# Personalized Medicine:
## Redefining Cancer Treatment

Matt Shaffer · W207 Final Project · 16 August 2017

## Workflow

1. A molecular pathologist selects a list of genetic variations of interest that he/she want to analyze

2. The molecular pathologist searches for evidence in the medical literature that somehow are relevant to the genetic variations of interest

3. Finally, this molecular pathologist spends a huge amount of time analyzing the evidence related to each of the variations to classify them

**Goal**

Replace step 3 by a machine learning model.

## Features

1. **Gene**

(the gene where this genetic mutation is located)

2. **Variation**
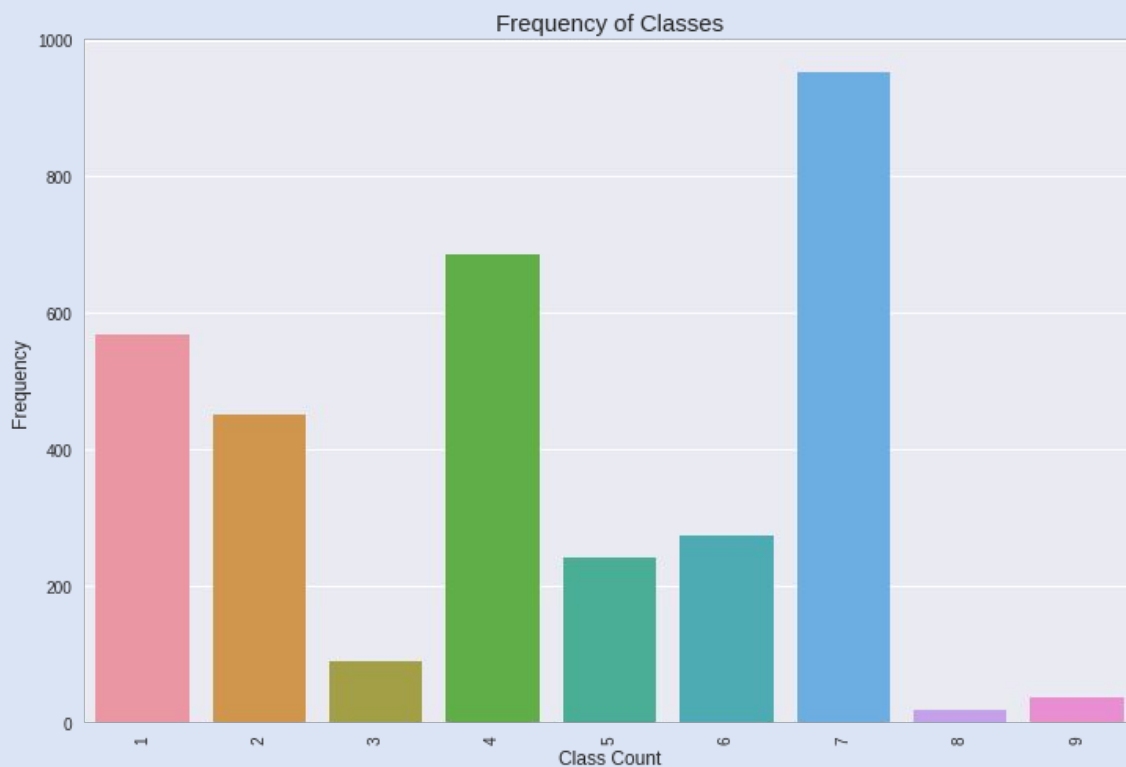
(the aminoacid change for this mutation)

3. **Class**

(1-9 the class this genetic mutation has been classified on)

4. **Text corpus**

(the clinical evidence used to classify the genetic mutation)

| | ID | Gene | Variation | Class | Text |
|---|---|---|---|---|---|
| 1108 | 1108 | FANCA | S858R | 4 | Fanconi anemia (FA) is an autosomal recessive ... |
| 1109 | 1109 | FANCA | S1088F | 1 | null |
| 1110 | 1110 | FANCA | Truncating Mutations | 1 | Abstract Fanconi anemia is characterized by c... |
| 1111 | 1111 | FANCA | H492R | 4 | Abstract Fanconi anemia (FA) is a genomic ins... |
| 1112 | 1112 | FANCA | Y510C | 4 | Abstract Fanconi anemia (FA) is a genomic ins... |
| 1113 | 1113 | FANCA | Deletion | 1 | Fanconi anemia (FA) is a genetic disease chara... |
| 1114 | 1114 | FANCA | L274P | 4 | Abstract Fanconi anemia (FA) is a genomic ins... |
| 1115 | 1115 | FANCA | W183A | 4 | Fanconi anemia (FA) is a recessively inherited... |
| 1116 | 1116 | FANCA | L210R | 4 | Abstract Fanconi anemia (FA) is a genomic ins... |



Frequency of Classes

# CLASSES

1. Likely Loss-of-function

2. Likely Gain-of-function

3. Neutral

4. Loss-of-function

5. Likely Neutral

6. Inconclusive

7. Gain-of-function
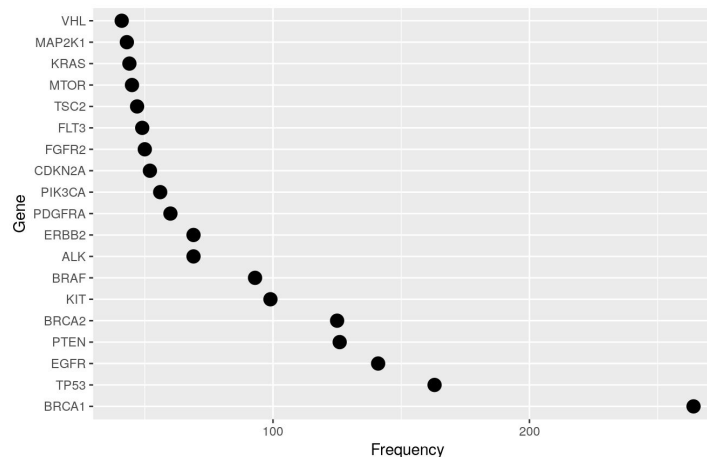
8. Likely Switch-of-function

9. Switch-of-function

# EDA

## Inconsistent Data

| | Class | Gene | ID | Variation |
|---|---|---|---|---|
| 3716 | NaN | RUNX2 | 395 | null522S |
| 3967 | NaN | PAX6 | 646 | null423L |
| 4017 | NaN | SHOX | 696 | null293R |
| 4540 | NaN | ITM2B | 1219 | null267R |
| 4749 | NaN | SH2D1A | 1428 | null129R |
| 4787 | NaN | FKRP | 1466 | null496R |
| 4859 | NaN | PNPO | 1538 | null262Q |
| 5428 | NaN | HSD3B2 | 2107 | null373C |
| 5644 | NaN | SELENON | 2323 | null462G |
| 5688 | NaN | KISS1R | 2367 | null399R |
| 5738 | NaN | IDUA | 2417 | null654G |
| 5862 | NaN | RAD50 | 2541 | null1313Y |
| 5907 | NaN | FHL1 | 2586 | null281E |
| 5952 | NaN | MOCS2 | 2631 | null189Y |
| 6094 | NaN | IKBKG | 2773 | null420W |
| 6899 | NaN | CTSK | 3578 | null330W |
| 7151 | NaN | DBT | 3830 | null483L |
| 7476 | NaN | NHP2 | 4155 | null154R |
| 8188 | NaN | FOXF1 | 4867 | null380R |

## Missing Values

```
X_test.loc[X_test['Text'].str.len() < 100]
```

| | ID | Gene | Variation | Text |
|---|---|---|---|---|
| 1623 | 1623 | AURKB | Amplification | null |

## Observations Disproportionately Represented



## Shared Text Corpus for Multiple Variations

| | Text | text_length | Gene | Variation |
|---|---|---|---|---|
| 3298 | Introduction Myelodysplastic syndromes (MDS) ... | 40127 | RUNX1 | Y113* |
| 3303 | Introduction Myelodysplastic syndromes (MDS) ... | 40127 | RUNX1 | P173S |
| 3305 | Introduction Myelodysplastic syndromes (MDS) ... | 40127 | RUNX1 | S70fsX93 |
| 3317 | Introduction Myelodysplastic syndromes (MDS) ... | 40127 | RUNX1 | A122* |
| 3316 | Introduction Myelodysplastic syndromes (MDS) ... | 73895 | RUNX1 | D171N |
| 3314 | Introduction Myelodysplastic syndromes (MDS) ... | 94151 | RUNX1 | G42R |

# TF-IDF

- Bigram– 2690998 tokens
- Trigrams – 33,126,986 tokens
- 10% of Vocabulary

# SVD



Explained Variance Vs. Number of Features

- 20 Features for dataset scaled to 10% of original
- 200 Final Model

[('128 MN1', 37587),
 ('problem genetic', 2394330),
 ('size inversely', 2649165),
 ('Invitrogen catalog', 669650),
 ('methylation mutational', 2115605),
 ('hotspot SNPs', 1850576),
 ('Endogenous MyD88', 500475),
 ('supplement 1B', 2732633),
 ('Recently NUP98', 921101),
 ('phase data', 2322665),
 ('resistance clinical', 2533397),
 ('site frequency', 2646231),
 ('underline 712', 2843950),
 ('pGBT9 TRP1', 2276264),
 ('Fgfr3 Viable', 533516),
 ('pY869 detected', 2282312),
 ('665752 JNJ38877605', 225464),
 ('Research Inc', 926139),
 ('key advance', 1985754),
 ('site Their', 2645338)]

[('day 43 005', 2358384),
 ('differently result combined', 2583154),
 ('man 191100 clinical', 4757053),
 ('ssa hdr', 7250466),
 ('10 dl hematocrit', 37545),
 ('medulloblastomas 68 demonstrated', 4851125),
 ('functioning different pathway', 3452747),
 ('analtech newark developed', 1070313),
 ('rasq61r control', 6362798),
 ('tumor tissue sts', 7876878),
 ('nk granulo monocytic', 5343760),
 ('block cdh1', 1433634),
 ('vivo study pk', 8146749),
 ('containing p53 dna', 2137352),
 ('deficiency promotes differentiated', 2392257),
 ('pvh1213 lysine', 6268396),
 ('project bi78d3 santa', 6135025),
 ('counted study property', 2214788),
 ('hdmec hemec transfected', 3747309),

| | Substitution | Insertion | Deletion |
|---|---|---|---|
| Original sequence | T G G **C** A G | T G G C A G | T G G ~~C A~~ G |
| Mutated sequence | T G G **T** A G | T G G **T A T** C A G | T G G G |

**tmVar normalization format:**

Substitution:
**<Sequence type>|SUB|<wild type>|<mutation position>|<mutant>**
e.g., "c.435C>G" --> "c|SUB|C|435|G"

Deletion:
**<Sequence type>|DEL|<mutation position>|<mutant>**
e.g., "c.104delT" --> "c|DEL|104|T"
e.g., "c.1544-?_2916+?" --> "c|DEL|1544-?_2916+?|"

Insertion:
**<Sequence type>|INS|<mutation position>|<mutant>**
e.g., "c.104insT" --> "c|INS|104|T"

Insertion+Deletion:
**<Sequence type>|INDEL|<mutation position>|<mutant>**
e.g., "c.2153_2155delinsTCCTGGTTTA" -->
"c|INDEL|2153_2155|TCCTGGTTTA"

Duplication:
**<Sequence type>|DUP|<mutation position>|<mutant>|<duplication times>**
e.g., "c.1285-1301dup" --> "c|DUP|1285_1301||"
e.g., "c.1978(TATC)(1-2)" --> "c|DUP|1978|TATC|1-2"

Frame shift:
**<Sequence type>|FS|<wild type>|<mutation position>|<mutant>|<frame shift position>**
e.g., "p.Val35AlafsX25" --> "p|FS|V|35|A|25"
e.g., "p.Ser119fsX" --> "p|FS|S|119||"

<Sequence type>:
c: DNA sequence
r: RNA sequence
g: Genome sequence
p: Protein sequence
m: Mitochondrial sequence

<wild type> / <mutant>:
A,T,C,G: DNA nucleotide
C,I,S,Q,M,N,P,K,D,T,F,A,G,H,L,R,W,V,E,Y,X: Amino acid

C630R

C 630 R

Cys 630 Arg

Cysteine 630 Arginine

| | | | |
|---|---|---|---|
| 2943 | C630R | | |
| 2944 | V648I | | |
| 2945 | I852M | | |
| 2946 | C620R | | |
| 2947 | C634Y | | |
| 2948 | V804G | | |
| 2949 | R886W | | |
| 2950 | F893L | | |
| 2951 | Y791F | | |
| 2952 | R177* | | |
| 2953 | Y113* | | |
| 2954 | R139G | | |
| 2955 | K83N | | |
| 2956 | R177Q | | |
| 2957 | R166Q | | |
| 2958 | P173S | | |
| 2959 | R201Q | | |
| 2960 | S70fsX93 | | |
| 2961 | W279* | | |
| 2962 | R174* | | |
| 2963 | D171G | | |
| 2964 | RUNX1-EVI1 Fusion | | |
| 2965 | TEL-RUNX1 Fusion | | |
| 2966 | H78Q | | |
| 2967 | G42R | | |
| 2968 | RUNX1-RUNX1T1 Fusion | | |
| 2969 | D171N | | |
| 2970 | A122* | | |
| 2971 | R80C | | |
| 2972 | K83E | | |

| | ID | Gene | Variation | Class |
|---|---|---|---|---|
| 138 | 138 | EGFR | L747_T751delinsP | 7 |
| 139 | 139 | EGFR | S752_I759del | 2 |
| 141 | 141 | EGFR | D770_P772dup | 7 |
| 144 | 144 | EGFR | N771_H773dup | 7 |
| 146 | 146 | EGFR | E746_T751insIP | 7 |
| 147 | 147 | EGFR | D770_N771insD | 7 |
| 149 | 149 | EGFR | K745_A750del | 7 |
| 165 | 165 | EGFR | D770_N771insNPG | 7 |
| 166 | 166 | EGFR | E746_A750del | 7 |
| 171 | 171 | EGFR | A859_L883delinsV | 2 |
| 174 | 174 | EGFR | A750_E758del | 7 |
| 175 | 175 | EGFR | V769_D770insGVV | 7 |
| 184 | 184 | EGFR | A750_E758delinsP | 7 |
| 187 | 187 | EGFR | L747_P753delinsS | 7 |

| Input | HGVS Committee | HGVS ClinVar/NCBI | HGVS Ensembl | HGVS Mutalyzer |
|---|---|---|---|---|
| m.8993T>G | m.8993T>G | NC_012920.1:m.8993T>G | MT:g.8993T>G | NC_012920.1:g.8993T>G |
| 8993G | m.8993T>G | NC_012920.1:m.8993T>G | MT:g.8993T>G | NC_012920.1:g.8993T>G |
| T8993G | m.8993T>G | NC_012920.1:m.8993T>G | MT:g.8993T>G | NC_012920.1:g.8993T>G |
| 8993d | m.8993_8993del | NC_012920.1:m.8993_8993del | MT:g.8993_8993del | NC_012920.1:g.8993_8993del |
| 8527 | m.8527A>G | NC_012920.1:m.8527A>G | MT:g.8527A>G | NC_012920.1:g.8527A>G |
| 8527A>G | m.8527A>G | NC_012920.1:m.8527A>G | MT:g.8527A>G | NC_012920.1:g.8527A>G |
| MT.6328C>T | m.6328C>T | NC_012920.1:m.6328C>T | MT:g.6328C>T | NC_012920.1:g.6328C>T |
| 8042_8043d | m.8042_8043del | NC_012920.1:m.8042_8043del | MT:g.8042_8043del | NC_012920.1:g.8042_8043del |
| 1494.1T | m.1494_1495insT | NC_012920.1:m.1494_1495insT | MT:g.1494_1495insT | NC_012920.1:g.1494_1495insT |
| 7472.XA | m.7472_7473insAA | NC_012920.1:m.7472_7473insAA | MT:g.7472_7473insAA | NC_012920.1:g.7472_7473insAA |

# Variant Types

| | 0 |
|---|---|
| 0 | EGFRvV |
| 1 | Hypermethylation |
| 2 | TRKAIII Splice Variant |
| 3 | Promoter Mutations |
| 4 | Deletion |
| 5 | Copy Number Loss |
| 6 | DNA binding domain deletions |
| 7 | Wildtype |
| 8 | DNA binding domain insertions |
| 9 | Epigenetic Silencing |
| 10 | MYC-nick |
| 11 | EGFRvIII |
| 12 | Overexpression |
| 13 | Truncating Mutations Upstream of Transactivati... |
| 14 | Amplification |
| 15 | Truncating Mutations in the PEST Domain |
| 16 | Single Nucleotide Polymorphism |
| 17 | Truncating Mutations |
| 18 | Promoter Hypermethylation |
| 19 | DNA binding domain missense mutations |
| 20 | EGFR-KDD |
| 21 | EGFRvII |
| 22 | EGFRvIV |



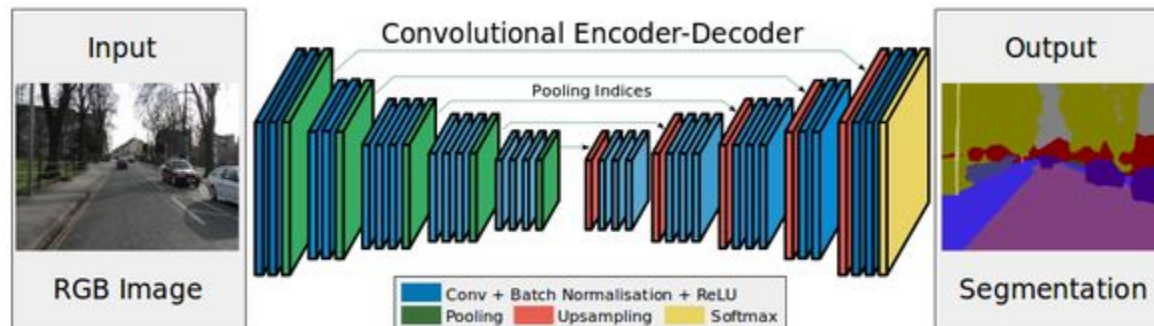Probability of classes in Training Set based on Variation Identifier Type

| Amino acid | 3-letter[132] | 1-letter[132] | Side chain class | Side chain polarity[132] | Side chain charge (pH 7.4)[132] | Hydropathy index[133] | Absorbance $\lambda_{max}$(nm)[134] | ε at $\lambda_{max}$ (mM$^{-1}$ cm$^{-1}$)[134] | MW (weight) | Occurrence in proteins (%)[135] |
|---|---|---|---|---|---|---|---|---|---|---|
| Alanine | Ala | A | aliphatic | nonpolar | neutral | 1.8 | | | 89.094 | 8.76 |
| Arginine | Arg | R | basic | basic polar | positive | −4.5 | | | 174.203 | 5.78 |
| Asparagine | Asn | N | amide | polar | neutral | −3.5 | | | 132.119 | 3.93 |
| Aspartic acid | Asp | D | acid | acidic polar | negative | −3.5 | | | 133.104 | 5.49 |
| Cysteine | Cys | C | sulfur-containing | nonpolar | neutral | 2.5 | 250 | 0.3 | 121.154 | 1.38 |
| Glutamic acid | Glu | E | acid | acidic polar | negative | −3.5 | | | 147.131 | 6.32 |
| Glutamine | Gln | Q | amide | polar | neutral | −3.5 | | | 146.146 | 3.9 |
| Glycine | Gly | G | aliphatic | nonpolar | neutral | −0.4 | | | 75.067 | 7.03 |
| Histidine | His | H | basic aromatic | basic polar | positive(10%) neutral(90%) | −3.2 | 211 | 5.9 | 155.156 | 2.26 |
| Isoleucine | Ile | I | aliphatic | nonpolar | neutral | 4.5 | | | 131.175 | 5.49 |
| Leucine | Leu | L | aliphatic | nonpolar | neutral | 3.8 | | | 131.175 | 9.68 |
| Lysine | Lys | K | basic | basic polar | positive | −3.9 | | | 146.189 | 5.19 |
| Methionine | Met | M | sulfur-containing | nonpolar | neutral | 1.9 | | | 149.208 | 2.32 |
| Phenylalanine | Phe | F | aromatic | nonpolar | neutral | 2.8 | 257, 206, 188 | 0.2, 9.3, 60.0 | 165.192 | 3.87 |
| Proline | Pro | P | cyclic | nonpolar | neutral | −1.6 | | | 115.132 | 5.02 |
| Serine | Ser | S | hydroxyl-containing | polar | neutral | −0.8 | | | 105.093 | 7.14 |
| Threonine | Thr | T | hydroxyl-containing | polar | neutral | −0.7 | | | 119.119 | 5.53 |
| Tryptophan | Trp | W | aromatic | nonpolar | neutral | −0.9 | 280, 219 | 5.6, 47.0 | 204.228 | 1.25 |
| Tyrosine | Tyr | Y | aromatic | polar | neutral | −1.3 | 274, 222, 193 | 1.4, 8.0, 48.0 | 181.191 | 2.91 |
| Valine | Val | V | aliphatic | nonpolar | neutral | 4.2 | | | 117.148 | 6.73 |

## Dense Network

```python
model = Sequential()
model.add(Dense(512, input_dim=input_shape, kernel_initializer='normal', activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(256, kernel_initializer='normal', activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(128, kernel_initializer='normal', activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(64, kernel_initializer='normal', activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(128, kernel_initializer='normal', activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(256, kernel_initializer='normal', activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(512, kernel_initializer='normal', activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(output_shape, kernel_initializer='normal', activation="softmax"))
model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
```

**(Similar Idea)**



Input — RGB Image

Convolutional Encoder-Decoder — Pooling Indices

Output — Segmentation

Conv + Batch Normalisation + ReLU
Pooling   Upsampling   Softmax

## Lessons

Feature engineering takes a long time.

Genetics is complicated

Text mining is hard with limited data

## Still to Try

Further exploration with models

Parsing external data sources:

Collect more text data using APIs

Collect more data on genes using APIs

Sampling methods to overcome data imbalance

Current best: 443 of 790
Score 0.82386