# שחזור עץ פילוגנטי על תאי נאופלזמות מיאלופרוליפרטיביות

# שמות חברי הקבוצה:

- .318510633 גל סזנה
- .208635334 נועה מרגוליס
- .326045515 עדי יפרואימסקי
  - .325245280 יואל מרקו •
  - .214935165 איתן סמסון • •

# :רקע

- Life histories of myeloproliferative neoplasms inferred from . "phylogenies
- פרטים: המאמר פורסם בכתב העת Nature, כרך 602, בתאריך 3 בפברואר 2022.
   המחברים כוללים את ניקולס וויליאמס ואחרים ממוסדות כמו מכון וולקאם סנגר והמחלקה להמטולוגיה באוניברסיטת קיימברידג'.
- תקציר רקע: מחלות מיילופרוליפרטיביות (MPNs) הן סוג של סרטן דם כרוני הנגרם ממוטציות גנטיות סומטיות בתאי גזע המטופויאטיים .(HSCs) המאמר עוסק בזיהוי מוטציות סומטיות ותיעוד ההיסטוריה הגנומית של מושבות מתאי גזע של חולי MPN .
   זוהו 580,133 מוטציות סומטיות לצורך שחזור אילנות יוחסין המטופויאטיים ובחינת ההיסטוריות והתרחבותית של השיבוטים.
- המחקר מצביע על כך שרכישת מוטציות מובילות (driver mutation) מוקדמות והתפתחותן לכל אורך החיים הן הבסיס למחלות מיילופרוליפרטיביות במבוגרים, דבר שמעלה הזדמנויות להתערבות מוקדמת ומציע מודל חדש להתפתחות סרטן.

# שאלת מחקר ומטרות:

# שאלת מחקר: ●

- כיצד מוטציות גנטיות מוקדמות ותהליכי התרחבות שיבוטיים משפיעים על
   התפתחותן של מחלות מיילופרוליפרטיביות לאורך החיים?
  - ? מתי בערך בחיים אדם עלול לקבל מוטציית דרייבר ⊙

#### :מטרות

- נרצה לנצל את הידע הנלמד בקורס בנושא בניית עצים פילוגנטיים על מנת לזהות
   פיצולים (מוטציות) בהתפחות התאים שמובילות לסרטן להבין את התזמון של
   רכישת מוטציות מובילות(driver mutations) ואת התרחבותן.
  - שיחזור תוצאות המאמר עליו אנחנו מתבססים.
    - הכרת עבודה עם דאטה ביולוגי.

#### :דאטה

- קישור למקור הדאטה סט.
- המידע נלקח מתוך קובץ נתונים גנומיים בפורמט (VCF (Variant Call Format), שמטרתו לתאר מוטציות גנטיות באזורים ספציפיים בגנום. במקרה זה, מדובר במוטציות סומטיות שנמצאו בתאי דם שנדגמו מ-12 פרטים. אלו מוטציות שנרכשו במהלך החיים ולא תורשתיות.
  - ראו הרחבה על הדאטה ונספח ניתוח הפורמט.

# היפותזה:

מוטציות מניעות שנרכשות בשלבים מוקדמים בחיים מובילות לדינמיקות התרחבות משתנות לפני האבחנה הקלינית.

הבנה זו תוכל לשפר את הגילוי המוקדם ולתרום להתאמה אישית של טיפולים.

#### שיטות חישוביות:

- שימוש באלגוריתמים של בניית עצים פילוגנטיים (MPBoot). •
- בשיעור למדנו שתי שיטות מרכזיות למימוש האלגוריתם הנבדלות זו מזו בקריטריון

  Neighbor joining השנייה UPGMA השנייה של הקודקודים האחת
  - מודלים בייזיאניים לחישוב קצבי התפשטות קלונית.
  - ניתוח סטטיסטי להשוואת מוטציות מניעות בין קלונים.

# שלבי ביצוע:

- 1. ניתוח ועיבוד ראשוני של הנתונים ובדיקת איכות.
- 2. בניית עץ פילוגנטי עבור כל הדאטה המצורף במאמר:
- השוואת גנים של חולים ברצפים הייעודיים לגנים בריאים ותיעוד היסטוריית המוטציות.
- הוא ללא NJ . נבדוק איזו שיטה של בניית עץ פילו גנטי תניב תוצאות טובות יותר. דרישה לאולטרמטריות.
  - 3. ניתוח התוצאות מהעץ ומסקנות.

# בביבליוגרפיה:

Williams, N. et al. (2022). Life histories of myeloproliferative neoplasms inferred .168–162 ,602 ,\*from phylogenies. \*Nature

# גיוון הקבוצה:

- 3 סטודנטים לביולוגיה חישובית, סטודנט למדעי המחשב חד חוגי, סטודנטית למדעי
   המחשב עם חטיבה בביולוגיה.
- . 3 עתודאים, 1 קצין במיל' (תותחנים), 1 קצינה בקבע (חיל הים) שילוב הזרועות השונים.
  - דתיים, 3 חילוניים.
    - .ג'ינג'י אחד •
    - גם בנים וגם בנות.
  - 4 משקפופרים 1 ללא משקפיים.
  - . גיוון אתני ישראל הראשונה עם ישראל השנייה

הקבוצה שלנו מגוונת ברקע האקדמי, המקצועי והאישי של חבריה.

אנו כוללים סטודנטים ממסלולים שונים כמו ביולוגיה חישובית ומדעי המחשב, לצד עתודאים, קצין במילואים וקצינה בקבע מחילות שונים. בקבוצה יש ייצוג של דתיים וחילוניים, בנים ובנות, וכן גיוון אתני המשקף את החברה הישראלית. אלמנטים אישיים כמו ג'ינג'י אחד ומשקפופרים מוסיפים צבע וגיוון נוסף.

המגוון שלנו מאפשר שיתוף פעולה ייחודי וחשיבה רחבה ויצירתית.

### הרחבה על הדאטה

כל המוטציות שנתעסק איתם במאמר הם מוטציות סומטיות ולא מוטציות מתאי נבט, מסוג (סניפ) החלפה של בסיס בודד או מסוג (Indel) הכנסה ומחיקה של בסיס בודד או מספר בסיסים. בתהליך עושים single-colony whole-genome sequencing על תאים מהדם או ממוח עצם. כלומר, התאים שנלקחו הופרדו לתאים בודדים והם הופרדו לתרביות שונות, ושם נתנו לכל אחת להתחלק ולהתפתח לקולונה.

כל קולונה מייצגת לנו גנום של הבן אדם בשלב כלשהו בחיים שלו, עם מידע כלשהו של מוטציות שהוא כן סבר ומוטציות שהוא לא סבר. במחקר נבקש להבין איזה מסלול של צבירת מוטציות מוביל לסרטן ומתי בערך אמורה להיות המוטציית דרייבר.

יש בידנו מוטציות שמתאפיינות עם סרטן ולכן נוכל להגיד שקולונה עם מסלול מוטציות שמוביל למוטציה עם סרטן הוא מסלול מוטציות שמוביל לסרטן. בשביל שני הכיוונים נרצה ליצור עץ פילוגנטי שמכיל את המסלולים של מוטציות מסוימות בגוף. אז מה צריך בשביל להפיק את העץ הפילוגנטי?

ראשית, צריך את המידע מהמחקר על איזה מוטציות כל קולונה צברה ואיזה היא לא צברה. יש לנו כבר מקומות שחשודים למוטציה שאנחנו יודעים עליהם. מכל קולונה לקחו כמה תאים ועשו להם WGS וקיבלו כמה גנומים מאותה קולונה ובדקו על הגנומים האלו במקומות החשודים במוטציה האם מופיעה אותה מוטציה או לא. ספרו את מספר הפעמים שהופיעה וכתבו "1" אם החליטו שיש, כתבו "0" אם החליטו שאין וכתבו "?" אם לא יכלו להסיק מסקנה חד משמעית. כל המידע הזה מאוגד לתוך קובץ מסוג VCF (הסבר על הבנת הקובץ בנספח א'), על הקובץ הזה עושים אנליזות בשביל להוציא את העץ הפילוגנטי.

#### נספח א' – ניתוח פורמט הדאטה

הקובץ מתחיל בתיאור של המידע שמוכל בתוכו, הקטע של התיאור של הצגת הנתונים מתחיל ב##. אחריו מגיע טמפלייט של מבנה הצגת שורה. כל שורה היא בדיקה של וריאנט סומטי ספציפי. אחריהם יש את השורות של הבדיקות עצמם (הנתונים).

```
של
                                                                                                                                                                            1. הצהרת
                             הקובץ
                ##INFO=<ID=SOMATIC,Number=0,Type=Flag,Description="Variant is Somatic">
               ##INFO=<ID=MUT_TYPE,Number=1,Type=String,Description="Mutation Type">
##INFO=<ID=GENE,Number=1,Type=String,Description="Vagrent annotated gene">
##INFO=<ID=PROTEIN_CHANGE,Number=1,Type=String,Description="Vagrent annotated protein change">
               ##INFO=<ID=VC,Number=1,Type=String,Description="Vagrent annotated most deletarious variant consequence">
##INFO=<ID=IS_IN_EXCLUDED_REGION,Number=1,Type=Integer,Description="In CNA/LOH or sex chromosome region: 0=Not Excluded,1=Excluded">
               ##FORMAT=<ID=GENOTYPE, Number=1, Type=String, Description="Genotype:0=absent,1=present,?=unknown">
               ##FORMAT=<ID=REF_COUNT,Number=1,Type=Integer,Description="Reads presenting a reference allele for this position">
               ##FORMAT=<ID=ALT_COUNT.Number=1,Type=Integer,Description="Reads presenting the specified alternative allele(s) for this position">
                                                                              2. פורמט המידע של הרצה של בדיקת קיום מוטציה על קולונות
                             germline
                                                                                            somatic
                                                                                                                                  :הוריאנט
                                                                                                                                                                      סוג
 ##TITETORMAL=VCFV4.1
 ##INFO=<ID=MUT_TYPE, Number=1, Type=String, Description="Mutation Type">
 ##INFO=<ID=GENE,Number=1,Type=String,Description="Vagrent annotated gene">
##INFO=<ID=PROTEIN_CHANGE,Number=1,Type=String,Description="Vagrent annotated protein change">
 ##INFO=<ID=VC,Number=1,Type=String,Description="Vagrent annotated most deletarious variant consequence">
##INFO=<ID=VC,Number=1,Type=String,Description="Vagrent annotated most deletarious variant consequence">
##INFO=<ID=VC,Number=1,Type=String,Description="In CNA/LOH or sex chromosome region: 0=Not Excluded,1=Excluded">
##FORMAT=<ID=GENOTYPE,Number=1,Type=String,Description="Genotype:0=absent,1=present,?=unknown">
 ##FORMAT=<ID=REF_COUNT,Number=1,Type=Integer,Description="Reads presenting a reference allele for this position">
 ##FORMAT=<ID=ALT_COUNT_Number=1,Type=Integer,Description="Reads presenting the specified alternative allele(s) for this position">
                             :המוטציה
                                                                                                                                                                    סוג
 ##fileformat=VCFv4.1
##INFO=<ID=SUMATIC,Number=0,Type=riag,Description= variant is Somatic >
 ##INFO=<ID=GENE,Number=1,Type=String,Description="Vagrent annotated gene">
 ##INFO=<ID=PROTEIN_CHANGE,Number=1,Type=String,Description="Vagrent annotated protein change"> ##INFO=<ID=VC,Number=1,Type=String,Description="Vagrent annotated most deletarious variant consequence">
 ##INFO=<ID=IS_IN_EXCLUDED_REGION,Number=1,Type=Integer,Description="In CNA/LOH or sex chromosome region: 0=Not Excluded,1=Excluded">
 ##FORNAT=<ID=REF_COUNT,Number=1,Type=Integer,Description="Reads presenting a reference allele for this position">
 ##FORMAT=<ID=ALT_COUNT_Number=1,Type=Integer_Description="Reads presenting the specified alternative allele(s) for this position">
                                                                                                 SNV/SNP – בודד
                                                                                                              הוספה ומחיקה – INDEL
                                                                                                                                      ג. הגן שנבחן בהרצה זו
               ##fileformat=VCFv4.1
               ##INFO=<ID=SOMATIC,Number=0,Type=Flag,Description="Variant is Somatic">
                ##TNEO-ZTD-MUT TVDE Numbon-1
               ##INFO=<ID=GENE, Number=1, Type=String, Description="Vagrent annotated gene">
                 #INFO=<ID=FKOTEIN_CHANGE,Number=I,Type=StrIng,DeStrIption= vagrent annotated protein change"
               ##INFO=<ID=VC,Number=1,Type=String,Description="Vagrent annotated most deletarious variant consequence">
               ##INFO=<ID=IS_IN_EXCLUDED_REGION,Number=1,Type=Integer,Description="In CNA/LOH or sex chromosome region: 0=Not Excluded,1=Excluded">
##FORMAT=<ID=IS_IN_EXCLUDED_REGION,Number=1,Type=Integer,Description="In CNA/LOH or sex chromosome region: 0=Not Excluded,1=Excluded">
##FORMAT=<ID=GENOTYPE,Number=1,Type=String,Description="Genotype:0=absent,1=present,?=unknown">
               ##FORMAT=<ID=REF_COUNT,Number=1,Type=Integer,Description="Reads presenting a reference allele for this position">
##FORMAT=<ID=ALT_COUNT,Number=1,Type=Integer,Description="Reads presenting the specified alternative allele(s) for this position">
```

ד. לכל מוטציה, החלבון שבו נעשה המוטציה, השינוי שנעשה ואיפה בדיוק בחלבון זה התרחש

```
##fileformat=VCFv4.1
##INFO=<ID=SOMATIC,Number=0,Type=Flag,Description="Variant is Somatic">
##INFO=<ID=NUT_TYPE,Number=1,Type=String,Description="Mutation Type">
##INFO=<ID=NUT_TYPE,Number=1,Type=String,Description="Wagrent annotated gene">
##INFO=<ID=ROTEIN_CHANGE,Number=1,Type=String,Description="Vagrent annotated protein change">
##INFO=<ID=ROTEIN_CHANGE,Number=1,Type=String,Description="Vagrent annotated protein change">
##INFO=<ID=IS_IN_EXCLUDED_REGION,Number=1,Type=Integer,Description="In CNA/LOH or sex chromosome region: @=Not Excluded,1=Excluded">
##FORMAT=<ID=ROTEIN_IN_EXCLUDED_REGION,Number=1,Type=Integer,Description="Genotype:@=absent,1=present,?=unknown">
##FORMAT=<ID=REF_COUNT,Number=1,Type=Integer,Description="Reads presenting a reference allele for this position">
##FORMAT=<ID=ALT_COUNT,Number=1,Type=Integer,Description="Reads presenting the specified alternative allele(s) for this position">
```

ה. לכל מוטציה ההשלכות שקרו בגללה, למשל missense, nonsense, 3 prime

ועוד utr

```
##fileformat=VCFv4.1
##INFO=<ID=SOMATIC,Number=0,Type=Flag,Description="Variant is Somatic">
##INFO=<ID=MUT_TYPE,Number=1,Type=String,Description="Mutation Type">
##INFO=<ID=GENE,Number=1,Type=String,Description="Vagrent annotated gene">
##INFO=<ID=POTFTN_CHANGE_Number=1_Type=String_Description="Vagrent annotated protein_change">
##INFO=<ID=VC,Number=1,Type=String,Description="Vagrent annotated most deletarious variant consequence">
##INFO=<ID=VC,Number=1,Type=String,Description="Vagrent annotated most deletarious variant consequence">
##INFO=<ID=VC,Number=1,Type=String,Description="Genotype:0=absent,1=present,?=unknown">
##FORMAT=<ID=GENOTYPE,Number=1,Type=String,Description="Genotype:0=absent,1=present,?=unknown">
##FORMAT=<ID=REF_COUNT,Number=1,Type=Integer,Description="Reads presenting a reference allele for this position">
##FORMAT=<ID=ALT_COUNT,Number=1,Type=Integer,Description="Reads presenting the specified alternative allele(s) for this position">
```

ו. גן CNA/LOH או גן בכרומוזום המין נכתוב 1 אחרת נכתוב 0 (ש-0 מסמן לכלול אותו) אותו באנליזות מסוימות ו-1 מסמן לא לכלול אותו)

3. הפורמט (מיוצג כשלישיה) של הצגת הנתונים עבור בדיקה על קולונה. (כל שלישיה זו קולונה).

א. האם המוטציה מופיעה (1 אם כן 0 אם לא: מספר הרידים שהייתה בהם מוטציה

: מספר הרידים שלא היה בהם מוטציה

```
##FORMAT=<ID=GENOTYPE,Number=1,Type=String,Description="Genotype:0=absent,1=present,?=unknown">
##FORMAT=<ID=REF_COUNT,Number=1,Type=Integer,Description="Reads presenting a reference allele for this position">
##FORMAT=<ID=ALT_COUNT,Number=1,Type=Integer,Description="Reads presenting the specified alternative allele(s) for this position">
```

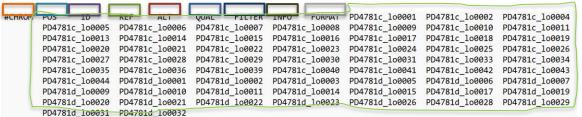
## 4. הכרומוזומים והאורכים שלהם:

```
##contig=<ID=1,assembly=NCBI37,length=249250621,species=Human>
##contig=<ID=2,assembly=NCBI37,length=243199373,species=Human>
##contig=<ID=3,assembly=NCBI37,length=198022430,species=Human>
##contig=<ID=4,assembly=NCBI37,length=191154276,species=Human>
##contig=<ID=5,assembly=NCBI37,length=180915260,species=Human>
##contig=<ID=6,assembly=NCBI37,length=171115067,species=Human>
##contig=<ID=7,assembly=NCBI37,length=159138663,species=Human>
##contig=<ID=8,assembly=NCBI37,length=146364022,species=Human>
##contig=<ID=9,assembly=NCBI37,length=141213431,species=Human>
##contig=<ID=10,assembly=NCBI37,length=135534747,species=Human>
##contig=<ID=11,assembly=NCBI37,length=135006516,species=Human>
##contig=<ID=12,assembly=NCBI37,length=133851895,species=Human>
##contig=<ID=13,assembly=NCBI37,length=115169878,species=Human>
##contig=<ID=14,assembly=NCBI37,length=107349540,species=Human>
##contig=<ID=15,assembly=NCBI37,length=102531392,species=Human>
##contig=<ID=16,assembly=NCBI37,length=90354753,species=Human>
##contig=<ID=17,assembly=NCBI37,length=81195210,species=Human>
##contig=<ID=18,assembly=NCBI37,length=78077248,species=Human>
##contig=<ID=19,assembly=NCBI37,length=59128983,species=Human>
##contig=<ID=20,assembly=NCBI37,length=63025520,species=Human>
##contig=<ID=21,assembly=NCBI37,length=48129895,species=Human>
##contig=<ID=22,assembly=NCBI37,length=51304566,species=Human>
##contig=<ID=X,assembly=NCBI37,length=155270560,species=Human>
##contig=<ID=Y,assembly=NCBI37,length=59373566,species=Human>
```

## 5. המושבות:

```
##SAMPLE=<ID=PD4781c_lo0001,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0001,Source=.>
##SAMPLE=<ID=PD4781c_lo0002,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0002,Source=.>
##SAMPLE=<ID=PD4781c lo0004, Description=., Accession=., Platform=hiseq, Protocol=WGS, SampleName=PD4781c lo0004, Source=.>
##SAMPLE=<ID=PD4781c lo0005, Description=., Accession=., Platform=hiseq, Protocol=WGS, SampleName=PD4781c lo0005, Source=.>
##SAMPLE=<ID=PD4781c_lo0006, Description=., Accession=., Platform=hiseq, Protocol=WGS, SampleName=PD4781c_lo0006, Source=.>
##SAMPLE=<ID=PD4781c_lo0007,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0007,Source=.>
##SAMPLE=<ID=PD4781c_lo0008, Description=., Accession=., Platform=hiseq, Protocol=WGS, SampleName=PD4781c_lo0008, Source=.>
##SAMPLE=<ID=PD4781c_lo0009,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0009,Source=.>
##SAMPLE=<ID=PD4781c_lo0010,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0010,Source=.>
##SAMPLE=<ID=PD4781c_lo0011,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0011,Source=.>
##SAMPLE=<ID=PD4781c lo0013, Description=., Accession=., Platform=hiseq, Protocol=WGS, SampleName=PD4781c lo0013, Source=.>
##SAMPLE=<ID=PD4781c_lo0014, Description=., Accession=., Platform=hiseq, Protocol=WGS, SampleName=PD4781c_lo0014, Source=.>
##SAMPLE=<ID=PD4781c_lo0015,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0015,Source=.>
##SAMPLE=<ID=PD4781c lo0016, Description=., Accession=., Platform=hiseq, Protocol=WGS, SampleName=PD4781c lo0016, Source=.>
##SAMPLE=<ID=PD4781c_lo0017,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0017,Source=.>
##SAMPLE=<ID=PD4781c_lo0018,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0018,Source=.>
##SAMPLE=<ID=PD4781c_lo0019,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0019,Source=.>
##SAMPLE=<ID=PD4781c_lo0020,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0020,Source=.>
##SAMPLE=<ID=PD4781c_lo0021, Description=., Accession=., Platform=hiseq, Protocol=WGS, SampleName=PD4781c_lo0021, Source=.>
##SAMPLE=<ID=PD4781c_lo0022,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0022,Source=.>
##SAMPLE=<ID=PD4781c_lo0023,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0023,Source=.>
##SAMPLE=<ID=PD4781c_lo0024, Description=., Accession=., Platform=hiseq, Protocol=WGS, SampleName=PD4781c_lo0024, Source=.>
##SAMPLE=<ID=PD4781c_lo0025,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0025,Source=.>
##$AMPLE=<ID=PD4781c_lo0026, Description=., Accession=., Platform=hiseq, Protocol=WGS, SampleName=PD4781c_lo0026, Source=.> ##$AMPLE=<ID=PD4781c_lo0027, Description=., Accession=., Platform=hiseq, Protocol=WGS, SampleName=PD4781c_lo0027, Source=.>
##SAMPLE=<ID=PD4781c lo0028, Description=., Accession=., Platform=hiseq, Protocol=WGS, SampleName=PD4781c lo0028, Source=.>
##SAMPLE=<ID=PD4781c_lo0029,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0029,Source=.>
##SAMPLE=<ID=PD4781c_lo0030,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0030,Source=.>
##SAMPLE=<ID=PD4781c_lo0031,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0031,Source=.>
##SAMPLE=<ID=PD4781c_lo0033,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0033,Source=.>
##SAMPLE=<ID=PD4781c_lo0034,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0034,Source=.>
##SAMPLE=<ID=PD4781c_lo0035,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0035,Source=.>
##SAMPLE=<ID=PD4781c lo0036, Description=., Accession=., Platform=hiseq, Protocol=WGS, SampleName=PD4781c lo0036, Source=.>
##SAMPLE=<ID=PD4781c_lo0039,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0039,Source=.>
##SAMPLE=<ID=PD4781c_lo0040, Description=., Accession=., Platform=hiseq, Protocol=WGS, SampleName=PD4781c_lo0040, Source=.>
##SAMPLE=<ID=PD4781c_lo0041,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0041,Source=.>
##SAMPLE=<ID=PD4781c_lo0042,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0042,Source=.>
##SAMPLE=<ID=PD4781c_lo0043,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0043,Source=.>
##SAMPLE=<ID=PD4781c_lo0044,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781c_lo0044,Source=.>
##SAMPLE=<ID=PD4781d_lo0001,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781d_lo0001,Source=.>
##SAMPLE=<ID=PD4781d lo0002, Description=., Accession=., Platform=hiseq, Protocol=WGS, SampleName=PD4781d lo0002, Source=.>
##SAMPLE=<ID=PD4781d_lo0003,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781d_lo0003,Source=.>
##SAMPLE=<ID=PD4781d_lo0005,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781d_lo0005,Source=.>
##SAMPLE=<ID=PD4781d_lo0006,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781d_lo0006,Source=.>
##SAMPLE=<ID=PD4781d_lo0007,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781d_lo0007,Source=.>
##SAMPLE=<ID=PD4781d_lo0009,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781d_lo0009,Source=.>
##SAMPLE=<ID=PD4781d_lo0010, Description=., Accession=., Platform=hiseq, Protocol=WGS, SampleName=PD4781d_lo0010, Source=.>
##SAMPLE=<ID=PD4781d_lo0011, Description=., Accession=., Platform=hiseq, Protocol=WGS, SampleName=PD4781d_lo0011, Source=.>
##SAMPLE=<ID=PD4781d_lo0014, Description=., Accession=., Platform=hiseq, Protocol=WGS, SampleName=PD4781d_lo0014, Source=.>
##SAMPLE=<ID=PD4781d_lo0015,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781d_lo0015,Source=.>
##SAMPLE=<ID=PD4781d_lo0017, Description=., Accession=., Platform=hiseq, Protocol=WGS, SampleName=PD4781d_lo0017, Source=.>
##SAMPLE=<ID=PD4781d_lo0019,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781d_lo0019,Source=.>
##SAMPLE=<ID=PD4781d_lo0020,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781d_lo0020,Source=.>
##SAMPLE=<ID=PD4781d_lo0021,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781d_lo0021,Source=.>
##SAMPLE=<ID=PD4781d lo0022, Description=., Accession=., Platform=hiseq, Protocol=WGS, SampleName=PD4781d lo0022, Source=.>
##SAMPLE=<ID=PD4781d_lo0023,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781d_lo0023,Source=.>
##SAMPLE=<ID=PD4781d_lo0026,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781d_lo0026,Source=.>
##SAMPLE=<ID=PD4781d_lo0028,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781d_lo0028,Source=.>
##SAMPLE=<ID=PD4781d_lo0029,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781d_lo0029,Source=.>
##SAMPLE=<ID=PD4781d_lo0031,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781d_lo0031,Source=.>
##SAMPLE=<ID=PD4781d lo0032,Description=.,Accession=.,Platform=hiseq,Protocol=WGS,SampleName=PD4781d lo0032,Source=.>
```

# 6. הטמפלייט של הרצה על מוטציה:



- א. מספר הכרומוזום שעליו נמצאת המוטציה
- ב. המיקום של הבסיס שבו מתחילה המוטציה
- ג. הזהות של הבן אדם שעליו אנחנו מריצים את זה
  - מה שאמור להיות בלי מוטציה.
  - מה שאמור להיות אחרי מוטציה ...
  - ו. עד כמה אנחנו בטוחים בתוצאה
    - וו. לא בטוח מה זה הFILTER
- ח. זה החלק של המידע של פורמט המידע שהסברתי בהרחבה ב2
  - ט. זה פורמט ההצגה של קולונה שהסברתי בהרחבה ב3
    - אלו השלשות שמייצגות את הקולונות שכתובות ב5