

Statistical Methods for Identifying Sequence Motifs Affecting Point Mutations

Yicheng Zhu¹, Teresa Neeman², Von Bing Yap³, and Gavin Huttley¹

¹Research School of Biology, The Australian National University, Canberra ACT 2601, Australia

²Statistical Consulting Unit, The Australian National University, Canberra ACT 2601, Australia

³Department of Statistics and Applied Probability, National University of Singapore, 117546, Singapore

ABSTRACT

Mutation processes differ between types of point mutation, genomic locations, cells, and biological species. For some point mutations, specific neighbouring bases are known to be mechanistically influential. Beyond these cases, numerous questions remain unresolved including: what are the sequence motifs that affect point mutations? how large are the motifs? and, do they vary between samples? We present new log-linear models that allow explicit examination of these questions along with sequence logo style visualisation to enable identifying specific motifs. We demonstrate the utility of these methods by analysing human germline and malignant melanoma mutations. We recapitulate the known CpG effect and identify numerous novel motifs, including a highly significant motif associated with A→G mutations. We show that major effects of neighbourhood on germline mutation lie within ± 2 of the mutating base. Models are also presented for contrasting the entire mutation spectra (the distribution of the different point mutations) and applied to the data. We show the spectra vary significantly between autosomes and X-chromosome, with a difference in T→C transition dominating. Analyses of malignant melanoma confirmed reported characteristic features of this cancer including strand asymmetry and markedly different neighbouring influences. The methods reported are made freely available as a Python library <https://bitbucket.org/gavin.huttley/mutationmotif>.

Keywords: log-linear model, context dependent mutation, germline mutation, somatic mutation

INTRODUCTION

Understanding the contributions of mutation processes to genetic diversity has broad relevance to topics ranging from estimating genetic divergence (Huttley, 2004; Schluter, 2009; Harris, 2015) to the aetiology of disease (Peltomaki and Vasen, 1997; Ying and Huttley, 2011; Nik-Zainal et al., 2012; Alexandrov et al., 2013a). While mutations occur on many scales, from single nucleotide point mutations to substantial genomic rearrangements, we restrict our attention here to point mutation processes. A multitude of mechanisms have been characterised that cause DNA lesions (Cooke et al., 2003; Helleday et al., 2014). Similarly, an array of processes repairing DNA lesions have also been described (Helleday et al., 2014). From examination of sequence composition alone it is apparent that mechanisms of mutagenesis (lesion formation and subsequent failure of DNA repair) differ between genomic locations (Francioli et al., 2015), between cell types (Nishino et al., 1996) and between species (Karlin et al., 1998). In evaluating natural systems, where only the starting and ending sequence states may be known, establishing the mechanistic origins remains a challenge. In mammals, an informative exception is the case of C→T point mutations. In this instance, a 3' G strongly implies a mechanism of 5-methyl-cytosine (5mC) deamination. This is due to the binding affinity of DNA methylases for the CpG sequence motif (Vinson and Chatterjee, 2012) and the greatly elevated mutation rate of 5mC (Coulondre et al., 1978). As the CpG example illustrates, predicting the contribution of a specific mechanism requires knowledge of a characteristic mutation sequence signature. Motivated by this, we focus here on development of a statistical method and associated visualisation approach for revealing signature sequence motifs associated with point mutations. We refer to these as mutation motifs.

Considerable evidence indicates that the influence of neighbouring bases on point mutations is a general phenomenon. Early studies on inherited, and thus germline, mutations in humans supported the hypermutability of the CpG dinucleotide as the dominant origin of C→T mutations (Cooper, 1995). Subsequent work further suggested that the remaining 11 point mutations are also affected by neighbouring bases (Krawczak et al., 1998). From analyses of mutations in human disease genes, (Krawczak et al., 1998) inferred the influence of neighbours are confined to the positions immediately flanking the mutated location. The work on human polymorphisms demonstrated these results applied more generally across the genome (Zhao and Boerwinkle, 2002). Recently, using trinucleotides where the mutated base is central, distinctive mutation signatures that discriminate human cancer types have been identified (Alexandrov et al., 2013a). These results demonstrate the influence of neighbouring bases generalises to somatic mutations. Early influential work on plant cpDNA completes the demonstration of the generality of neighbouring influences across the tree-of-life (Morton et al., 1997). While Krawczak et al. (1998) and Zhao and Boerwinkle (2002) identified the influence of neighbours is proportional to distance, the work of Alexandrov et al. (2013a) was focused on the immediate flanking bases.

The influence of neighbouring bases on mutagenesis can have multiple causes. The chemical properties of DNA alone can confer a neighbourhood influence on mutation susceptibility. Adjacent pyrimidines are vulnerable to a dimerisation in the presence of UV light (Brown, 2002, p 426) with TpT being most susceptible. As the influence of DNA methylase preference for CpG dinucleotides demonstrates, DNA binding properties of macromolecules are a further likely source of neighbouring base influences. With numerous DNA-protein binding interactions central to DNA repair processes, any affinity to specific sequence motifs of these molecules may result in those motifs being under-represented in mutated sequences.

Analysis techniques for estimation of neighbouring base influences on mutation draw on different approaches. Krawczak et al. (1998) quantified neighbouring base influence by contrasting observed base frequencies against an equiprobable frequency distribution via a Euclidean distance. Zhao and Boerwinkle (2002) used just the base frequencies per position except beyond $\pm 10\text{bp}$ where averages across position ranges were used. In both these approaches, the background sequence distribution is assumed to be random occurrence of bases. These approaches therefore potentially obscure the real signal by confounding it with the non-random occurrence of bases characteristic of DNA sequences.

The distinctive mutagenic biology of cancer has motivated development of methods to identify specific mutation signatures across all point mutations. The related methods of Alexandrov et al. (2013b) and Shiraishi et al. (2015) tackle the problem of resolving the signatures of different mutational processes. As these signatures consist of all point mutations, they are a composite of distinct underlying mutational processes operating across all point mutations. The different mutation signatures may, therefore, contain component(s) that are identical. The methods are therefore not well suited to examining the influence of neighbouring bases on single point mutation directions.

More recently, Aggarwala and Voight (2016) examined the influence of sequence neighbourhood by using a probability of polymorphism that was conditioned on the sequence context. They identified 7-mer contexts as accounting for a median of 81% of the variability in polymorphism rate across point mutations. Their results indicated inclusion of higher order (three-way and higher) interactions accounted for as much as 50% of this predictive power. However, k-mers exhibit a non-random distribution within the human genome (Karlin, 1998; Chor et al., 2009). Moreover, variation in sequence composition is correlated with variation in substitution rate (e.g. Hodgkinson and Eyre-Walker, 2011). These suggest by averaging across all occurrences of the sequence context, the results of Aggarwala and Voight (2016) reflect the relationship between genomic location and polymorphism rate rather than the influence of neighbours.

Detection of functional sequence motifs is a related problem to which information theoretic techniques have been extensively applied. Mutual information (MI) per position in a sequence alignment is computed by subtracting the positions Shannon's entropy from entropy of the uniform distribution (Shannon, 2001). Coupling of this metric with the sequence logo visualisation approach has led to its widespread application for discovery of functional motifs (Schneider and Stephens, 1990). The display used the MI statistic to define a stack of colour coded letters, representing the sequence states, with each letters height scaled proportional to its contribution to the total MI (Schneider and Stephens, 1990). While MI has many appealing properties for measuring information, it shares the restriction of comparison to the independent equiprobable distribution. As removing the constraint of equal frequencies can lead to cases of taking the

logarithm of a negative number, which is not defined, MI is not appropriate for examination of most DNA sequences as the equiprobable property typically does not hold.

Many of the developed techniques are confounded by common signatures of genome DNA sequences – nucleotides do not occur with equal frequency or randomly. For the genomes of many organisms, such as vertebrates, there is also considerable within genome variation in k-mer frequencies (Karlin, 1998; Chor et al., 2009). These factors will contribute substantial noise to any statistic that does not account for them. Most available methods do not distinguish contributions from independent positions compared with joint contributions from multiple positions. For instance, are mutations affected by the sequence of bases present at two positions (Zhang and Mathews, 1995)? General statistical techniques have been developed to characterise genome sequence signatures that incorporate such between position interactions (Karlin et al., 1998). While this statistic has been modified for inferring neighbouring influences from sequence substitutions (Nevarez et al., 2010), the approach has limitations compared with more general alternatives. Log-linear models allow flexible parameterisations for hierarchical hypothesis testing of categorical data and have been previously applied to examination of neighbouring influences (Huttley et al., 2000). Their generality allows for controlling of potential confounding differences, such as differences in sample size and nucleotide composition. The support for comparing hypotheses in a hierarchical manner enables identifying parametrically succinct models that explain the data well.

In this study, we develop log-linear approaches for examination of mutation processes. Our work is distinguished from previous methods by conditioning on the mutation event, rather than the sequence context, and employs a control distribution that is matched for genomic location. We present hierarchical hypothesis tests for evaluating whether: (i) neighbouring bases influence mutation direction, (ii) neighbouring base influences are equal between samples, and (iii) the spectrum of mutations (the relative abundance of the 12 point mutations) are equal between samples. A sequence logo inspired visualisation approach is also presented. We demonstrate application of the models by applying them to data previously reported to exhibit distinctive mutation processes; namely, germline mutations in different sequence classes (e.g. transcribed, untranscribed) and chromosome classes (e.g. autosome and sex-chromosome), and somatic mutations in cancer. Data were human SNPs obtained from Ensembl for both germline and somatic mutations. In addition to replicating the well known CpG effect, our results indicate that neighbourhood size can be quite large and, as we demonstrate for the A→G transition mutation, the influence of neighbours does not decay monotonically with distance. We further show, that both independent and dependent position influences contribute to mutational process. Through formal testing of equivalence between samples, we demonstrate significant differences between sequence classes, chromosome classes and between melanoma and germline mutations. Software implementing all these methods, released under the GPL open source license, is made available <https://bitbucket.org/gavin.huttley/mutationmotif>.

MATERIALS AND METHODS

Data sampling

We sampled SNP data and flanking sequences from Ensembl (Flicek et al., 2013) release 79 using PyCogent's Ensembl querying capabilities (Knight et al., 2007). The Ensembl variation database records whether a variant is classified as somatic. We sampled germline SNPs using that flag and required the Ensembl record indicated the SNP was validated, had an inferred ancestral allele and that its flanking sequence matched the reference genome. For each such filtered SNP, we recorded the alleles, strand, sequence class (exonic, intronic or intergenic), genomic coordinates and 300bp of flanking sequence either side of the SNP location.

Sampling somatic SNPs involved both the COSMIC (Forbes et al., 2015) and Ensembl databases. Complete mutant export data was obtained from COSMIC, which included SNP identifiers and the primary pathology from which a SNP had been reported. Flanking sequence was derived by obtaining the Ensembl records for the SNP identifiers, ensuring the record was flagged as somatic and then following the same procedure as for the germline variants. We restricted our attention to SNPs obtained from malignant melanoma.

Determining base counts

For each mutation direction (e.g. C→T) we obtained base counts from paired mutated and reference base locations from the same genomic fragments. Neighbour positions were indexed relative to the position of the chosen location. For a mutated base, the chosen location was the annotated site of the SNP (Fig 1).

For the reference base, the chosen location was derived from the same genomic fragment by randomly selecting from among the positions that had the same starting base as that affected by the mutation (e.g. a random choice of a position with a C in the case of a C→T mutation), but excluded the SNP location. For each mutation direction, for each sampled genomic segment, a 5bp long sequences with the chosen location at centre was extracted and the bases observed per position relative to the chosen location were recorded. As the total number of possible neighbourhoods was 256, a single file was written with counts for each of the possible neighbourhoods for both the mutated and reference locations.

C→T	T	G	A	G	C	C	G	G	G	C	A
	-5	-4	-3	-2	-1	0	1	2	3	4	5
Random C	C	T	G	G	G	C	A	T	G	A	G
	-1	0	1	2	3	4	5	-5	-4	-3	-2

Figure 1. Specifying neighbourhoods. The neighbourhood of a position at which a C→T mutation occurred is compared with the neighbourhood of a randomly selected occurrence of C from the same sequence. The location of the C→T SNP is the central position for the mutated base and is assigned the index 0. The C at position 4 was randomly chosen as the reference location and the sequence is shifted so it is centred on this position (see ‘Determining base counts’ for fuller explanation).

Log-linear modelling of neighbour effects

We first demonstrate the general approach of applying log-linear models for understanding neighbour influences on mutation by focusing on the influence of a single neighbouring position. We then consider the extension of comparing neighbour contributions between samples. Both of these analyses are concerned with the independent contribution of bases at a position to mutation status.

For a single position, we evaluate whether *base* and mutation *status* occur independently using a straightforward log-linear model. Under the most saturated log-linear model, the log of the expected frequency f_{ij} for *base* i and mutation *status* j can be expressed as

$$\ln f_{ij} = \lambda + \lambda_i^{base} + \lambda_j^{status} + \lambda_{ij}^{base:status} \quad (1)$$

where λ represents the intercept (i.e., common to all counts), λ_i^{base} , the contribution to the frequency of being *base* i , λ_j^{status} the contribution to the frequency of being mutation *status* j , and the interaction between *base* and *status* $\lambda_{ij}^{base:status}$. The latter expresses the degree of non-independence between *base* and mutation *status*. The number of levels for each factor are: *base*, 4 levels (A, C, G, T); and mutation *status*, 2 levels (M and R). The fit of a log-linear model is measured as the deviance (D). We specify the null hypothesis that bases occur independent of mutation status by setting $\lambda_{ij}^{base:status} = 0$. The alternate is the fully saturated model. The difference in D between the null and alternate, nested models, is taken as χ^2 with degrees of freedom equal to the difference in the number of free parameters. In this instance, the degrees of freedom is 3.

When comparing groups, e.g. autosome versus X-chromosome, we add another factor (λ_j^{group}) to the log-linear model (2). The fully parameterised version of this log-linear model requires addition of further 3 interaction parameters: 2 two-way interactions and the three-way interaction parameter $\lambda_{ij}^{base:status:group}$. This parameter represents the influence of group on the *base : status* interaction. We therefore evaluate the null hypothesis of no difference between samples by setting all $\lambda_{ij}^{base:status:group} = 0$ and compare this against the fully saturated model. If the group factor has only 2 levels, then the degrees of freedom for the resulting D is 3.

$$\begin{aligned} \ln f_{ij} = & \lambda + \lambda_i^{base} + \lambda_j^{status} + \lambda_j^{group} \\ & + \lambda_{ij}^{base:status} + \lambda_{ij}^{base:group} + \lambda_{ij}^{group:status} \\ & + \lambda_{ij}^{base:status:group} \end{aligned} \quad (2)$$

We now extend this approach to consider the simultaneous influence on mutation status of bases at multiple positions. To illustrate, consider the two neighbours following the base C in Fig 1. There are

sixteen possible dinucleotides at the 1,2 positions. The goal of this model is to establish whether the dinucleotides at these two positions affect mutation status of C after taking account of the independent contributions of these positions. In order to achieve this, our two-position interaction model extends the independent contribution model (1), adding factors for the additional position and then interaction terms between the parameters. The fully saturated two-position interaction model is

$$\ln f_{ijk} = \lambda + \lambda_i^{base_1} + \lambda_j^{base_2} + \lambda_k^{status} + \lambda_{ik}^{base_1:status} + \lambda_{jk}^{base_2:status} + \lambda_{ij}^{base_1:base_2} + \lambda_{ijk}^{base_1:base_2:status} \quad (3)$$

where λ^{base_1} and λ^{base_2} represent the base contributions at positions one and two. In addition to including factors for the independent contributions of the two positions on mutation status, the $\lambda^{base_1:base_2}$ accounts for non-independent occurrence of bases at the positions, a key property of DNA sequences. The null hypothesis of no interaction between dinucleotides and mutation status is specified by setting all $\lambda^{base_1:base_2:status} = 0$ and comparing this against the fully saturated model. The resulting D has 9 degrees of freedom. For a given mutation direction, we perform this analysis for all possible combinations of pairs of sites.

These approaches are further extended to consider interactions amongst three positions, amongst four positions and for comparison of these effects amongst groups.

Log-linear model of mutation spectra

For analysis of mutation spectra, we evaluate the null hypothesis that the distribution of mutations is the same between groups. The opportunity for a specific mutation direction is affected by the total occurrence of the starting base. This quantity can be difficult to ascertain, such as in cancers where there may be major genomic rearrangements (e.g. deletions) relative to a reference group. To avoid this uncertainty, we restrict the analysis to point mutations from a specific base, comparing the relative counts of each of the 3 possible mutations between groups. This is a test of independence between ending base and group.

For a specific base, the log of the expected frequency is defined as

$$\ln f_{ij} = \lambda + \lambda_i^{direction} + \lambda_j^{group} + \lambda_{ij}^{direction:group} \quad (4)$$

where the factor $\lambda^{direction}$ represents the counts of the 3 different point mutation directions, λ^{group} the counts in the different groups, and $\lambda^{direction:group}$ the interaction between these factors. We specify the null hypothesis of equivalent proportions between the groups by setting $\lambda^{direction:group} = 0$. For two groups, comparing against the fully saturated model, the D has 2 degrees of freedom.

Visualisation

Sequence logo's display motifs using the mutual information as the letter stack height, and the fraction contributed to the mutual information (MI) by individual bases is derived from their individual terms in the MI calculation. We adopt a similar approach here. Instead of using MI, we use relative entropy (RE). The log likelihood ratio D is converted to RE by dividing by twice the sample size. RE from a log-linear analysis specifies the letter stack height. We use the terms in the RE equation to determine the proportion of the stack height attributable to a specific base. We differ from the conventional sequence logo approach by distinguishing between bases that are under or over represented in the mutated class, relative to the unmutated class. Under-represented bases are indicated by a 180° rotation.

Interpretation of the logo is straightforward. A higher RE value indicates that a position(s) has a greater influence on mutation. Support for concluding a stack height reflects a meaningful influence on mutation derives from the p-value, from the log-linear model, that the null hypothesis is correct. The magnitudes and orientations of letters further conveys meaning in that ordinary letter orientation is indicative of over representation in the mutated group while inverted orientation indicates under representation. We note here that we make a choice to use residuals from the mutated class for display. Using residuals from the unmutated class would generate an image with the opposite letter orientations.

For multi-position models (e.g. 3), the stack height is equal between the indicated positions. For the two-position model, the characters for the nucleotide pair at the two positions share the same proportion

and orientation. For the more complicated analyses involving contrasting neighbour effects between groups, the reference category is the one provided first to the software.

Differences in mutation spectra are visualised using a grid with rows corresponding to the starting base and columns the base resulting from the mutation. Each row corresponds to a single log-linear test for equivalent distribution of the possible point mutations from the base indicated by the row label (see Log-linear model of mutation spectra). The RE for each row is computed from the deviance of the corresponding spectra test. Letter heights for each base are scaled proportional to the corresponding term in the RE equation. The sum of letter heights in a row is the total RE for that test. Bases over-represented in the reference group are oriented in the conventional manner while under-represented bases are rotated 180°. In the spectra analysis, the largest base in the grid is the dominant mutation product difference between the groups.

Availability of data and materials

MutationMotif is a Python 2.7 compatible library for performing the statistical analyses outlined in this work that will be made freely available under the GPL. The project homepage is at <https://bitbucket.org/gavin.huttley/mutationmotif> and the version employed for the reported work is available in Zenodo (DOI 10.5281/zenodo.53215). It draws on R (Ihaka and Gentleman, 1996) for log-linear modelling, via the glm function, using the rpy2 Python binding to R. Sequence logo's are drawn using custom Python code included in MutationMotif. Other dependencies include PyCogent (Knight et al., 2007), pandas, numpy, matplotlib and scitrack.

The scripts performing the data sampling and applying the analyses reported in this work will be made freely available under the GPL at <https://bitbucket.org/gavin.huttley/analysemutations> and the version employed for the reported work is available in Zenodo (DOI 10.5281/zenodo.53220). AnalyseMutations includes the counts data required by MutationMotif and the complete set of results contained in this work. These counts data were produced from data sampled from the Ensembl and COSMIC databases, as described in Data sampling. Because the data files from which the compact counts files were produced are so large, they are available separately in Zenodo (DOIs 10.5281/zenodo.53158 <https://zenodo.org/record/53158> and 10.5281/zenodo.53164 <https://zenodo.org/record/53164>) under the Creative Commons Attribution-Share Alike license. Data files are typically gzip compressed standard formats; tab delimited text files, fasta formatted sequence files, serialised data is stored as json or pickle (Python's native serialised format).

RESULTS

Overview of notation and neighbour effect log-linear models

The notation $X \rightarrow Y$ refers to a point mutation from X to Y, $X \rightarrow Y^*$ refers to a point mutation and its strand symmetric counterpart, e.g. $C \rightarrow T^*$ is $C \rightarrow T$ or $G \rightarrow A$.

The log-linear model of neighbour influence evaluates the null hypothesis that a neighbouring base(s) flanking a specific point mutation is the same as that flanking a random occurrence of the starting base. For instance, does the distribution of bases at sites flanking $C \rightarrow T$ mutations differ from that flanking all C's? As the frequency of bases varies between genomic locations (Bernardi, 2000), the mutated and reference locations need to be matched to avoid possible confounding. We achieve this matching by deriving a reference location proximal to each mutated location. The sampling process is shown in Fig 1. We sampled 300bp of flanking genomic sequence each side of a SNP and within this segment chose, at random, another occurrence of the starting base affected by the mutation event. Unless stated otherwise, we limited our analysis of neighbouring influence to ± 2 bp either side of the mutated position, resulting in 256 possible neighbourhoods. For any given mutation direction, counts of these different neighbourhoods are obtained from both the sample centred on the mutated base and the sample centred on a random occurrence of the starting base. These counts are used to construct the contingency tables for the log-linear analysis. This approach achieves the objectives of controlling for compositional variation across the genome and controlling for the non-random occurrence of bases. See 'Determining base counts' for more detail on this procedure.

The log-linear models used to examine the effect of neighbours on point mutation include parameters that represent an interaction between neighbouring base(s) and mutation status (see). The contribution of this parameter to model fit is measured as a Deviance which, along with the residual degrees-of-freedom, is to calculate the corresponding p-value for the null hypothesis. We convert the Deviance to relative

entropy (hereafter, RE) as this measures the information content of the data under the model in a manner that is robust to sample size, allowing comparisons among analyses.

As we are concerned with whether flanking positions individually or jointly affect mutation process we describe the influence of neighbouring bases as independent or dependent/joint effects respectively. The influence of a base at a single neighbouring position on a point mutation will be referred to as an “independent” effect. The case when bases at two or more neighbouring positions influence a point mutation will be referred to as a “dependent” interactive effect or the joint influence of multiple bases. We note here that in the case of a dependent effect the actual positions are not necessarily contiguous. The number of positions involved in a dependent effect is referenced as the “order” of the interaction. An independent effect, the influence of a single position on mutation, is a first order effect while the joint influence of two positions on mutation is a second order effect. Flanking locations are indexed relevant to the mutated position. The immediate flanking 5’ base is at position -1 while the immediate flanking 3’ base is at position $+1$ (see Fig 1). A series of positions are indicated by the relative indices in parentheses e.g. $(-2, -1)$ are two positions 5’ to the mutated base.

Log-linear models recapitulate the CpG effect and reveal higher order effects

In the analyses we report below, we focus principally on analyses of intergenic autosomal data. We also sampled SNPs from introns and exons. We relegate all results from analysis of other genomic regions to supplementary material as the results are substantively the same as those from the intergenic sequence class.

We benchmarked our method by examining the influence of neighbouring bases on C→T point mutations in the autosomal intergenic sample. (As none of the strand symmetry tests were significant for the intergenic autosomal mutations, we limit our discussion to the “plus” strand directions only.) We expected the influence of methylation induced deamination at CpG to reveal a strong G effect at the $+1$ position (Cooper and Youssoufian, 1988). This prediction was confirmed in the results of the hypothesis test (Table 1) and visually in the mutation motif logo (Fig 2 b). The analysis established that while all positions made highly significant independent contributions to mutation (Table 1) the magnitude of their influence was small compared to that at the $+1$ position and only one of these was evident in the mutation logo, that of A at the -1 position (Fig 2 b). (Results from the equivalent analysis of autosomal exon data are shown in Fig S1.)

Position(s)	Deviance	df	p-value
-2	1574.2	3	0.0
-1	18674.9	3	0.0
+1	346848.0	3	0.0
+2	2174.5	3	0.0
$(-2, -1)$	1603.1	9	0.0
$(-2, +1)$	555.5	9	0.0
$(-2, +2)$	352.7	9	1.7×10^{-70}
$(-1, +1)$	2341.3	9	0.0
$(-1, +2)$	315.1	9	1.6×10^{-62}
$(+1, +2)$	1965.0	9	0.0
$(-2, -1, +1)$	939.7	27	0.0
$(-2, -1, +2)$	523.0	27	2.7×10^{-93}
$(-2, +1, +2)$	264.6	27	7.3×10^{-41}
$(-1, +1, +2)$	467.8	27	6.5×10^{-82}
$(-2, -1, +1, +2)$	273.9	81	9.1×10^{-23}

Table 1. Log-linear analysis of C→T autosomal intergenic mutations. Position(s) are relative to the index position (see Figure 1). Deviance is from the log-linear model, with df degrees-of-freedom and corresponding p-value obtained from the χ^2 distribution. p-values listed as 0.0 are below the limit of detection.

Specific combinations of bases at multiple positions also significantly affected C→T mutations. All higher order interactions were statistically significant. A feature of the second and third order joint effects

was that bases immediately proximal to each other or to the mutated position had the strongest association: $(-2, -1)$, $(-1, +1)$, $(+1, +2)$ second order interactions (Table 1 and Fig 2 c), and the $(-2, -1, +1)$ third order interaction (Fig 2 d).

Despite the highly significant associations between combinations of positions and interactions, the independent position contributions dominated. All effect orders were significantly associated with mutation status even when using the sequential Holm-Šidák correction for 15 tests (Holm, 1979). These results reflect the enormous statistical power resulting from the large sample sizes, e.g. over 1 million C→T intergenic SNPs. Contrasting the magnitudes of these different effects by displaying the maximum RE value from each effect order (RE_{max} , Fig 2 a) provide a useful indicator of their relative influence; $RE_{max}(1)$ is the maximum RE score for first position effects across all positions (e.g. +1 in this case), $RE_{max}(2)$ the maximum RE score from combinations of two positions, and so on for the higher orders. This display established that the 3'-G influence dominates all other neighbouring base effects on C→T mutation. Furthermore, contrasting these values between the point mutations (Table 2) affirms that neighbourhood has the strongest effect on C→T mutations (Fig S2).

A→G mutations are also strongly affected by neighbours

The A→G transition mutation exhibited the next strongest influence of neighbouring bases (Table 2). As for C→T, all effect orders were highly significant after correcting for 15 tests (Table 3). All positions showed significant first order influences, but the $-2, -1, +1$ positions were particularly strong (Fig 3 b). Two of these, $(-2, -1)$, also exhibited a prominent second order interaction (Fig 3 c) while all three contributed the strongest third order interaction (Fig 3 d). For A→G mutations, our analysis indicated that while first order effects dominated, higher order effects were important factors affecting this mutation direction (Fig 3 a). (Results from the equivalent analysis of autosomal exon data are shown in Fig S3.)

Direction	$RE_{max}(1)$	Pos.(1)	$RE_{max}(2)$	Pos.(2)	$RE_{max}(3)$	Pos.(3)
A→C	0.0039	-1	0.0016	(+1, +2)	0.0012	(-2, -1, +1)
A→G	0.0188	+1	0.0030	(-2, -1)	0.0007	(-2, -1, +1)
A→T	0.0095	+1	0.0051	(-1, +1)	0.0023	(-1, +1, +2)
C→A	0.0091	+1	0.0044	(-1, +1)	0.0015	(-1, +1, +2)
C→G	0.0054	-2	0.0025	(+1, +2)	0.0010	(-1, +1, +2)
C→T	0.0860	+1	0.0006	(-1, +1)	0.0002	(-2, -1, +1)

Table 2. Summary of neighbourhood contributions to plus strand mutations with an autosomal intergenic location. $RE_{max}(\#)$ is the maximum RE for order # and Pos.(#) the corresponding position(s). All point mutations had at least one significant test after correcting for 15 tests (see Table 1) using the Holm-Šidák procedure.

Transversion mutations are affected by neighbours

All transversion mutations had significant neighbourhood influences but to a lesser extent than that evident for transition mutations (Table 2). The transversion mutations showed $RE_{max}(1)$ that were 20-fold less than for the C→T mutations. However, higher order effects were typically more pronounced for transversions than transitions. The A→T and C→A transversion mutations showed the greatest influence of neighbours at all levels. The dominant influences were immediately adjacent to the mutating base except for C→G, where position -2 had the strongest effect.

The size of the neighbourhood

Our analyses above indicated first order effects exerted the strongest influence on mutations. Accordingly, we limited our examination of neighbourhood size to first order effects and sampled intergenic autosomal SNPs with a flank size of ± 10 bp for an analysis. After correcting for multiple tests, all 20 flanking positions were significant for all point mutations (Table S1). This suggests a neighbourhood size ≥ 10 . The tendency for even very distant positions to be highly significant in this analysis likely reflects the enormous sample sizes employed for this analysis and does not necessarily reflect the magnitude of a positions influence. Therefore, for each mutation we estimated the most distant position with a RE that was $\geq 10\%$ of $RE_{max}(1)$. For the transition mutations, the neighbourhood size was restricted to positions

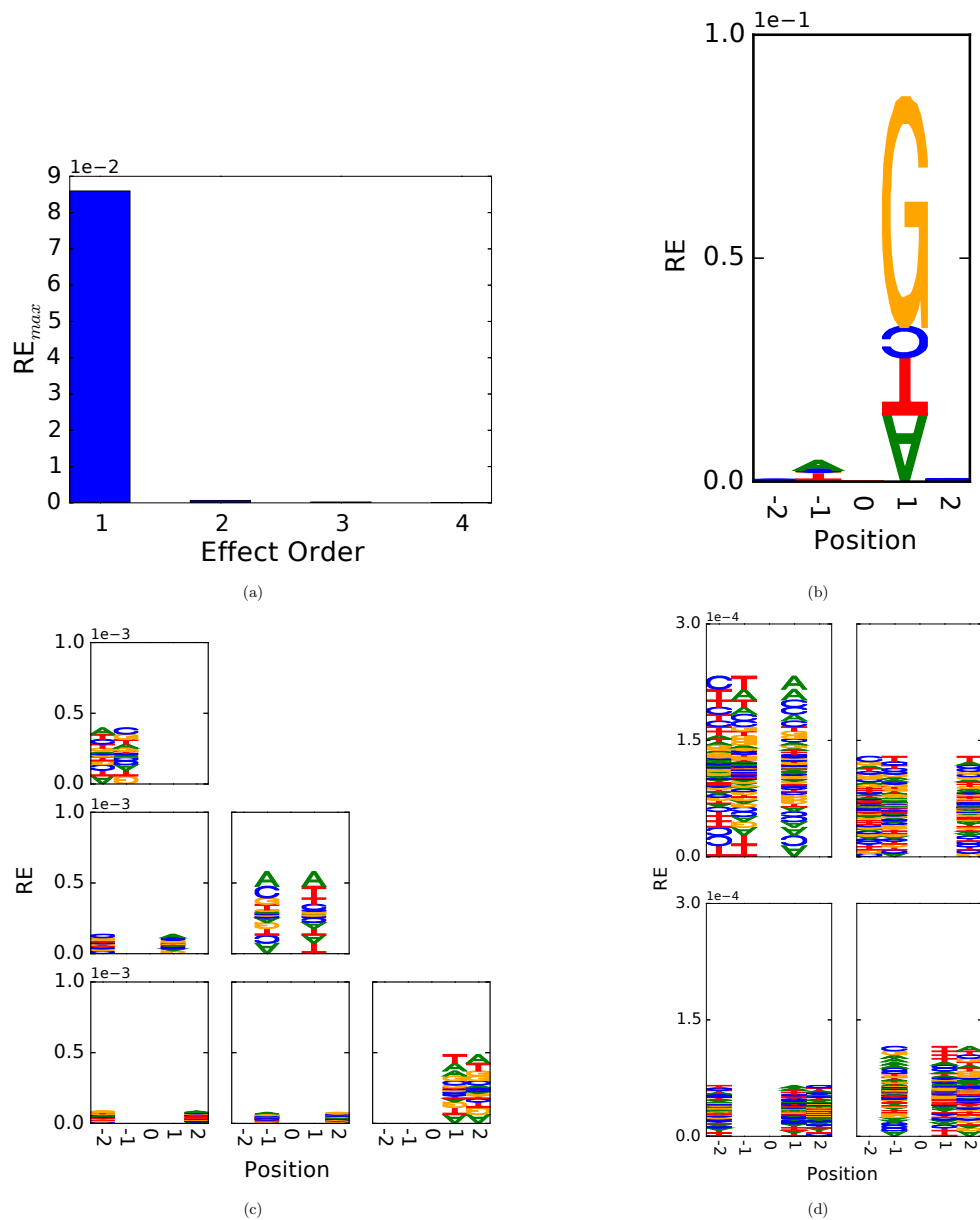


Figure 2. Flanking influences on C→T mutations. (a) First order effects are the dominant neighbourhood influence, RE_{max} (y-axis) is the maximum RE from the possible evaluations for a motif length (x-axis), (b) Single position effects, (c) Two-way effects, and (d) Three-way effects. For b-d, the y-axis is RE and the x-axis is the position index relative to the mutated base.

within $\pm 2bp$ (Fig S4) whilst for transversion mutations, the neighbourhood size was within $\pm 4bp$ (Table S1).

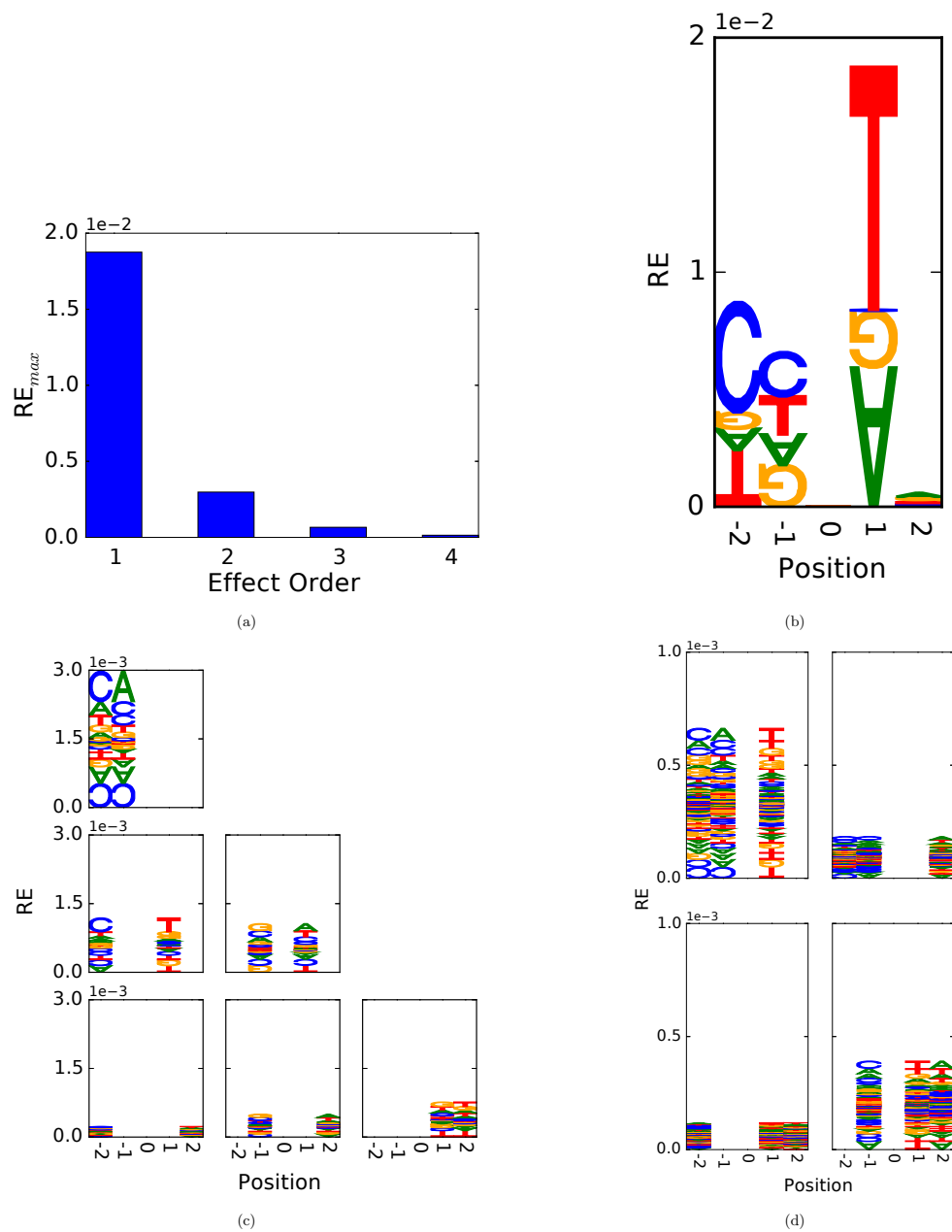


Figure 3. Flanking influences on A→G mutation in autosomal intergenic sequences. (a) First order effects are the dominant neighbourhood influence, (b) Single position effects, (c) Two-way effects, and (d) Three-way effects. For b-d, the y-axis is RE and the x-axis is the position index relative to the mutated base.

Position(s)	Deviance	df	p-value
-2	26528.3	3	0.0
-1	20038.7	3	0.0
+1	57037.8	3	0.0
+2	1802.0	3	0.0
(-2, -1)	9058.8	9	0.0
(-2, +1)	3615.8	9	0.0
(-2, +2)	701.1	9	0.0
(-1, +1)	3233.2	9	0.0
(-1, +2)	1516.8	9	0.0
(+1, +2)	2329.1	9	0.0
(-2, -1, +1)	2018.3	27	0.0
(-2, -1, +2)	561.1	27	0.0
(-2, +1, +2)	362.2	27	2.4×10^{-60}
(-1, +1, +2)	1191.2	27	0.0
(-2, -1, +1, +2)	426.5	81	2.3×10^{-48}

Table 3. Log-linear analysis of A→G autosomal intergenic mutations. Position(s) are relative to the index position (see Figure 1). Deviance is from the log-linear model, with df degrees-of-freedom and corresponding p-value obtained from the χ^2 distribution. p-values listed as 0.0 are below the limit of detection.

Some germline point mutations exhibited different neighbouring effects between sequence classes

The operation of transcription coupled DNA repair processes suggested a possible difference in neighbour effect may exist between transcribed and untranscribed sequences. This predicts a difference in mutation profile between intergenic and intronic sequences. Our analysis of neighbourhood contributions to mutation established that for first order effects, every point mutation was significantly different between the sequence classes (Table S2). For second order effects, only the transition mutations showed significant differences. The biggest difference between the regions was for A→T*. While these effects were highly significant, their $RE_{max}(1)$ were ≈ 100 fold lower than the overall influence of neighbourhood on intergenic A→T.

Neighbouring effects differ between chromosome classes

Differences in germline biology between males and females predict distinct mutation profiles between sequences located on the autosomes and X-chromosome (Huttley et al., 2000). Our test of the hypothesis of no difference in flanking base effect between autosome and X-chromosome mutations in intergenic sequences was rejected for first order influences on several of the point mutations, after correcting for 15 tests using the Holm-Šidák procedure (Holm, 1979) (Table S3). Interestingly, A→G* and C→T* showed comparable differences in flanking base effect between the chromosome classes (Deviances ≈ 26.0 and ≈ 25.4 respectively). In all cases, the effect exists at the same position as that identified as $RE_{max}(1)$ in the intergenic analysis (Table 2). While the transition mutations were the most statistically significant, their RE lay within the range of the other point mutations (Table S3) indicating their significance reflects greater abundance and thus a greater rate.

Analysis of germline mutation spectra

Our log-linear model for analysis of mutation spectra compares counts of point mutations from the same starting base between groups. By considering only mutations from a single base between different locations, differences in the abundance of the starting base between groups are controlled for. This approach can be applied to groups representing different strands, different genomic regions or different biological materials (e.g. germline and somatic).

Our analysis of germline mutation spectra indicated point mutations were uniformly strand symmetric but different between sequence categories. No sequence category exhibited strand asymmetry in mutation spectra for autosomal data. Significant differences in autosomal mutation spectra were evident between

intergenic and intronic regions. The major differences were for transversion mutations, specifically C→A and its strand complement (Table S4).

Significant differences between chromosome classes were evident (Fig 4 and Table S5). For the intergenic sequence class, T→C* transition mutations were in strong excess on autosomes compared with X-chromosome (Fig 4). Comparable results were evident for intronic sequences (Table S6).

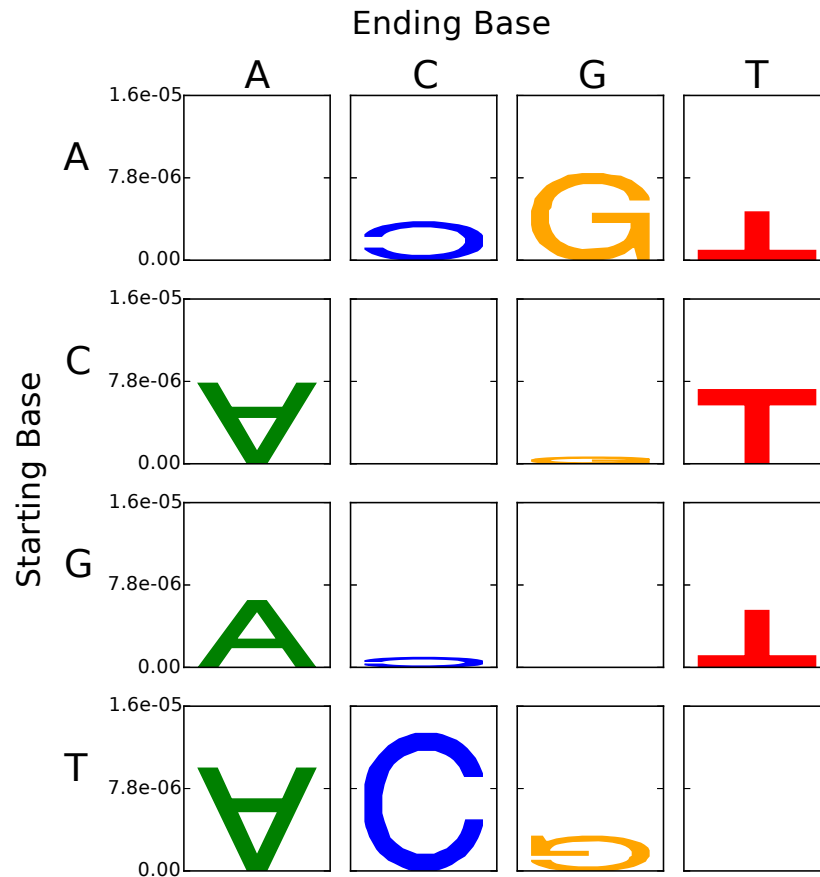


Figure 4. Significant differences in mutation spectra between autosomal and X-chromosomal intergenic sequence regions. Starting base, Ending Base correspond to X, Y respectively in X→Y. The y-axis is RE from the spectra hypothesis test and letter heights are as for the mutation motif logo. Letters in the normal orientation indicate an excess of that mutation direction in autosomal relative to the X-chromosomal mutations. Inverted letters indicate a deficit in autosomal relative to the X-chromosomal mutations.

Melanoma mutations exhibit strikingly different neighbour effects and spectra

Mutation processes in malignant melanoma are known to be distinctive and to include strand asymmetric mutation processes within genes (Plesance et al., 2010). Our analysis confirm that the profile of point mutations in the malignant melanoma sample was strikingly different to the germline mutations (Tables S10, S11). The grid of all point mutations (Fig 6) demonstrates that neighbouring influences were most pronounced for C→T point mutations and much stronger influence of neighbouring bases on transversion mutations. The neighbour effects were also significantly strand asymmetric (Table S7). Of particular note wristically distinctive for melanoma. Only substitutions affecting C were significantly different in spectra between strands with the C→T direction being over abundant on the + strand (Fig 5, Table S8).

DISCUSSION

While it has long been appreciated that sequence neighbourhoods affect point mutations, statistical methods for disentangling how neighbourhood contributes have been limited. Here we addressed this

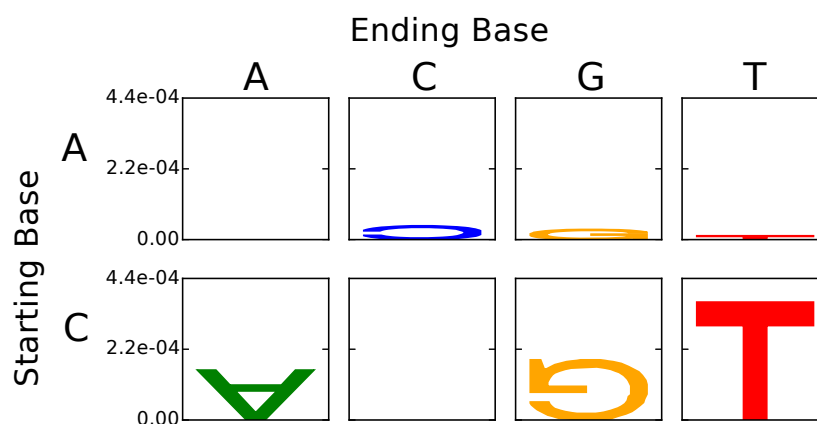


Figure 5. Strand asymmetry in malignant melanoma. Only mutations from C were statistically significant. Starting base, Ending Base correspond to X, Y respectively in $X \rightarrow Y$. The y-axis is RE from the spectra hypothesis test and letter heights are as for the mutation motif logo. Letters in the normal orientation indicate an excess of that mutation direction on the + strand. Inverted letters indicate a deficit on the + strand.

Direction	$RE_{max}(1)$	Pos.(1)	$RE_{max}(2)$	Pos.(2)	$RE_{max}(3)$	Pos.(3)
A \rightarrow C	0.0167	-1	0.0101	(-1, +1)	0.0078	(-2, +1, +2)
A \rightarrow G	0.0135	-1	0.0118	(-1, +1)	0.0051	(-1, +1, +2)
A \rightarrow T	0.0110	-1	0.0039	(-2, +1)	0.0033	(-2, -1, +1)
C \rightarrow A	0.0319	-1	0.0102	(-1, +1)	-	-
C \rightarrow G	0.0264	+1	0.0035	(-1, +1)	0.0041	(-2, -1, +1)
C \rightarrow T	0.0788	-1	0.0130	(-1, +1)	0.0006	(-2, -1, +1)
G \rightarrow A	0.0918	+1	0.0090	(-1, +1)	0.0009	(-1, +1, +2)
G \rightarrow C	0.0254	-1	0.0028	(-2, +1)	0.0043	(-1, +1, +2)
G \rightarrow T	0.0242	+1	0.0078	(+1, +2)	0.0052	(-1, +1, +2)
T \rightarrow A	0.0123	+1	0.0042	(+1, +2)	0.0044	(-1, +1, +2)
T \rightarrow C	0.0135	+1	0.0244	(-1, +1)	0.0057	(-1, +1, +2)
T \rightarrow G	0.0137	+1	0.0118	(-1, +1)	0.0074	(-2, +1, +2)

Table 4. Summary of neighbourhood contributions to mutations in malignant melanoma. $RE_{max}(\#)$ is the maximum RE for order # and Pos.(#) the corresponding position(s). All point mutations had at least one significant test after correcting for 15 tests (see Table 1) using the Holm-Šidák procedure. Non-significant results are indicated by ‘-’.

using a novel determination of the reference distribution and log-linear models. This methodological combination is robust to complexity in the genomic background of nucleotide composition. It further enables hierarchical hypothesis testing for establishing the significance and relative importance of neighbourhood effects. We illustrated utility of the models by applying them to analyses of mutations from samples reported to exhibit distinctive properties. Our analyses recapitulated well-known effects in terms of neighbour dependence and in terms of differences between genomic regions and somatic and germline, supporting the accuracy of the methods. The results revealed previously unreported neighbourhood effects that extends beyond immediate flanking positions. Analyses of mutation spectra complemented the neighbourhood analyses, confirming known features of point mutations in malignant melanoma and identifying novel differences in germline point mutation abundance between sex-chromosomes and autosomes.

The hypermutability of C \rightarrow T in CpG dinucleotides is the exemplar of context dependent mutation and a gold standard that a method of analysis should correctly recover. We established that the conventional sequence logo analysis approach did not recapitulate the dominant influence of a 3'-G (Fig 7). As this

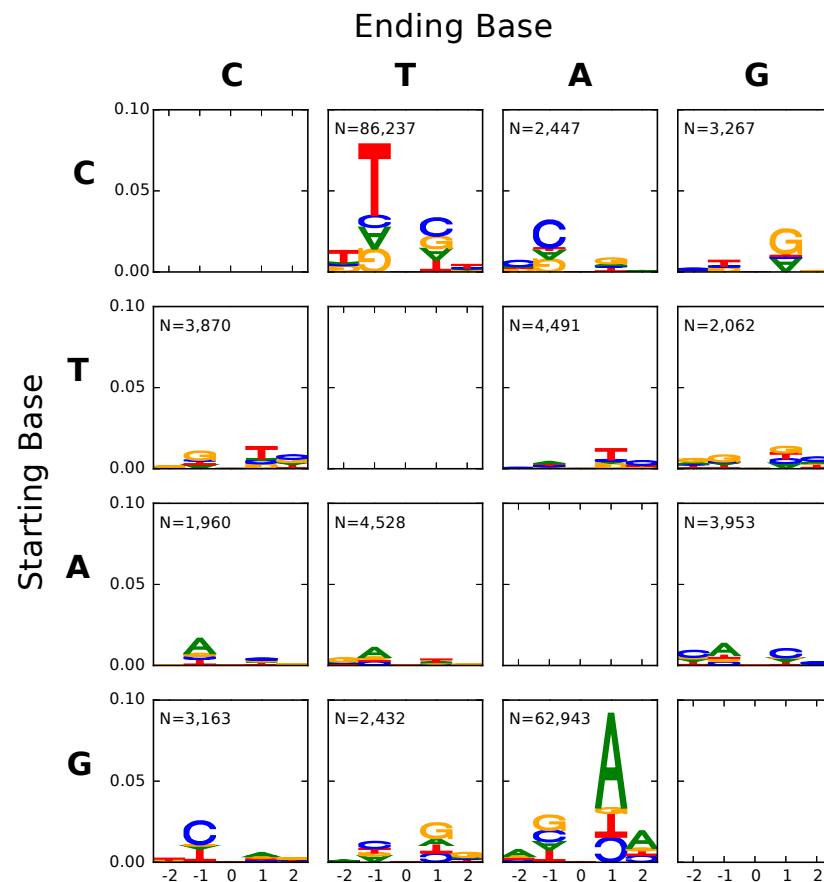


Figure 6. Panel of first order effects from all 12 point mutations from the malignant melanoma sample. Starting base, Ending Base correspond to X, Y respectively in $X \rightarrow Y$. The y-axis is RE and the x-axis is the position index relative to the mutated base. N refers to the number of SNPs from which the logo was derived.

method shares the assumption of equiprobable bases with that of Krawczak et al. (1998), the failure suggests the Euclidean distance approach will also be flawed. In contrast, as shown in Table 1 and Fig 2, our analysis successfully recapitulated this known effect. The RE_{max} values (Table 2 b) further affirm $C \rightarrow T$ as most strongly affected by neighbouring bases.

In order to sensibly interpret the results of our analyses we de-emphasise the importance of statistical significance and focus instead on effect magnitude. Due to the very large number of inferred mutations, our analyses possess very high power to detect small effects. This is illustrated by the very small p-values associated with, for example, third order effects for the $C \rightarrow T$ mutation (Table 1). Yet, the magnitude of these effects is relatively small in comparison with the first order effects (Fig 2 a). Consequently, and in addition to considering whether effects are statistically significant according to standard criteria, we contrast RE statistics to establish relative importance.

Our analysis identified numerous novel properties of neighbouring sequence influence on point mutation in the germline. First, all mutations were significantly affected by neighbouring bases with transition mutations showing a larger influence of neighbours than transversions. Interestingly, as illustrated by the $A \rightarrow G^*$ mutations, these influences did not decay monotonically with distance from the mutation (Fig 3 b). This point mutation further illustrated that multiple neighbouring positions can influence mutation outcome. Comparing RE values to that for $C \rightarrow T$ indicates that the first order neighbourhood effects of other point mutations were $\sim 5-20$ fold less, with those values corresponding to $A \rightarrow G$ and $A \rightarrow C$ mutations respectively (Table 2). Second, all mutations were significantly affected by higher order effects (interactions between adjacent bases). These were evident in a manner such that

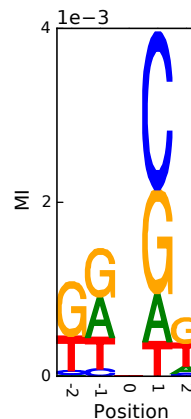


Figure 7. The CpG effect on C→T is not revealed by applying the conventional sequence logo method to autosomal intergenic mutations. MI is mutual information.

bases contiguous with each other and the mutated location showed the largest RE. This may reflect the importance of interactions amongst adjacent bases (base-stacking) in affecting DNA stability (Karlin and Burge, 1995; Yakovchuk et al., 2006). For all point mutations, the RE terms from first order effects were markedly stronger than those for higher order effects. These results were replicated in our analysis of intronic SNPs (Table S9).

The evidence for neighbouring influence on mutation raised the important question of how far these effects of flanking sequence extend? While there was strong statistical significance of positions as far as 10bp from the mutating base (Table S1), considering the relative magnitude of RE values indicated a very rapid decay away from the mutated position. In particular, that the magnitude of the effect decayed below an order of magnitude within 2 bases for transition mutations. This trend is illustrated by the mutation motif logo displays (Fig S4). While transversion mutations exhibited a slower decay in effect magnitude, and hence a larger neighbourhood, these reflect the smaller $RE_{max}(1)$ of transversions which constitute a less stringent cut-off.

Our results regarding the importance of higher order interactions indicate that considering 3-mers accounts for the majority of deviance, contradicting previous conclusions (Aggarwala and Voight, 2016). The deviances from the first order effects of A→G* and C→T* transition mutations accounted for 81% and 98% of the total deviance respectively in the autosomal intergenic sample. Inclusion of second order effects increased both these to > 96% (Tables 1, 3). Across all point mutations in the autosomal intergenic sample, combining first and second order effects accounted for a median 91% of the total deviance of the 5-mer model. These differences are further illustrated by the motif [C/T]CAAT[C/G/T]N, identified by Aggarwala and Voight (2016) as exhibiting an odds ratio of ~ 6000 for enrichment in mutated sequences. Our results (Table 3 and Fig 3d) also identified this motif as highly significant, with a p-value below the limit of detection. However, this is a third-order interaction and the RE for this specific combination of sites is 28 fold less than the strongest first order effect and accounts for only 1.5% the total deviance. We speculate the estimates of Aggarwala and Voight (2016) reflect the specificity of 7-mers for genomic regions with distinctive mutation rates and not the specific influence of neighbours on mutation.

The profile of somatic mutations is expected to exhibit differences to germline mutations due to requisite defects in DNA repair systems. As illustrated by Nik-Zainal et al. (2012), such defects are characteristic of cancers. Of the characterised cancers, malignant melanoma exhibit the most distinctive mutation signatures. Included in the distinctiveness of malignant melanoma is a striking strand asymmetry (Pleasant et al., 2010). This putatively derives from UV light induced formation of pyrimidine dimers. In transcribed regions, nucleotide excision repair processes coupled to transcription-coupled repair mechanism, results in efficient repair of transcribed strand lesions. As a consequence, mutations are expected to accumulate on the non-transcribed strand. Evidence supporting this, with more C→T mutations on non-transcribed strand than on the transcribed strand, has been reported (Pleasant et al., 2010).

Our analysis demonstrated that point mutations in melanoma were dependent on neighbours in a

manner strikingly different from that of germline processes discussed thus far (Fig 6 and Table 4). While C→T mutations were again the point mutation most affected by neighbouring bases, the motif was markedly different to that from the germline process with a 5'-T showing the greatest influence. This difference indicates that 5mC deamination plays a less prominent role in C→T. Since melanoma arises in part due to defect(s) in DNA repair the distinctive mutation motifs in melanoma indicate either a very effective masking of sequence neighbourhood effects on lesion formation, or that the DNA repair mechanisms inactivated in melanoma are strongly affected by sequence neighbourhood. Our melanoma analysis also strongly supported strand asymmetry of mutations, with the effect most pronounced for C→T.

A major asset to the log-linear modelling framework is the ease of extension to enable comparisons between samples. The utility of this is illustrated above in comparing somatic to germline processes. The appeal of this capability, however, is much broader as it further allows evaluation of the processes that contribute to within genome heterogeneity in sequence composition. We have illustrated this application here by considering genomic regions for which the incidence of mutation processes are known to differ (X-chromosome versus autosomes) or where DNA repair processes are known to differ (transcribed versus untranscribed regions).

The notion that there is a systematic tendency for mutations to originate in males has been known since Haldane (Haldane, 1935, 1946, 1948). The most popular hypothesis to account for male biased evolution is the mutation-through-DNA-replication hypothesis (Li et al., 2002; Webster et al., 2005). Other, non-replication based, differences in mutation between the sexes have also been proposed (Huttley et al., 2000). Included in these is evidence for elevated methylation of DNA in the male germline. This suggests the relative contribution of 5mC derived lesions will be greater on the autosomes compared to X-chromosome as the latter spends less time (on average) in males. Our analyses for differences in neighbour influences did lend support to existence of distinct 5mC affecting mutation processes operating between the X-chromosome and autosomes (Table S3), including a reduced magnitude of the +1 influence on the X-chromosome. However, this was not the strongest difference in neighbourhood effect between the chromosomal classes; A→G showed the strongest statistical significance while C→G showed the greatest RE. The spectra analyses further emphasised the importance of differences in A→G* point mutations (Fig 4). These results therefore indicate more extensive point mutation differences between these chromosome classes than previously appreciated and suggest a corresponding diversity in mutational processes between male and female germlines.

That differences in operation of DNA repair processes may affect mutation is predicted by the localised influence of transcription coupled DNA repair. This process is known to operate in a manner that is strand asymmetric. Differences in base parity – the frequency of A should equal that of T, G should equal C – support an effect of transcription on point mutation mutation (Touchon et al., 2003). Significant differences in neighbour effects for all point mutations were evident between intergenic and intron regions. However, our analysis of strand symmetry for neighbour effects was not significant for intron sequences for any point mutation. This suggests a distinctive mutation profile arising from transcription, rather than the influence of transcription coupled DNA repair.

As formulated, the neighbour analysis do not evaluate the relative abundance of mutations between samples. For this purpose, we introduce what we termed the mutation spectrum analysis. As the opportunity for mutation is affected by the frequency of the starting base, and base frequency differs between genomic locations, we perform spectrum analysis for each nucleotide separately. The null hypothesis is a very simple one, that the 3 possible point mutations from a starting base occur in equal frequency between samples. As such, this spectrum approach does not consider neighbouring base contributions at all and is therefore complementary to it.

For each of the above analyses comparing groups we also undertook mutation spectrum analyses. There were no significant strand differences for autosomal data. Comparisons between the X-chromosome and autosomes revealed highly significant differences in composition for all bases (Fig 4). The most pronounced difference was an excess of A→G* transition mutations on autosomes. Similarly, all point mutations showed significantly different mutation spectra between intergenic and intronic regions (Table S4). In this case, however, the dominant differences were an excess of transversions creating A/T base pairs in intergenic regions while introns were characterised by an excess of C/G base pair creating mutations.

CONCLUSION

The methods we present enable characterising mutational processes affecting samples. For the neighbour analyses, the critical properties of the methods we present derive from the specification of the reference distribution and utilisation of the well established log-linear modelling framework. This combination has considerable potential for detailed interrogations of mutation properties and should improve our understanding the mechanism of mutations, both germline and somatic. Our application of the method generated mutation motifs consistent with well known effects. We further revealed a pronounced influence of flanking bases on all point mutation processes. From germline mutations we have identified a striking dependence of the A→G transition on multiple positions. The mechanistic basis of this mutation motif is unknown.

The neighbourhood and spectral analyses examine complementary aspects of mutational process. The former examines the contribution of neighbouring bases to the mutation outcome from a starting base and the latter considers the breakdown of mutations from a single base. While the p-values from the hypothesis tests are sensitive to sample size, a property that may be proportional to mutation rate, neither approach explicitly considers the rate of mutation.

As with all methods that seek to characterise data arising from unobserved processes, there are challenges of interpretation. In both the neighbour and spectral analysis approaches, the data are a composite of mutation events with potentially diverse etiological histories. As a consequence, differences between samples will potentially reflect multiple mechanistic differences. Regardless of these issues, analyses that use measures of genetic distance, such as phylogenetics, cannot rationally rely on models of sequence divergence that assume mutations affect nucleotides independent of their neighbours. Instead, models that accommodate neighbour effects to at least ± 2 positions will need to be developed in order to reasonably capture the neighbourhood influences described here.

COMPETING INTERESTS

The authors declare that they have no competing interests.

AUTHOR'S CONTRIBUTIONS

YZ, GAH conceived the project and designed the experiments. YZ, GAH, VBY and TN specified the statistical analyses. YZ and GAH performed the research. All authors wrote the manuscript.

SUPPORTING INFORMATION

S1 File. Supplementary figures and tables.

ACKNOWLEDGEMENTS

We thank Jeremy Widman for allowing us to use his Python implementation of logo drawing code for visualisation. We thank Ben Kaehler and Stephen Haslett for their comments on versions of this work.

REFERENCES

- Aggarwala, V. and Voight, B. F. (2016). An expanded sequence context model broadly explains variability in polymorphism levels across the human genome. *Nature Genetics*, 48(4):349–355.
- Alexandrov, L. B., Nik-Zainal, S., Wedge, D. C., Aparicio, S. A., Behjati, S., Biankin, A. V., Bignell, G. R., Bolli, N., Borg, A., Børresen-Dale, A.-L., et al. (2013a). Signatures of mutational processes in human cancer. *Nature*.
- Alexandrov, L. B., Nik-Zainal, S., Wedge, D. C., Campbell, P. J., and Stratton, M. R. (2013b). Deciphering signatures of mutational processes operative in human cancer. *Cell reports*, 3(1):246–259.
- Bernardi, G. (2000). Isochores and the evolutionary genomics of vertebrates. *Gene*, 241(1):3–17.
- Brown, T. (2002). *Genomes*. Wiley-Liss.
- Chor, B., Horn, D., Goldman, N., Levy, Y., and Massingham, T. (2009). Genomic DNA k-mer spectra: models and modalities. *Genome Biol*, 10(10):R108.
- Cooke, M. S., Evans, M. D., Dizdaroglu, M., and Lunec, J. (2003). Oxidative dna damage: mechanisms, mutation, and disease. *FASEB J*, 17(10):1195–214.

- Cooper, D. N. (1995). The nature and mechanisms of human gene mutation. *The metabolic and molecular bases of inherited disease*, pages 259–291.
- Cooper, D. N. and Youssoufian, H. (1988). The CpG dinucleotide and human genetic disease. *Human genetics*, 78(2):151–155.
- Coulondre, C., Miller, J. H., Farabaugh, P. J., and Gilbert, W. (1978). Molecular basis of base substitution hotspots in *Escherichia coli*. *Nature*, 274(5673):775–780.
- Flicek, P., Amode, M. R., Barrell, D., Beal, K., Billis, K., Brent, S., Carvalho-Silva, D., Clapham, P., Coates, G., Fitzgerald, S., et al. (2013). Ensembl 2014. *Nucleic Acids Res*, page gkt1196.
- Forbes, S. A., Beare, D., Gunasekaran, P., Leung, K., Bindal, N., Boutselakis, H., Ding, M., Bamford, S., Cole, C., Ward, S., Kok, C. Y., Jia, M., De, T., Teague, J. W., Stratton, M. R., McDermott, U., and Campbell, P. J. (2015). Cosmic: exploring the world's knowledge of somatic mutations in human cancer. *Nucleic Acids Res*, 43(D1):D805–D811.
- Francioli, L. C., Polak, P. P., Koren, A., Menelaou, A., Chun, S., Renkens, I., van Duijn, C. M., Swertz, M., Wijmenga, C., van Ommen, G., Slagboom, P. E., Boomsma, D. I., Ye, K., Guryev, V., Arndt, P. F., Kloosterman, W. P., de Bakker, P. I. W., and Sunyaev, S. R. (2015). Genome-wide patterns and properties of de novo mutations in humans. *Nat Genet*, 47(7):822–6.
- Haldane, J. (1946). The mutation rate of the gene for haemophilia, and its segregation ratios in males and females. *Annals of eugenics*, 13(1):262–271.
- Haldane, J. (1948). Croonian lecture: The formal genetics of man. *Proceedings of the Royal Society of London B: Biological Sciences*, 135(879):147–170.
- Haldane, J. B. (1935). The rate of spontaneous mutation of a human gene. *Journal of Genetics*, 31(3):317–326.
- Harris, K. (2015). Evidence for recent, population-specific evolution of the human mutation rate. *Proceedings of the National Academy of Sciences*, 112(11):3439–3444.
- Helleday, T., Eshtad, S., and Nik-Zainal, S. (2014). Mechanisms underlying mutational signatures in human cancers. *Nat Rev Genet*, 15(9):585–598.
- Hodgkinson, A. and Eyre-Walker, A. (2011). Variation in the mutation rate across mammalian genomes. *Nat Rev Genet*, 12(11):756–766.
- Holm, S. (1979). A simple sequentially rejective multiple test procedure. *Scandinavian Journal of Statistics*, 6(2):65–70.
- Huttley, G. A. (2004). Modeling the impact of DNA methylation on the evolution of BRCA1 in mammals. *Mol Biol Evol*, 21(9):1760–1768.
- Huttley, G. A., Jakobsen, I. B., Wilson, S. R., and Easteal, S. (2000). How important is dna replication for mutagenesis? *Molecular biology and evolution*, 17(6):929–937.
- Ihaka, R. and Gentleman, R. (1996). R: a language for data analysis and graphics. *J. Comput. and Graph. Statistics*, 5:299–314.
- Karlin, S. (1998). Global dinucleotide signatures and analysis of genomic heterogeneity. *Curr Opin Microbiol*, 1(5):598–610.
- Karlin, S. and Burge, C. (1995). Dinucleotide relative abundance extremes: a genomic signature. *Trends Genet*, 11(7):283–90.
- Karlin, S., Campbell, A. M., and Mrázek, J. (1998). Comparative DNA analysis across diverse genomes. *Annual review of genetics*, 32(1):185–225.
- Knight, R., Maxwell, P., Birmingham, A., Carnes, J., Caporaso, J. G., Easton, B. C., Eaton, M., Hamady, M., Lindsay, H., Liu, Z., et al. (2007). Pycogent: a toolkit for making sense from sequence. *Genome Biol*, 8(8):R171.
- Krawczak, M., Ball, E. V., and Cooper, D. N. (1998). Neighboring-nucleotide effects on the rates of germ-line single-base-pair substitution in human genes. *The American Journal of Human Genetics*, 63(2):474–488.
- Li, W.-H., Yi, S., and Makova, K. (2002). Male-driven evolution. *Current opinion in genetics & development*, 12(6):650–656.
- Morton, B. R., Oberholzer, V. M., and Clegg, M. T. (1997). The influence of specific neighboring bases on substitution bias in noncoding regions of the plant chloroplast genome. *J Mol Evol*, 45(3):227–31.
- Nevarez, P. A., DeBoever, C. M., Freeland, B. J., Quitt, M. A., and Bush, E. C. (2010). Context dependent substitution biases vary within the human genome. *BMC bioinformatics*, 11(1):462.
- Nik-Zainal, S., Alexandrov, L. B., Wedge, D. C., Van Loo, P., Greenman, C. D., Raine, K., Jones, D.,

- Hinton, J., Marshall, J., Stebbings, L. A., et al. (2012). Mutational processes molding the genomes of 21 breast cancers. *Cell*, 149(5):979–993.
- Nishino, H., Buettner, V. L., Haavik, J., Schaid, D. J., and Sommer, S. S. (1996). Spontaneous mutation in big blue transgenic mice: analysis of age, gender, and tissue type. *Environ Mol Mutagen*, 28(4):299–312.
- Peltomaki, P. and Vasen, H. (1997). Mutations predisposing to hereditary nonpolyposis colorectal cancer: database and results of a collaborative study. the international collaborative group on hereditary nonpolyposis colorectal cancer. *Gastroenterology*, 113(4):1146–1158.
- Pleasant, E. D., Cheetham, R. K., Stephens, P. J., McBride, D. J., Humphray, S. J., Greenman, C. D., Varela, I., Lin, M.-L., Ordóñez, G. R., Bignell, G. R., et al. (2010). A comprehensive catalogue of somatic mutations from a human cancer genome. *Nature*, 463(7278):191–196.
- Schluter, D. (2009). Evidence for ecological speciation and its alternative. *Science*, 323(5915):737–741.
- Schneider, T. D. and Stephens, R. M. (1990). Sequence logos: a new way to display consensus sequences. *Nucleic Acids Res*, 18(20):6097–6100.
- Shannon, C. E. (2001). A mathematical theory of communication. *ACM SIGMOBILE Mobile Computing and Communications Review*, 5(1):3–55.
- Shiraishi, Y., Tremmel, G., Miyano, S., and Stephens, M. (2015). A Simple Model-Based Approach to Inferring and Visualizing Cancer Mutation Signatures. *PLOS Genetics*, 11(12):e1005657.
- Touchon, M., Nicolay, S., Arneodo, A., d'Aubenton Carafa, Y., and Thermes, C. (2003). Transcription-coupled ta and gc strand asymmetries in the human genome. *FEBS Lett*, 555(3):579–82.
- Vinson, C. and Chatterjee, R. (2012). Cg methylation. *Epigenomics*, 4(6):655–663.
- Webster, M. T., Smith, N. G., Hultin-Rosenberg, L., Arndt, P. F., and Ellegren, H. (2005). Male-driven biased gene conversion governs the evolution of base composition in human alu repeats. *Molecular biology and evolution*, 22(6):1468–1474.
- Yakovchuk, P., Protozanova, E., and Frank-Kamenetskii, M. D. (2006). Base-stacking and base-pairing contributions into thermal stability of the dna double helix. *Nucleic Acids Res*, 34(2):564–74.
- Ying, H. and Huttley, G. (2011). Exploiting CpG Hypermethylability to Identify Phenotypically Significant Variation Within Human Protein-Coding Genes. *Genome Biology and Evolution*, 3:938.
- Zhang, X. and Mathews, C. K. (1995). Natural DNA precursor pool asymmetry and base sequence context as determinants of replication fidelity. *J Biol Chem*, 270(15):8401–4.
- Zhao, Z. and Boerwinkle, E. (2002). Neighboring-nucleotide effects on single nucleotide polymorphisms: a study of 2.6 million polymorphisms across the human genome. *Genome research*, 12(11):1679–1686.