

Pan-Cancer TF Footprinting & Noncoding Mutation Data Integration

Suffian Azizan

Note: This is a Github README doc that is dynamically generated from the qmd notebook.

1 Background

This pan-cancer study makes use of restricted ATAC-seq and DNase-seq datasets generated from healthy and cancerous human tissue samples available in the public databases (i.e. TCGA and BLUEPRINT) to generate TF footprinting data from the open chromatin regions. The TF footprinting data is then combined with the noncoding mutation data from the same samples obtained via variant calling to identify TF binding sites (TFBS) that carry variant alleles (or potentially mutations) that may modulate TF footprint scores (proxy for TF binding activity).

2 Data Preprocessing

TF footprints are determined and scored using the TOBIAS program published by Looso Lab (Max Plank Institute). A customized TOBIAS pipeline was run on individual samples so TOBIAS's internal intersample normalization was not applied. The raw footprint scores were collated for all samples in this pan-cancer study and combined into a large dense matrix containing unique footprint sites of a particular motif of interest across all samples. As we have 1360 motifs of interest in this study, we have 1360 large data tables to process.

For the purpose of demonstration, a subset of TF footprint data of the open chromatin regions in only breast cancer samples from TCGA database was used here to save storage space and lower processing overhead. Additionally, only 1 TF motif will be presented here in the analysis workflow.

First, load up required Python packages.

```
import os
import textwrap
import pandas as pd
import seaborn as sns
import statsmodels.api as sm
import matplotlib.pyplot as plt
from natsort import import index_natsorted
```

2.1 Loading up the footprint data table

Now, we can import the data tables as Pandas dataframes. The data are stored in tab-delimited text files (.tsv). The first column defines the chromosomal location, the second and third column contain the genomic coordinates of the TF footprint, the fourth column retains the strandedness of the TF binding site (TFBS), the fifth column contains the TFBS score (similarity score with the tested motif PWM), and the rest of the columns carry the actual TOBIAS-calculated TF footprint scores for individual samples.

The data is loaded into a Pandas dataframe and the first 5 rows are displayed.

```
# import the data
filepath = '../demo-data/E2F2_E2F2_HUMAN.H11M0.0.B_BRCA-subtype-vcf-filtered-matrix.txt'
df_fpscore = pd.read_csv(filepath, sep='\t')
```

Aside: specify a formatter function to wrap long text in the data tables.

```
#replace underscore with whitespace
func_underscore_replace = lambda x: x.replace("_", " ")
#wrap text
def func_wrap(x):
    if isinstance(x, str):
        return textwrap.fill(x, width=10)
    else:
        return x
```

TFBS chr	TFBS start	TFBS end	TFBS strand	TFBS score	98JKPD8 LumA score	AN- AB5F7 Basal score	S6R691V Her2 score	PU24GB8 LumB score	2GAMBDQ Nor- mal- like score
chr1	10628	10638	+	7.401890	0.000000	0.000000	0.000000	0.000000	0.000000
chr1	181224	181234	+	7.998660	0.000000	0.000000	0.000000	0.000000	0.000000
chr1	779214	779224	-	7.796470	0.000000	0.000000	0.000000	0.000000	0.000000
chr1	998754	998764	+	8.560290	0.137600	0.140350	0.128420	0.145000	0.105100

Rename the `_score` column to `_fps` to avoid confusion with the mutation score column in the mutation data table later, and drop the `TFBS_strand` and `TFBS_score` columns as they are not needed for now.

```
# drop the column "TFBS_strand" and "TFBS_score"
df_fpscore = df_fpscore.drop(columns=["TFBS_strand", "TFBS_score"])
# rename columns in the dataframe
df_fpscore = df_fpscore.rename(columns={"TFBS_chr": "Chromosome", "TFBS_start": "Start",
"TFBS_end": "End", "2GAMBDQ_Normal-like_score": "2GAMBDQ_Norm_fps"})
# for all column names that end with the string 'score', replace the string with 'fps'
df_fpscore = df_fpscore.rename(columns=lambda x: x.replace('score', 'fps') if
x.endswith('score') else x)
```

Chromo- some	Start	End	98JKPD8 LumA fps	ANAB5F7 Basal fps	S6R691V Her2 fps	PU24GB8 LumB fps	2GAMBDQ Norm fps
chr1	10628	10638	0.000000	0.000000	0.000000	0.000000	0.000000
chr1	181224	181234	0.000000	0.000000	0.000000	0.000000	0.000000
chr1	779214	779224	0.000000	0.000000	0.000000	0.000000	0.000000
chr1	998754	998764	0.137600	0.140350	0.128420	0.145000	0.105100
chr1	998768	998778	0.182240	0.167080	0.169110	0.183020	0.127590
chr1	1019693	1019703	0.622590	0.379450	0.489290	0.644200	0.492220
chr1	1041096	1041106	0.096050	0.097560	0.090640	0.127590	0.110160
chr1	1164827	1164837	0.110270	0.149950	0.073960	0.067850	0.080790
chr1	1206766	1206776	0.067360	0.063970	0.041980	0.069200	0.072820
chr1	1309989	1309999	0.123050	0.144050	0.141260	0.124540	0.096460

2.2 Loading up the mutation data table

This mutation data is generated from the output of bcftools variant calling pipeline. First, load up an example data file to see how the data is structured.

```
vcfpath = '../demo-data/2GAMBDQ_E2F2_E2F2_HUMAN.H11M0.0.B_AF-per-site-with-indels.txt'
# load up the vcf file with indels and multiallelic sites split into separate rows
df_vcf = pd.read_csv(vcfpath, sep="\t")
```

#[1]CHROM	[2]POS	[3]REF	[4]ALT	[5]AF
chr1	10629	GGCGCGC	GGCGC	1.0
chr1	998764	C	G	0.5
chr1	998764	CGGAGGG	CGGAGGGGAG GG	0.3125
chr1	1041101	CGGAGCGGGG CGGGAGCGGG GCGGGAGCGG GG	CGGAGCGGGG CGGGAGCGGG G	0.5
chr1	1164837	C	T	0.8125

The file above corresponds to just one of the sample IDs in this pan-cancer study. To load up all the mutation data for all samples, we need to load up all the files in the directory. Let's put the dataframes in one dictionary object using a loop.

```
# create a vcf load function for the query vcfs
def load_vcf(vcf_path):
    # load up the vcf file with indels and multiallelic sites split into separate rows
    df_vcf = pd.read_csv(vcf_path, sep="\t")
    # rename columns in the dataframe
    df_vcf = df_vcf.rename(columns={"#[1]CHROM": "Chromosome", "[2]POS": "Start", "[3]REF":
    "ref_allele", "[4]ALT": "alt_allele", "[5]AF": "AF"})
    # add a column next to the "start" column called "end" with the same value as the
    "start" column
    df_vcf.insert(2, "End", df_vcf["Start"])
    return df_vcf

# now put the paths in a list
paths = [
    "../demo-data/2GAMBDQ_E2F2_E2F2_HUMAN.H11M0.0.B_AF-per-site-with-indels.txt",
    "../demo-data/98JKPD8_E2F2_E2F2_HUMAN.H11M0.0.B_AF-per-site-with-indels.txt",
    "../demo-data/ANAB5F7_E2F2_E2F2_HUMAN.H11M0.0.B_AF-per-site-with-indels.txt",
    "../demo-data/PU24GB8_E2F2_E2F2_HUMAN.H11M0.0.B_AF-per-site-with-indels.txt",
    "../demo-data/S6R691V_E2F2_E2F2_HUMAN.H11M0.0.B_AF-per-site-with-indels.txt"
]

# create a list of IDs
ids = [ "ANAB5F7_basal", "98JKPD8_lumA", "PU24GB8_lumB", "S6R691V_her2", "2GAMBDQ_norm"]

# create a pair dictionary using nested dict comprehension
# for each id in the list of ids, iterate through the list of paths and check if the id is
in the path; this means there is no need to order the list of ids according to the order
of the paths

path_id_dict = {id: load_vcf(path) for id in ids for path in paths if id.split("_")[0]
in path}
```

2.3 Using PyRanges for dataframe merging

As we are dealing with genomic regions where the data is related to interval values (start and end coordinates) spanning across two columns, it is not possible to do dataframe overlap or join using Pandas. We will need to use a specialized Python package called PyRanges to handle genomic coordinates.

Import the package and then first convert the footprint dataframe into a PyRanges object.

```
import pyranges as pr

gr_fpscore = pr.PyRanges(df_fpscore)
```

Chromosome	Start	End	98JKPD8 LumA fps	ANAB5F7 Basal fps	S6R691V Her2 fps	PU24GB8 LumB fps	2GAMBDQ Norm fps
chr1	10628	10638	0.0	0.0	0.0	0.0	0.0
chr1	181224	181234	0.0	0.0	0.0	0.0	0.0
chr1	779214	779224	0.0	0.0	0.0	0.0	0.0
chr1	998754	998764	0.1376	0.14035	0.12842	0.145	0.1051
chr1	998768	998778	0.18224	0.16708	0.16911	0.18302	0.12759

Do the same for the mutation dataframes in the dictionary. Loop through it and save them in a new dictionary.

```
# load up vcf_dfs into pyranges
grs = {}
for name,vcf in path_id_dict.items():
    gr_vcf = pr.PyRanges(vcf)
    grs[name] = gr_vcf

print(grs["ANAB5F7_basal"].head(n=10))
```

```
+-----+-----+-----+-----+
| Chromosome | Start   | End     | +3     |
| (category) | (int64) | (int64) | ...    |
+-----+-----+-----+-----+
| chr1       | 10629   | 10629   | ...    |
| chr1       | 181234  | 181234  | ...    |
| chr1       | 779216  | 779216  | ...    |
| chr1       | 998764  | 998764  | ...    |
| ...        | ...     | ...     | ...    |
| chr1       | 1019700 | 1019700 | ...    |
| chr1       | 1041101 | 1041101 | ...    |
| chr1       | 1164837 | 1164837 | ...    |
| chr1       | 1206770 | 1206770 | ...    |
+-----+-----+-----+-----+
```

Unstranded PyRanges object has 10 rows and 6 columns from 1 chromosomes.
For printing, the PyRanges was sorted on Chromosome.
3 hidden columns: ref_allele, alt_allele, AF

2.4 Merging the PyRanges objects

Now, we can merge the PyRanges objects using the join function. The how argument is set to left to retain all the rows in the left dataframe (i.e. the footprint dataframe) and the suffix argument is set to `_[sample ID]_varsite_pos` to add a suffix to the columns in the right dataframe (i.e. the mutation dataframe) to avoid column name clashes.

```

count = 0
for key, val in grs.items():

    if count == 0:
        overlap = gr_fpscore.join(val, how='left', suffix=f"_{key}_varsite_pos",
preserve_order = True)
    else:
        overlap = filtered_gr.join(val, how='left', suffix=f"_{key}_varsite_pos",
preserve_order = True)

    # drop the column "End" column
    overlap = overlap.drop([f"End_{key}_varsite_pos"])

    # cluster the pyRanges object by genomic range; overlapping regions wil share the same
    id. This will add a new column called "Cluster"
    overlap = overlap.cluster(slack=-1)

    # cast back into a dataframe and filter by the AF column's max value (by Cluster);
    this returns a filtered dataframe
    filtered_df = overlap.df.loc[overlap.df.groupby('Cluster')['AF'].idxmax()]

    # cast back into a dataframe and rename metadata columns
    filtered_df = filtered_df.rename(columns={f"Start_{key}_varsite_pos": f"{key}_
varsite_pos", "ref_allele": f"{key}_REF", "alt_allele": f"{key}_ALT", "AF": f"{key}_
_AF"})

    # replace all the -1 values in column 'Start_varsites', 'ref_allele' and 'alt_allele',
    and AF with 0
    # Define a dictionary mapping column names to values to replace
    replace_dict = {f"{key}_varsite_pos": {-1: None}, f"{key}_REF": {str(-1): None},
f"{key}_ALT": {str(-1): None}, f"{key}_AF": {-1: 0}}
    filtered_df = filtered_df.replace(replace_dict)

    # drop cluster column
    filtered_df = filtered_df.drop(columns=["Cluster"])

    # cast back into pyrange object
    filtered_gr = pr.PyRanges(filtered_df)

    # increment count
    count += 1

```

Note that during the overlap process, the use of cluster function is to ensure that overlapping regions will share the same ID. This is important as we will need to filter the overlapping regions by the maximum allele frequency (AF) value so that only 1) unique chromosome regions are returned, and 2) regions with multiallelic sites, only the site with the highest AF value is returned.

Now, we can clean up the PyRanges object and convert it back to a Pandas dataframe. This file will be the basis for all downstream analyses.

```

final_df = filtered_gr.df

# create a column called 'region_id'
final_df["region_id"] = final_df["Chromosome"].astype(str) + ":" +
final_df["Start"].astype(str) + "-" + final_df["End"].astype(str)

# for all column name ending with the string '_fps', split the string, take the second

```

```
final_df = final_df.rename(columns=lambda x: x.split('_')[0] + '_' + x.split('_')[1]
[0].lower() + x.split('_')[1][1:] + '_fps' if x.endswith('_fps') else x)
```

[illegible]

3 Merged Dataframe Data Analysis

First load up the merged dataframe.

The motif ID of the current TF data: E2F2 E2F2 HUMAN.H11M0.0.B

Chromosome	Start	End	8		9		10		11		12		13		14		15		16		17		18		19		20		21		22		23		24		25		26		27		28		29		30		31		32		33		34		35		36		37		38		39		40		41		42		43		44		45		46		47		48		49		50		51		52		53		54		55		56		57		58		59		60		61		62		63		64		65		66		67		68		69		70		71		72		73		74		75		76		77		78		79		80		81		82		83		84		85		86		87		88		89		90		91		92		93		94		95		96		97		98		99		100		101		102		103		104		105		106		107		108		109		110		111		112		113		114		115		116		117		118		119		120		121		122		123		124		125		126		127		128		129		130		131		132		133		134		135		136		137		138		139		140		141		142		143		144		145		146		147		148		149		150		151		152		153		154		155		156		157		158		159		160		161		162		163		164		165		166		167		168		169		170		171		172		173		174		175		176		177		178		179		180		181		182		183		184		185		186		187		188		189		190		191		192		193		194		195		196		197		198		199		200		201		202		203		204		205		206		207		208		209		210		211		212		213		214		215		216		217		218		219		220		221		222		223		224		225		226		227		228		229		230		231		232		233		234		235		236		237		238		239		240		241		242		243		244		245		246		247		248		249		250		251		252		253		254		255		256		257		258		259		260		261		262		263		264		265		266		267		268		269		270		271		272		273		274		275		276		277		278		279		280		281		282		283		284		285		286		287		288		289		290		291		292		293		294		295		296		297		298		299		300		301		302		303		304		305		306		307		308		309		310		311		312		313		314		315		316		317		318		319		320		321		322		323		324		325		326		327		328		329		330		331		332		333		334		335		336		337		338		339		340		341		342		343		344		345		346		347		348		349		350		351		352		353		354		355		356		357		358		359		360		361		362		363		364		365		366		367		368		369		370		371		372		373		374		375		376		377		378		379		380		381		382		383		384		385		386		387		388		389		390		391		392		393		394		395		396		397		398		399		400		401		402		403		404		405		406		407		408		409		410		411		412		413		414		415		416		417		418		419		420		421		422		423		424		425		426		427		428		429		430		431		432		433		434		435		436		437		438		439		440		441		442		443		444		445		446		447		448		449		450		451		452		453		454		455		456		457		458		459		460		461		462		463		464		465		466		467		468		469		470		471		472		473		474		475		476		477		478		479		480		481		482		483		484		485		486		487		488		489		490		491		492		493		494		495		496		497		498		499		500		501		502		503		504		505		506		507		508		509		510		511		512		513		514		515		516		517		518		519		520		521		522		523		524		525		526		527		528		529		530		531		532		533		534		535		536		537		538		539		540		541		542		543		544		545		546		547		548		549		550		551		552		553		554		555		556		557		558		559		560		561		562		563		564		565		566		567		568		569		570		571		572		573		574		575		576		577		578		579		580		581		582		583		584		585		586		587		588		589		590		591		592		593		594		595		596		597		598		599		600		601		602		603		604		605		606		607		608		609		610		611		612		613		614		615		616		617		618		619		620		621		622		623		624		625		626		627		628		629		630		631		632		633		634		635		636		637		638		639		640		641		642		643		644		645		646		647		648		649		650		651		652		653		654		655		656		657		658		659		660		661		662		663		664		665		666		667		668		669		670		671		672		673		674		675		676		677		678		679		680		681		682		683		684		685		686		687		688		689		690		691		692		693		694		695		696		697		698		699		700		701		702		703		704		705		706		707		708		709		710		711		712		713		714		715		716		717		718		719		720		721		722		723		724		725		726		727		728		729		730		731		732		733		734		735		736		737		738		739		740		741		742		743		744		745		746		747		748		749		750		751		752		753		754		755		756		757		758		759		760		761		762		763		764		765		766		767		768		769		770		771		772		773		774		775		776		777		778		779		780		781		782		783		784		785		786		787		788		789		790		791		792		793		794		795		796		797		798		799		800		801		802		803		804		805		806		807		808		809		810		811		812		813		814		815		816		817		818		819		820		821		822		823		824		825		826		827		828		829		830		831		832		833		834		835		836		837		838		839		840		841		842		843		844		845		846		847		848		849		850		851		852		853		854		855		856		857		858		859		860		861		862		863		864		865		866		867		868		869		870		871		872		873		874		875		876		877		878		879		880		881		882		883		884		885		886		887		888		889		890		891		892		893		894		895		896		897		898		899		900		901		902		903		904		905		906		907		908		909		910		911		912		913		914		915		916		917		918		919		920		921		922		923		924		925		926		927		928		929		930		931		932		933		934		935		936		937		938		939		940		941		942		943		944		945		946		947		948		949		950		951		952		953		954		955		956		957		958		959		960		961		962		963		964		965		966		967		968		969		970		971		972		973		974		975		976		977		978		979		980		981		982		983		984		985		986		987		988		989		990		991		992		993		994		995		996		997		998		999		1000		1001		1002		1003		1004		1005		1006		1007		1008		1009		1010		1011		1012		1013		1014		1015		1016		1017		1018		1019		1020		1021		1022		1023		1024		1025		1026		1027		1028		1029		1030		1031		1032		1033		1034		1035		1036		1037		1038		1039		1040		1041		1042		1043		1044		1045		1046		1047		1048		1049		1050		1051		1052		1053		1054		1055		1056		1057		1058		1059		1060		1061		1062		1063		1064		1065		1066		1067		1068		1069		1070		1071		1072		1073		1074		1075		1076		1077		1078		1079		1080		1081		1082		1083		1084		1085		1086		1087		1088		1089		1090		1091		1092		1093		1094		1095		1096		1097		1098		1099		1100		1101		1102		1103		1104		1105		1106		1107		1108		1109		1110		1111		1112		1113		1114		1115		1116		1117		1118		1119		1120		1121		1122		1123		1124		1125		1126		1127</	
------------	-------	-----	---	--	---	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	-----	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	------	--	--------	--

```
afps matrix = afps df.filter(regex='AF$|fps$|id$').copy()
```

This matrix is in the wide format so it should be converted into a long format for easier wrangling.

```
# sort the dataframe by region_id naturally
afps_mtx_lpv = afps_mtx_lpv.reindex(index=index_natsorted(afps_mtx_lpv['region_id']))
afps_mtx_lpv = afps_mtx_lpv.reset_index(drop=True)
```

	region id	sample id	AF	FPS
0	chr1:10628-10638	2GAMBDQ_norm	1.000000	0.000000
1	chr1:10628-10638	98JKPD8_lumA	0.958333	0.000000
2	chr1:10628-10638	ANAB5F7_basal	1.000000	0.000000
3	chr1:10628-10638	PU24GB8_lumB	1.000000	0.000000
4	chr1:10628-10638	S6R691V_her2	1.000000	0.000000

Number of rows in the wide form: 2972
Number of rows in the long form: 14860

3.1 Scaling the TF footprint scores

As the TF footprint scores are not normalized across samples, has a range from 0 to +Inf, and the fact that the mutation data come in the form of allelic frequency (AF) probabilistic values (i.e. values between 0 and 1), the TF footprint scores should be scaled between 0 to

1.

```
# use MinMaxScaler to scale the raw fps values to range between 0 and 1
from sklearn.preprocessing import MinMaxScaler
# scale the FPS values to a range of 0-1
# Initialize a MinMaxScaler
scaler = MinMaxScaler()

# copy df
fps_df_scaled = afps_matrix.filter(regex='_fps$|_id$').copy()

# set the index to 'region_id'
fps_df_scaled = fps_df_scaled.set_index('region_id')

# Fit the MinMaxScaler to the 'FPS' column and transform it
fps_df_scaled = pd.DataFrame(scaler.fit_transform(fps_df_scaled),
                             columns=fps_df_scaled.columns, index=fps_df_scaled.index)

# rename columns by adding '_scaled' to the column names
fps_df_scaled = fps_df_scaled.add_suffix('_scaled')
```


	98JKPD8 lumA fps scaled	ANAB5F7 basal fps scaled	S6R691V her2 fps scaled	PU24GB8 lumB fps scaled	2GAMBDQ norm fps scaled
region_id					
chr1:10628-10638	0.000000	0.000000	0.000000	0.000000	0.000000
chr1:181224-181234	0.000000	0.000000	0.000000	0.000000	0.000000
chr1:779214-779224	0.000000	0.000000	0.000000	0.000000	0.000000
chr1:998754-998764	0.034220	0.034358	0.040032	0.041878	0.027917
chr1:998768-998778	0.045321	0.040901	0.052716	0.052858	0.033891

Now, convert the scaled dataframe into a long format.

```
# reset index
fps_df_scaled_long = fps_df_scaled.reset_index()
# convert to long format
fps_df_scaled_long = fps_df_scaled_long.melt(id_vars=["region_id"], var_name="variable",
value_name="value")

# split the variable column into sample_id and type columns using reverse split string
method, which returns a dataframe of columns based on the number of splits (n=x); this
can directly be assigned to new columns in the original dataframe
# Split the 'variable' column into three parts
fps_df_scaled_long[['part1', 'part2', 'part3']] =
fps_df_scaled_long['variable'].str.rsplit('_', n=2, expand=True)

# Assign part1 to 'sample_id' and concatenate the other parts to form 'type'
fps_df_scaled_long['sample_id'] = fps_df_scaled_long['part1']
fps_df_scaled_long['type'] = fps_df_scaled_long['part2'].str.upper() + '_' +
fps_df_scaled_long['part3']

# Drop the unnecessary columns
fps_df_scaled_long = fps_df_scaled_long.drop(['variable', 'part1', 'part2', 'part3'],
axis=1)

# now pivot the dataframe to create new columns based on the type column
fps_df_scaled_lpv = fps_df_scaled_long.pivot(index=['region_id', 'sample_id'],
columns='type', values='value').reset_index()

# remove the index name and rename the columns to match the type values
fps_df_scaled_lpv = fps_df_scaled_lpv.rename_axis(None, axis=1)

# sort the dataframe by region_id naturally
fps_df_scaled_lpv =
fps_df_scaled_lpv.reindex(index=index_natsorted(fps_df_scaled_lpv['region_id']))
fps_df_scaled_lpv = fps_df_scaled_lpv.reset_index(drop=True)
```

	region id	sample id	FPS scaled
0	chr1:10628-10638	2GAMBDQ_norm	0.000000
1	chr1:10628-10638	98JKPD8_lumA	0.000000
2	chr1:10628-10638	ANAB5F7_basal	0.000000
3	chr1:10628-10638	PU24GB8_lumB	0.000000
4	chr1:10628-10638	S6R691V_her2	0.000000

Number of rows in the scaled matrix: 2972

Number of rows in the scaled matrix in the long form: 14860

The distribution of the unscaled and scaled FPS datasets can be plotted using Seaborn's `displot` function.

