessary condition for immune recognition of mutated DNA then, by extension, it must also be a necessary condition for predicting DB, assuming immunotherapy operates through an adaptive immune response. Focussing on transcripts with a SBS4 and/or SBS87 dominant variant (a signature-dominant variant accounts for $\geq 50\%$ of the signature's mutations, see Sec. A.4, Supplementary Material), we asked whether their transcription is also a sufficient condition for predicting DB.

We, therefore, did a subanalysis of paired RNA and DNA samples ($m = 36$ patients, Discovery dataset). To ensure that the subanalysis does not introduce a bias in the statistical analysis, we first verified that durable benefit ($p = 0.67$, Fisher exact test) and both SBS4 and SBS87 signatures ($p = 0.30$, $p = 0.28$ respectively, K-S test) did not differ from the entire discovery population. Next, comparing the number of transcripts containing an SBS4 or SBS87-dominant singlet versus transcripts containing any other singlet shows that the RNA abundance is similar in both groups ($p = 0.71$, K-S test, Fig. S2a, Supplementary Material). Compared to DB, the distribution of transcripts harbouring a smoking-associated signature SBS4-dominant variant was different from patients with non-DB ($p = 0.0036$, K-S, see Fig. S2b, Supplementary Material). However, this difference could be attributed to the number of SBS4 DNA mutations (Fig. S2c, Supplementary Material). For the thiopurine chemotherapy associated SBS87 signature, no difference in absolute number ($p = 0.11$, K-S, Fig. S2d, Supplementary Material) nor relative to the number of corresponding variants ($p = 0.48$, K-S, Fig. S2e, Supplementary Material) was found in the number of mutated transcripts between patients with and without DB. Differential gene expression of the mutated RNA of the aforementioned eight significant genes ($m = 22$ patients with paired DNA, RNA, and ≥ 1 mutations in any of these eight genes) detected no difference in the amount of mutated RNA between patients with and without DB ($q > 0.05$ for all genes, B-HK-S test). Together, these results show no evidence that distinguishes SBS4 and SBS87 from other mutations at the transcription level in durable beneficiaries. While these negative findings could point to the differences (i) being manifested further downstream immune

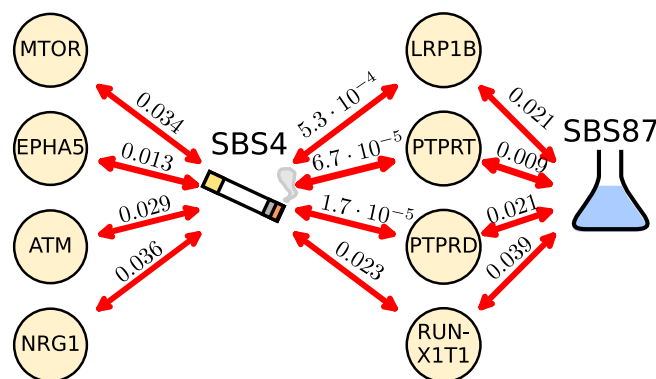


Figure 2. Signatures SBS4 and SBS87 correlate with mutations in genes canonically mutated in cancer (discovery set). Correlations were assessed using a Benjamini-Hochberg corrected Kendall τ (corrected p-values along the arrows) and the correlation strengths are indicated by the arrow line width.

processes or (ii) the capacity to bind to a broader spectrum of HLA alleles³⁶, we do not rule out that we lacked insufficient variant coverage at the RNA level (which was on average 19.8× for all and 19.1× for SBS4/SBS87 dominant variants) to detect subtle differences.

Conclusion

We have identified tobacco smoking signature SBS4 and the recently identified thiopurine chemotherapy exposure-associated signature SBS87^{9,27} as factors that are predictive of benefit from immunotherapy. Both signatures are linked to mutations in genes involved in DNA repair and/or the function of tumor suppressor genes that have a role in cellular immunological interactions and inflammatory cytokines. In contrast, none of the doublet base substitution, indel, or copy number alteration signatures were associated with durable benefit from immunotherapy. RNA analysis of the two signatures found no evidence that distinguishes these from other mutations in terms of immunotherapy efficacy.

When combined, these two signatures can help to select advanced NSCLC patients who may benefit from (combination) immunotherapy using information available prior to treatment initiation. These signatures were not tested in the context of other therapeutic interventions. As such, further investigations are required to validate if this signature is therapy-specific or if it may have prognostic value. The advantage of our predictor is that it inherits the mechanistic grounding of mutational signatures. However, more research is needed to establish if our approach reliably generalises to the (underrepresented) treatment-naive and non-smoking patient populations, or if additional adjustments are needed.

Methods

Cohort assembly. We retrospectively compiled clinical records and data of tumor and matched normal tissue (collected prior to treatment initiation) of immunotherapy-treated advanced NSCLC patients. To this end, patient characteristics, whole genome sequencing (WGS) and total RNA data derived from fresh frozen biopsies were requested from Hartwig Medical Foundation (HMF)³⁷. This cohort, containing metastatic NSCLC patients, formed the discovery dataset. The external validation was extracted from Ref.⁴ and consisted of whole exome sequencing (WES) of formalin-fixed paraffin-embedded (pre-immunotherapy) NSCLC samples of tumors and matched normal tissue. A more detailed description of cohort assembly can be found in the Supplementary Material.

DNA processing. Whole genome sequencing, variant calling, and purity estimation were performed by HMF^{37,38}. Since the whole genome sequencing reads were mapped to GRCh37, we used crossmap³⁹ to perform a liftover (or, remapping) from GRCh37 to GRCh38 using Ensembl's corresponding chain file.

Mutation deconvolution. COSMIC v3.3 (June 2022) mutational signatures, \mathbf{H} , were used to deconvolute mutations (see Supplementary Table 1 for a list of analysed signatures). For substitutions and short indels, these signatures \mathbf{H} describe the nucleotide alteration distribution^{10,40}. Release v3.3 adds the recently developed copy number signatures which capture the copy number × zygosity × length distribution^{20,21}.

We first used SigProfilerMatrixGenerator on amino acid-changing mutations to extract mutation spectra of single base substitutions (singlets) with two flanking bases (SBS-96), doublet base substitutions (doublets, DB-78), and insertion deletions (indels, ID-83)⁴¹. The same package was used to compute genome-wide copy number alterations (CN-48)²¹. Each mutation spectrum, \mathbf{X} , is a positive m -by- n matrix (i.e., $\mathbf{X} \in \mathcal{R}_+^{m \times n}$) consisting of m samples, and n mutation/channels ($n = 96, n = 78, n = 83, n = 48$ for SBS-96, DB-78, ID-83, and CN-48, respectively) counting the number of mutations per channel. The positive mutational signature matrix \mathbf{H} relates the p signatures (rows) to the corresponding n mutation type/channels (columns), $\mathbf{H} \in \mathcal{R}_+^{p \times n}$. (Throughout this paper we adhere to the convention that m, n , and p indicate the number of patients, number of mutation types/channels, and number of signatures [except for p -values which will be clear from the context], respectively.) Using the spectrum \mathbf{X} and mutational signatures \mathbf{H} , signature attribution \mathbf{W} , a positive m -by- p matrix (i.e., $\mathbf{W} \in \mathcal{R}_+^{m \times p}$) such that $\mathbf{X} \approx \mathbf{WH}$, was computed by non-negative matrix factorisation using a coordinate descent solver with an error tolerance of 10^{-6} for no more than 10^4 iterations.

Signature attributions that refer to possible sequencing artefacts⁴⁰ or with zero variance in either durable benefit stratum in the discovery set were excluded from the analysis. Mutational signatures of singlets, doublets, indels, and TMB (i.e., total mutation count) were obtained from non-synonymous mutations and normalised by (exome) coverage size in megabases (Mb) and rounded to two decimals to retain three significant digits. For the Hartwig WGS, the exome coverage size in individual samples was unknown. Therefore, a size of 47.9 Mb was taken⁴². For the discovery cohort (Miao et al.⁴), we used the size as indicated in their Supplementary Table 1. Since the copy number variants span large portions of the genome (both exonic and intronic regions), we report the total, whole genome, copy number attributions.

RNA processing. Singlets were traced back to transcripts to study how the mutational signatures manifest at the transcription level, as a surrogate for protein expression. Out of the $m = 101$ patients in the discovery cohort, raw (total) RNA sequencing data of $m = 40$ patients were available (no RNA was available in the validation cohort). Briefly, raw sequencing data were trimmed, aligned, and converted into transcripts per million (TPM). After quality control, two inferior-quality samples and two samples with insufficient follow-up were excluded, leaving a total of $m = 36$ samples for analysis. The amount of transcripts containing a variant was re-estimated to account for differences in tumor content. RNA per signature was obtained by pooling transcripts containing the $\geq 50\%$ dominant mutations of the given signature. A more detailed description of the method can be found in the Supplementary Information.

Statistical analysis. Here and in the following, all tests were two-sided.

Univariate analysis. Differences in mutational signature distributions were determined by a Kolmogorov–Smirnov (K–S) test. When more than one mutational signature was considered at a time, the Benjamini–Hochberg (B–H) correction was applied to control for false positive discoveries. Correlations between signatures and mutation counts per gene were evaluated using Kendall τ rank correlation (with B–H correction) to account for ties (both measures were derived from amino-acid sequence changing mutations only). Significantly correlated genes were subsequently annotated as (i) tumor suppressor gene using the tumor suppressor gene database website version 2.0⁴³ (accessed 1st September 2022) and (ii) as oncogenes when present in the Cancer Gene Census COSMIC v.96⁴⁴ (accessed 16th September 2022) and tested for enrichment. For differential gene expression on the transcripts with a variant, we used a non-parametric B–HK–S test because for some genes no mutated transcripts were measured (these zeros were not possible to analyse with DESeq2). We use q to denote the multiple testing corrected p -values and a significance level $\leq 5\%$ was considered statistically significant.

Efficacy classifier. Patients were labelled as durable benefit [progression-free survival (PFS) $\geq \frac{1}{2}$ year] or non-durable benefit (PFS $< \frac{1}{2}$ year), whenever censoring permitted unambiguous label assignment. All patients were included in PFS analysis. In view of the limited overall survival (OS) data available in the Discovery cohort, we did not consider OS as an alternative endpoint (to PFS) for our classifier. To predict DB, we used a naive Bayes classifier: a classic supervised machine learning method^{45,46} that works particularly well with few samples, even when its conditional independence assumption is violated⁴⁷. Features were modelled with a zero-inflated exponential distribution. Results reported on the discovery cohort were obtained by leave-one-out cross-validation while inference on the validation set was done after training on the entire discovery cohort.

Estimates a and ninety-five per cent confidence intervals [$a - b, a + c$] of the average precision, area under the receiver operating characteristic curve (ROC AUC), F_1 score, sensitivity, and specificity were estimated by bootstrapping for 1000 iterations and denoted as a_{-b}^{+c} . Head-to-head model comparisons were evaluated using a paired permutation test (PPT), while performance estimates of different sets were compared using an unpaired permutation test (UPT).

For the PFS analysis (including all patients), we visualised the (out-of-fold) predicted versus actual PFS outcome using the Kaplan–Meier method. To quantify agreement, hazard ratios and significance were evaluated using Cox regression (for which we tested appropriateness). To compare hazard ratios h_r , corresponding coefficients (i.e., $-\ln h_r$), were compared with a regression coefficient test⁴⁸.

Data ethics statement. The discovery set was collected as part of pan-cancer studies CPCT-02, DRUP and WIDE described in³⁷, was evaluated and approved by the medical ethical committees of University Medical Center Utrecht and the Netherlands Cancer Institute and was executed in accordance with the relevant guidelines and regulations. All data contained in the discovery set was obtained through written informed consent of all the subjects for the purpose of whole genome sequencing and data sharing for cancer research.

Data availability

The discovery data that support the findings of this study are available from Hartwig Medical Foundation (data request # 094) but restrictions apply to the availability of these data, which were used under license for the current study, and so are not publicly available. Data are however available from the corresponding author on reasonable request and with permission of Hartwig Medical Foundation. The validation set can be extracted from Ref.⁴. Code and analysis notebooks are publicly available on <https://gitlab.com/hylkedonker/genomic-scars-predictor-nsclc-immunotherapy> under the MIT license.

Received: 24 October 2022; Accepted: 28 March 2023

Published online: 21 April 2023

References

- Camidge, D. R., Doebele, R. C. & Kerr, K. M. Comparing and contrasting predictive biomarkers for immunotherapy and targeted therapy of NSCLC. *Nat. Rev. Clin. Oncol.* **16**, 341–355 (2019).
- Sholl, L. M. *et al.* The promises and challenges of tumor mutation burden as an immunotherapy biomarker: A perspective from the International Association for the Study of Lung Cancer Pathology Committee. *J. Thorac. Oncol.* **15**(9), 1409–1424 (2020).
- Duchemann, B. *et al.* Current and future biomarkers for outcomes with immunotherapy in non-small cell lung cancer. *Transl. Lung Cancer Res.* **10**(6), 2937 (2021).
- Miao, D. *et al.* Genomic correlates of response to immune checkpoint blockade in microsatellite-stable solid tumors. *Nat. Genet.* **50**(9), 1271–1281 (2018).
- Paz-Ares, P. *et al.* Pembrolizumab (pembro) plus platinum-based chemotherapy (chemo) for metastatic NSCLC: Tissue TMB (tTMB) and outcomes in KEYNOTE-021, 189, and 407. *Ann. Oncol.* **30**, v917–v918 (2019).
- Alexandrov, L. B. *et al.* Signatures of mutational processes in human 475 cancer. *Nature* **500**(7463), 415–421 (2013).
- Kucab, J. E. *et al.* A compendium of mutational signatures of environmental agents. *Cell* **177**(4), 821–836 (2019).
- Kim, Y.-A. *et al.* Mutational signatures: From methods to mechanisms. *Annu. Rev. Biomed. Data Sci.* **4**(1), 189–206 (2021).
- Brady, S. W., Gout, A. M. & Zhang, J. Therapeutic and prognostic insights from the analysis of cancer mutational signatures. *Trends Genet.* **1**, 1–10 (2021).
- Alexandrov, L. B. *et al.* The repertoire of mutational signatures in human cancer. *Nature* **578**(7793), 94–101 (2020).
- Vöhringer, H. *et al.* Learning mutational signatures and their multidimensional genomic properties with TensorSignatures. *Nat. Commun.* **12**(1), 1–16 (2021).
- Nik-Zainal, S. *et al.* The genome as a record of environmental exposure. *Mutagenesis* **30**(6), 763–770 (2015).
- Brash, D. E. UV signature mutations. *Photochem. Photobiol.* **91**(1), 15–26 (2015).

14. Boot, A. *et al.* In-depth characterization of the cisplatin mutational signature in human cell lines and in esophageal and liver tumors. *Genome Res.* **28**(5), 654–665 (2018).
15. Chan, K. *et al.* An APOBEC3A hypermutation signature is distinguishable from the signature of background mutagenesis by APOBEC3B in human cancers. *Nat. Genet.* **47**(9), 1067–1072 (2015).
16. Wang, S. *et al.* APOBEC3B and APOBEC mutational signature as potential predictive markers for immunotherapy response in non-small cell lung cancer. *Oncogene* **37**(29), 3924–3936 (2018).
17. Chong, W. *et al.* Association of clock-like mutational signature with immune checkpoint inhibitor outcome in patients with melanoma and NSCLC. *Mol. Ther. Nucleic Acids* **23**, 89–100 (2021).
18. Nik-Zainal, S. *et al.* Mutational processes molding the genomes of 21 breast cancers. *Cell* **149**(5), 979–993 (2012).
19. Anagnostou, V. *et al.* Multimodal genomic features predict outcome of immune checkpoint blockade in non-small-cell lung cancer. *Nat. Cancer* **1**(1), 99–111 (2020).
20. Steele, C. D. *et al.* Undifferentiated sarcomas develop through distinct evolutionary pathways. *Cancer Cell* **35**(3), 441–456 (2019).
21. Steele, C. D. *et al.* Signatures of copy number alterations in human cancer. *Nature* **1**, 1–8 (2022).
22. The Cancer Genome Atlas Research Network. Comprehensive molecular profiling of lung adenocarcinoma. *Nature* **511**, 543–50 (2014).
23. The Cancer Genome Atlas Research Network. Comprehensive genomic characterization of squamous cell lung cancers. *Nature* **489**(7417), 519 (2012).
24. Degasperi, A. *et al.* Substitution mutational signatures in whole-genome-sequenced cancers in the UK population. *Science* **376**(6591), 9283 (2022).
25. Chen, J.-M., Férec, C. & Cooper, D. N. Patterns and mutational signatures of tandem base substitutions causing human inherited disease. *Hum. Mutat.* **34**(8), 1119–1130 (2013).
26. Rizvi, N. A. *et al.* Mutational landscape determines sensitivity to PD1 blockade in non-small cell lung cancer. *Science* **348**(6230), 124–128 (2015).
27. Li, B. *et al.* Therapy-induced mutations drive the genomic landscape of relapsed acute lymphoblastic leukemia. *Blood* **135**(1), 41–55 (2020).
28. Petljak, M. *et al.* Characterizing mutational signatures in human cancer cell lines reveals episodic APOBEC mutagenesis. *Cell* **176**(6), 1282–1294 (2019).
29. Donker, H. C. *et al.* Decoding circulating tumor DNA to identify durable benefit from immunotherapy in lung cancer. *Lung Cancer* **170**, 52–57 (2022).
30. Norum, J. & Nieder, C. Tobacco smoking and cessation and PD-L1 inhibitors in non-small cell lung cancer (NSCLC): A review of the literature. *ESMO Open* **3**(6), e000406 (2018).
31. Vickers, A. J., Van Calster, B. & Steyerberg, E. W. Net benefit approaches to the evaluation of prediction models, molecular markers, and diagnostic tests. *BMJ* **1**, 352 (2016).
32. Vickers, A. J., Van Calster, B. & Steyerberg, E. W. A simple, step-by-step guide to interpreting decision curve analysis. *Diagn. Progn. Res.* **3**(1), 1–8 (2019).
33. Talluri, R. & Shete, S. Using the weighted area under the net benefit curve for decision curve analysis. *BMC Med. Inf. Decis. Mak.* **16**(1), 1–9 (2016).
34. Stelzer, G. *et al.* The GeneCards suite: From gene data mining to disease genome sequence analyses. *Curr. Protoc. Bioinform.* **54**(1), 1–30 (2016).
35. Safran, M. *et al.* The genecards suite. In: *Practical Guide to Life Science Databases*, 27–56 (2021).
36. Litchfield, K. *et al.* Escape from nonsense-mediated decay associates with anti-tumor immunogenicity. *Nat. Commun.* **11**(1), 1–11 (2020).
37. Priestley, P. *et al.* Pan-cancer whole-genome analyses of metastatic 558 solid tumours. *Nature* **575**(7781), 210–216 (2019).
38. Cameron, D. L. *et al.* GRIDSS, PURPLE, LINX: Unscrambling the tumor genome via integrated analysis of structural variation and copy number. *BioRxiv* <https://doi.org/10.1101/781013> (2019).
39. Zhao, H. *et al.* CrossMap: A versatile tool for coordinate conversion between genome assemblies. *Bioinformatics* **30**(7), 1006–1007 (2014).
40. COSMIC. *COSMIC-Mutational Signatures*. <https://cancer.sanger.ac.uk/signatures>. Accessed 17 May 2022.
41. Bergstrom, E. N. *et al.* SigProfilerMatrixGenerator: A tool for visualizing and exploring patterns of small mutational events. *BMC Genom.* **20**(1), 1–12 (2019).
42. Coffey, A. J. *et al.* The GENCODE exome: Sequencing the complete human exome. *Eur. J. Hum. Genet.* **19**(7), 827–831 (2011).
43. Zhao, M. *et al.* TSGene 2.0: An updated literature-based knowledgebase for tumor suppressor genes. *Nucleic Acids Res.* **44**(D1), D1023–D1031 (2016).
44. Sondka, Z. *et al.* The COSMIC Cancer Gene Census: Describing genetic dysfunction across all human cancers. *Nat. Rev. Cancer* **8**(11), 696–705 (2018).
45. Koller, D. & Friedman, N. *Probabilistic Graphical Models: Principles and Techniques* (MIT Press, 2009).
46. Murphy, K. P. *Probabilistic Machine Learning: An Introduction* (MIT Press, 2022).
47. Zhang, H. The optimality of naive Bayes. *Aa* **1**(2), 3 (2004).
48. Paternoster, R. *et al.* Using the correct statistical test for the equality of regression coefficients. *Criminology* **36**(4), 859–866 (1998).

Acknowledgements

We thank Fenneke Zwierenga en Benthe Muntinghe for proofreading. This publication and the underlying research are partly facilitated by Hartwig Medical Foundation and the Center for Personalized Cancer Treatment (CPCT) which have generated, analysed and made available data for this research.

Author contributions

Conceptualization: H.C.D., B.v.E., M.T., G.A.L., H.J.M.G.; Methodology: H.C.D.; Software: H.C.D.; Validation: H.C.D.; Formal analysis: H.C.D., B.v.E.; Investigation: H.C.D.; Data Curation: H.C.D., B.v.E.; Writing Original Draft: H.C.D.; Writing Review & Editing: H.C.D., B.v.E., M.T., G.A.L., L.C.L.T.v.K., E.S., T.J.N.H., H.J.M.G.; Visualization: H.C.D., B.v.E., T.J.N.H.; Supervision: H.J.M.G.

Funding

This material is based upon work supported by the Google Cloud Research Credits program with the Award GCP19980904.

Competing interests

HCD: None to declare; BvE: None to declare; MT: None to declare; GAL: None to declare; ES: Honoraria/speakers fee: Bio-Rad, Roche, Agena Bioscience, Illumina, Lilly; Consulting or Advisory Role: MSD/Merck, Astellas, Bayer, BMS, Agena Bioscience, Janssen Cilag (Johnson & Johnson), Novartis, Roche, AstraZeneca, Amgen, Lilly; Research Funding: Biocartis, Bio-Rad, Roche, Agena Bioscience, AstraZeneca, InVitaE/Archer (all paid to UMCG); Travel, Accommodations, Expenses: Roche Molecular Diagnostics, Bio-Rad. LCLTvK: Grants, non-financial support from Roche, advisory board presence for AstraZeneca, Novartis, Merck, Janssen-Cilag, Bayer, BMS, nanoString and Pfizer, grants and non-financial support from Invitae, non-financial support from Biocartis, grants from Bayer, non-financial support from nanoString. TJNH: Advisory/consultancy fees from AstraZeneca, Bristol-Myers-Squibb, Merck Sharp Dohme, Roche, and research grants/funding from AstraZeneca, Hoffmann-La Roche. HJMG: Consulting or Advisory Role: Novartis, Lilly, Roche/Genentech.

Additional information

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1038/s41598-023-32499-3>.

Correspondence and requests for materials should be addressed to B.E.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2023