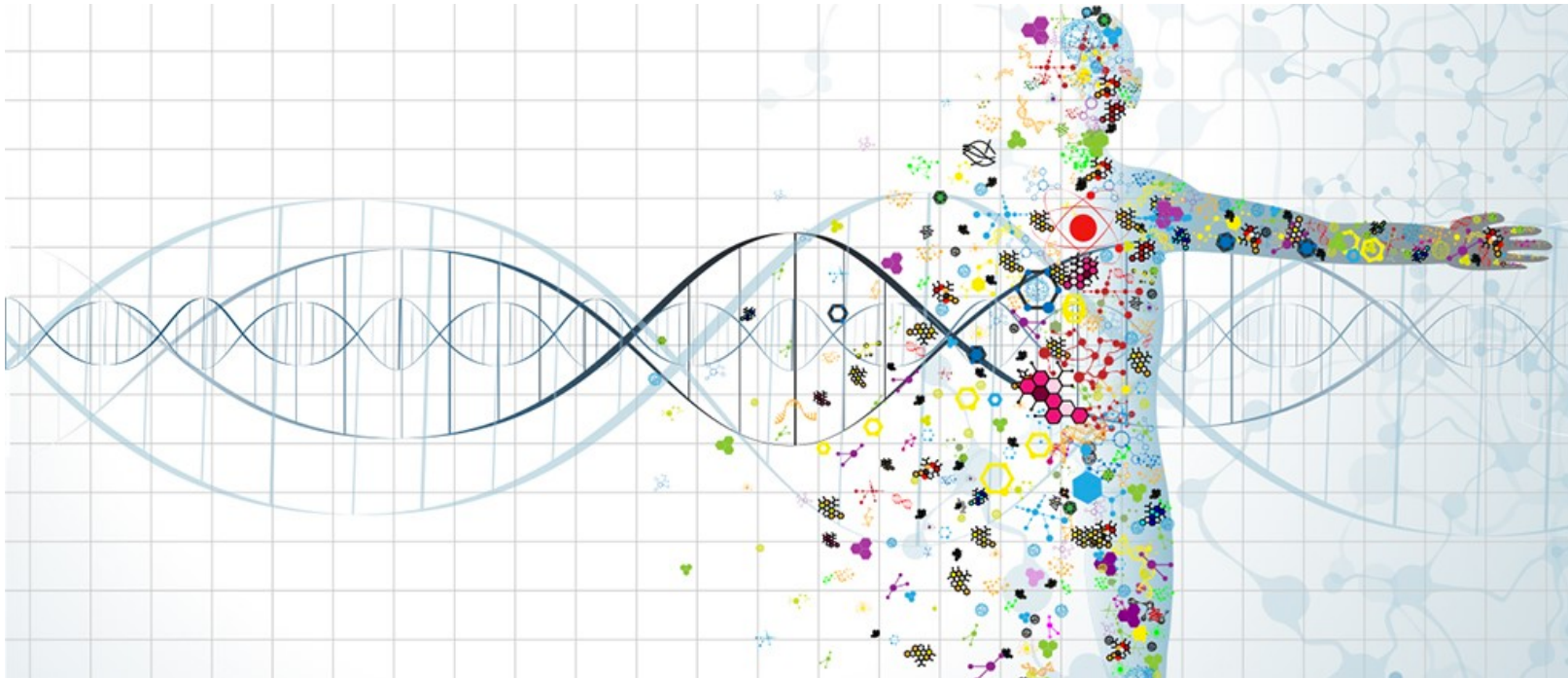
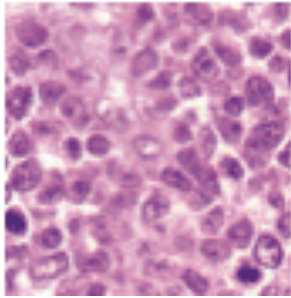


# Predicting Lymphoma Subtype from Patient Chromosome Copy Number Data



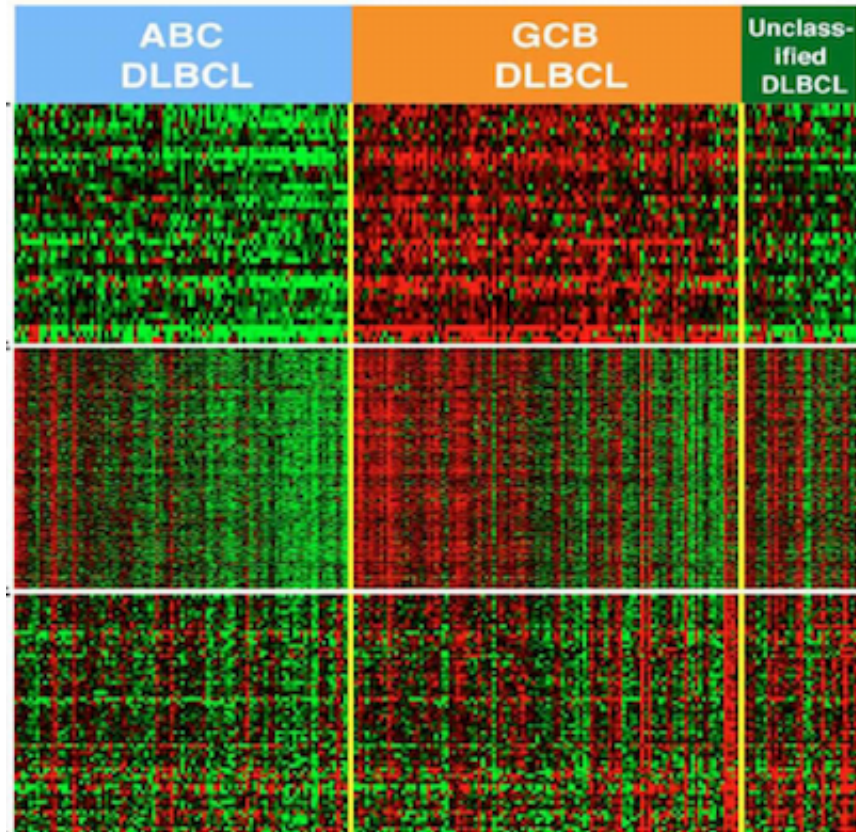
# Diffuse Large B-Cell Lymphoma (DLBCL)

DLBCL

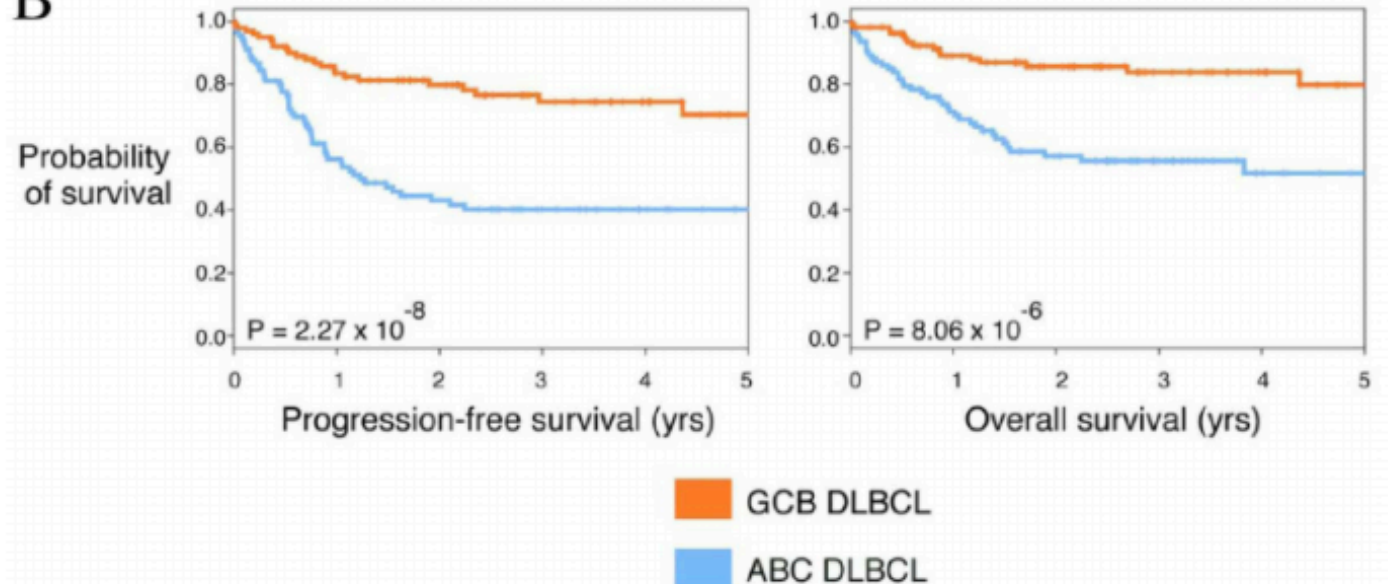


We currently divide DLBCL into three major subtypes: ABC, GCB, and Unclassified.

The cells from these lymphomas look identical under the microscope, but have unique genetic signatures, respond differently to drugs, and the patients have different survival outcomes.



B



# Goal

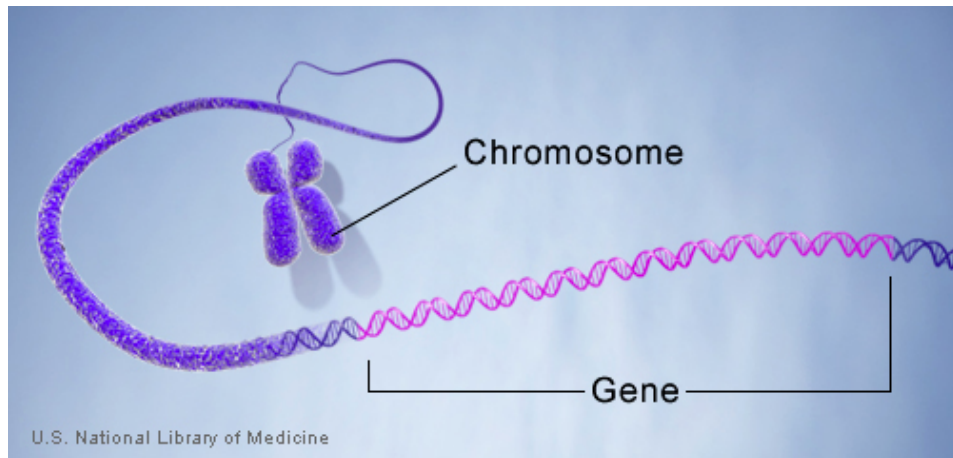
Use chromosome copy number data to predict DLBCL subtype

- Features = copy number data
- Target = DLBCL subtype

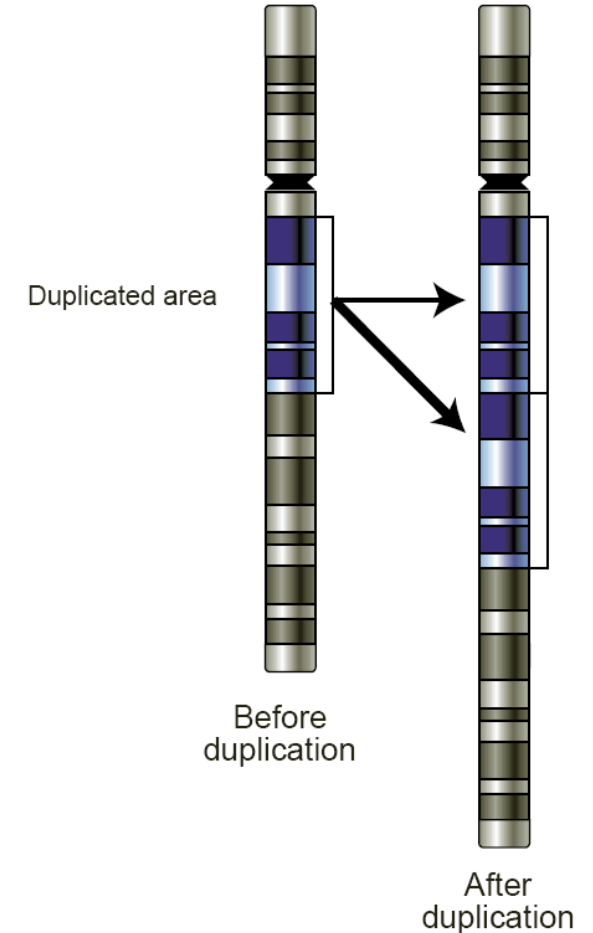
# Example of data from one patient

539 Patients

Each patient has their own file



Chromosome	Start	Stop	Copy Number
2 chr1	51723	149236241	0
3 chr1	149236241	149251876	2
4 chr1	149251876	249234376	0
5 chr2	2772	243099444	0
6 chr3	50333	197906005	0
7 chr4	2269	144920433	0
8 chr4	144920433	145044529	2
9 chr4	145044529	161880120	0
10 chr4	161880120	161884742	-1
11 chr4	161884742	183319181	0
12 chr4	183319181	183399086	-1
13 chr4	183399086	191030138	0
14 chr5	5520	158300100	1
15 chr5	158300100	158446952	-1
16 chr5	158446952	180729790	1
17 chr6	139649	23743214	0
18 chr6	23743214	23746575	-2
19 chr6	23746575	170992522	0
20 chr7	33259	109437087	0.5
21 chr7	109437087	109453934	-1.5
22 chr7	109453934	126046316	0.5
23 chr7	126046316	126050494	-2
24 chr7	126050494	159129708	0.5
25 chr8	36386	7236424	0
26 chr8	7240368	7787836	1
27 chr8	7787836	11009995	0
28 chr8	11009995	11011334	-2
29 chr8	11011334	12237105	0
30 chr8	12237105	12403096	2





# Macroscopic look at the data

```
#Read in all .txt patient chromosome copy number files into
for filename in os.listdir('.'):
    if filename.endswith('.txt'):
        filename_df = pd.read_table(filename,\
names = ['chr', 'start', 'stop', 'copynumber'], index_c
filename_df.drop(['start', 'stop'], axis=1, inplace=T
filename_df = filename_df.groupby('chr').sum()
filename_df = filename_df.T
filename_df['Sample'] = re.sub('\.txt$', '' , str(fi
filename_df.set_index('Sample', inplace=True)
df1 = df1.append(filename_df)
df1 = df1.iloc[:,0:24]
```

```
#Merge another data frame to include the DLBCL Subtype in th
df2 = pd.read_excel('DLBCL_Subtype.xlsx')
df2.set_index('Sample', inplace = True)
df1 = df1.join(df2, how='left')
df1.reset_index(inplace=True)
del df1['Sample']
df1.dropna(inplace=True)
df1['DLBCL Subtype Map'] = df1['DLBCL Subtype']\
.map({'ABC' : 1, 'GCB' : 0, 'Unclass' : 2})
```

+ 539  
patient files



	chr1	chr10	chr11	chr12	chr13	chr14	chr15	chr16	chr17	chr18
0	2.0	0.0	0.0	0.0	0.0	1.0	-2.0	-2.0	2.0	0.0
1	-4.0	1.0	0.0	-2.0	-4.0	-2.0	-2.0	0.0	3.0	3.0
2	-1.0	-1.0	-3.0	2.0	-2.0	-5.0	-3.0	-2.0	-3.0	3.0
3	-6.0	-2.0	0.0	-2.0	-1.0	1.0	-2.0	-2.0	-2.0	6.0
4	-3.0	-1.0	-5.0	2.0	-2.0	-4.5	-6.0	-1.0	5.0	15.0

	...	chr4	chr5	chr6	chr7	chr8	chr9	chrX	chrY	...
0	...	0.0	7.0	-3.0	0.0	0.0	-2.0	4.0	0.0	
1	...	-6.0	0.0	1.0	-1.0	3.0	-3.0	-1.0	0.0	
2	...	-1.0	2.0	-10.0	1.0	-3.0	-2.0	5.0	1.0	
3	...	-5.0	-1.0	-5.0	6.0	0.0	-2.0	1.0	-1.0	
4	...	-2.0	2.0	-5.0	-6.0	-2.0	-4.0	6.0	-2.0	

	DLBCL Subtype	DLBCL Subtype Map
0	ABC	1
1	ABC	1
2	ABC	1
3	ABC	1
4	ABC	1

[5 rows x 26 columns]  
(539, 26)

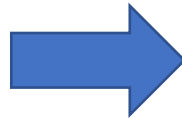
Final output:  
539 rows, 26 columns

# Microscopic look at the data

```
1 df2 = pd.DataFrame()
2
3 for filename in os.listdir('.'):
4     if filename.endswith('.txt'):
5         filename_df = pd.read_table(filename, names =
6             ['chr', 'start', 'stop', 'copy_number'], hea
7
8         filename_df['start'] = filename_df['start'].a
9         filename_df['stop'] = filename_df['stop'].ast
10
11         filename_df['copynumber_master'] = filename_d
12             filename_df['start'] + '_' + filename_df[
13
14         filename_df.drop(['start', 'stop', 'chr'], axi
15         filename_df = filename_df.groupby('copynumber
16         filename_df = filename_df.T
17         filename_df['Sample'] = re.sub('\.txt$', '' ,
18         filename_df.set_index('Sample', inplace = Tru
19         df2 = df2.append(filename_df)
20         df2.fillna(0, inplace = True)
```

```
1 df = pd.read_excel('DLBCL_Subtype.xlsx')
2 df.set_index('Sample', inplace=True)
3 df2 = df2.join(df, how = 'left')
4 df2.reset_index(inplace=True)
5 del df2['Sample']
6 df2.dropna(inplace=True)
```

+ 539  
patient files



	chrY_9390879_28809923	chrY_9432191_9443799	chrY_9439059_9441690
0	0.0	0.0	0.0
1	0.0	0.0	0.0
2	0.0	0.0	0.0
3	0.0	0.0	0.0
4	0.0	0.0	0.0

	chrY_9441690_22222983	chrY_9441690_28809923	chrY_9443799_9787623
0	0.0	0.0	0.0
1	0.0	0.0	0.0
2	0.0	0.0	0.0
3	0.0	0.0	0.0
4	0.0	0.0	0.0

	chrY_9758611_13934470	chrY_9774900_22066345	chrY_9787623_28809923
0	0.0	0.0	0.0
1	0.0	0.0	0.0
2	0.0	0.0	0.0
3	0.0	0.0	0.0
4	0.0	0.0	0.0

Final output:  
539 rows, 54,620 columns

# Models

```
ABC      278
GCB      154
Unclass   107
Name: DLBCL Subtype, dtype: int64
```

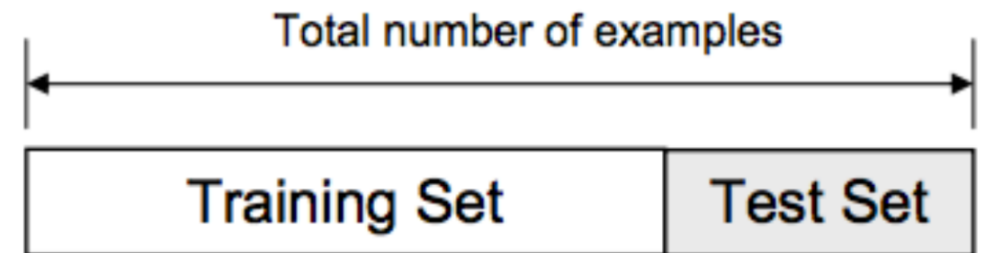
```
The overall null accuracy rate for the whole data set is: 52%
```

Use train, test split to narrow down the best model using accuracy score

- Logistic regression
- K-nearest neighbors
- Random forests

Double check with cross-validation

Apply bagging and boosting

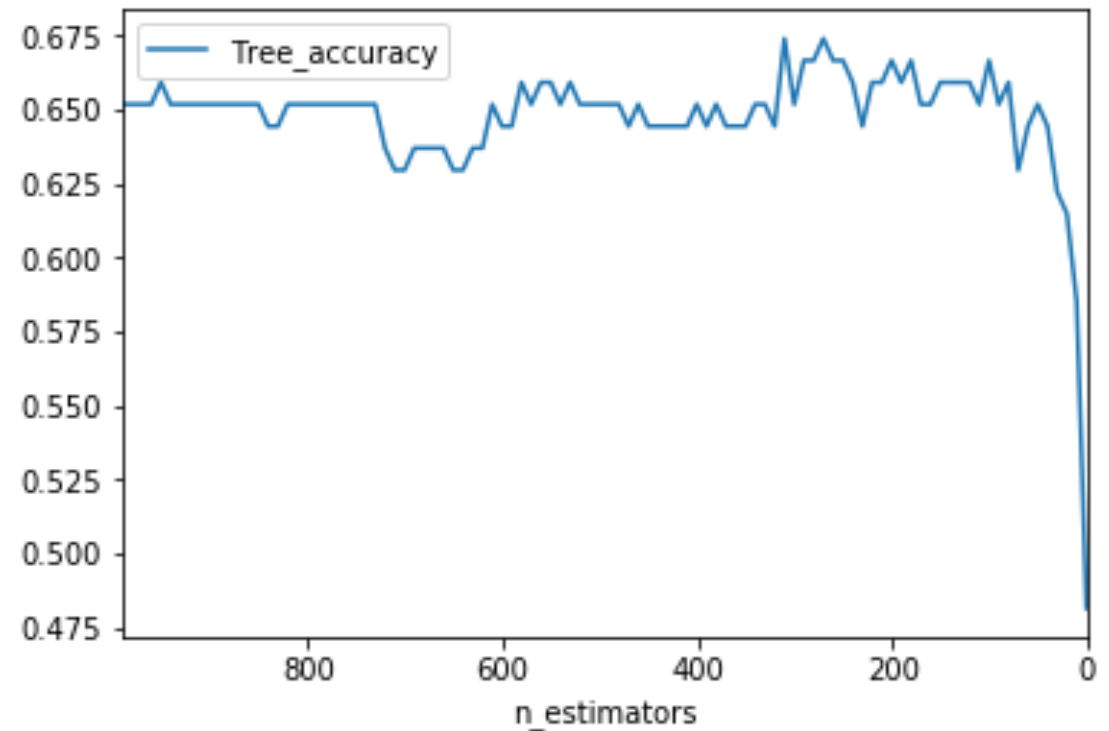


# Macroscopic view (24 features)

Random forest is best model

Randomized and grid searches for optimal number of trees and depth

```
1 from sklearn.ensemble import RandomForestClassifier
2
3 depth_range = range(1,1000,10)
4
5 tree_accuracy_list = []
6
7 def tree_score(depth_range):
8     for i in depth_range:
9         treeclass = RandomForestClassifier(n_estimators
10         X_train, X_test, y_train, y_test = train_test_s
11         treeclass.fit(X_train, y_train)
12         y_pred = treeclass.predict(X_test)
13         tree_accuracy = accuracy_score(y_test, y_pred)
14         tree_accuracy_list.append(tree_accuracy)
15
16 tree_score(depth_range)
17 tree_dict = {'n_estimators': depth_range, 'Tree_accuracy
18 df_tree = pd.DataFrame(tree_dict).set_index('n_estimator
19
20 df_tree.plot(y='Tree_accuracy')
21
22 df_tree.sort_values('Tree_accuracy', ascending = False).
```



~65-67% for highest accuracy scores

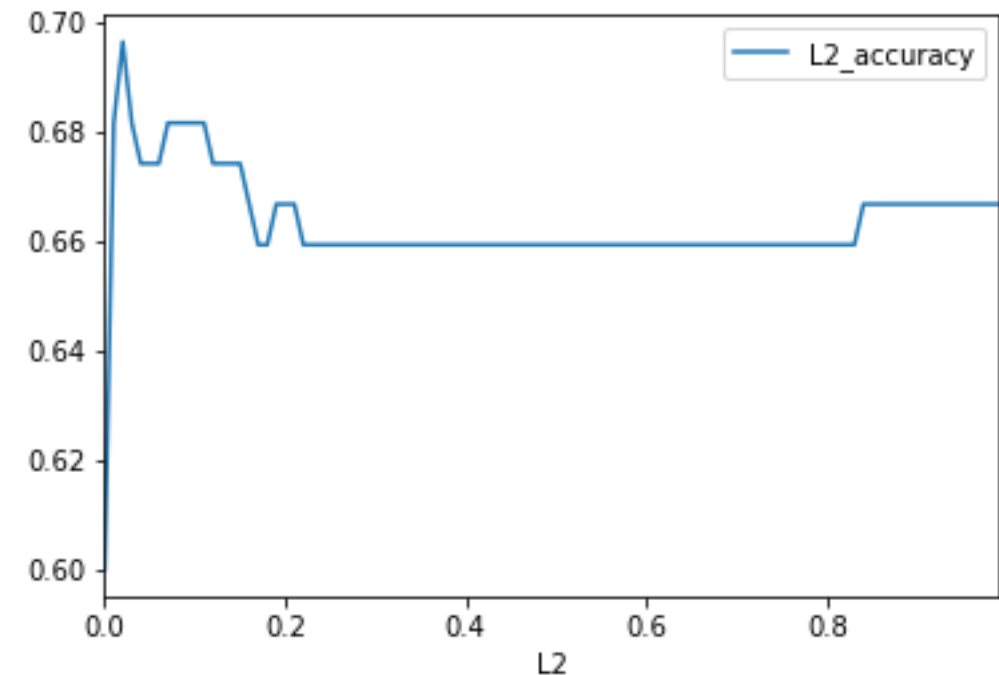


# Microscopic view (54,620 features)

Logistic regression is best model

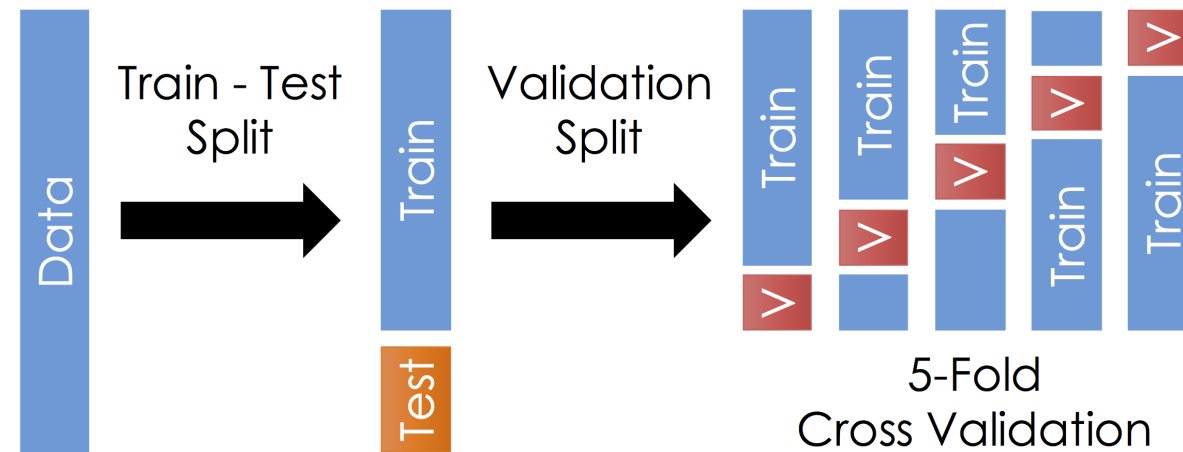
Randomized and grid searches for optimal number regularization method and “C” value

```
1 L2_range = np.arange(.0001, 1, 0.01)
2
3 L2_accuracy_list = []
4
5 def L2_logreg_accuracy(L2_range):
6     for i in L2_range:
7         logreg = LogisticRegression(C = i, random_state = 42)
8         logreg.fit(X_train, y_train)
9         y_pred = logreg.predict(X_test)
10        L2_accuracy = accuracy_score(y_test, y_pred)
11        L2_accuracy_list.append(L2_accuracy)
12
13 L2_logreg_accuracy(L2_range)
14 L2_dict = {'L2' : L2_range, 'L2_accuracy' : L2_accuracy_list}
15 L2_df = pd.DataFrame(L2_dict).set_index('L2')
16 L2_df.plot(y = 'L2_accuracy')
17 L2_df.sort_values('L2_accuracy', ascending = False).head()
```



~66-68% for highest accuracy scores

# Cross-validation for both forms of data



## Macroscopic view of data

```
1 from sklearn.model_selection import RandomizedSearch
2
3 clf = RandomForestClassifier()
4
5 depth_range = range(1, 100, 5)
6 n_range = range(1, 1000, 10)
7
8 param_dist = dict(n_estimators = n_range, max_depth
9
10 rand = RandomizedSearchCV(clf, param_dist, cv = 5, s
11                          n_iter = 100, random_sta
12
13 rand.fit(X, y)
14
15 print(rand.best_score_)
16 rand.best_params_
```

0.643784786642

{'max\_depth': 66, 'n\_estimators': 231}

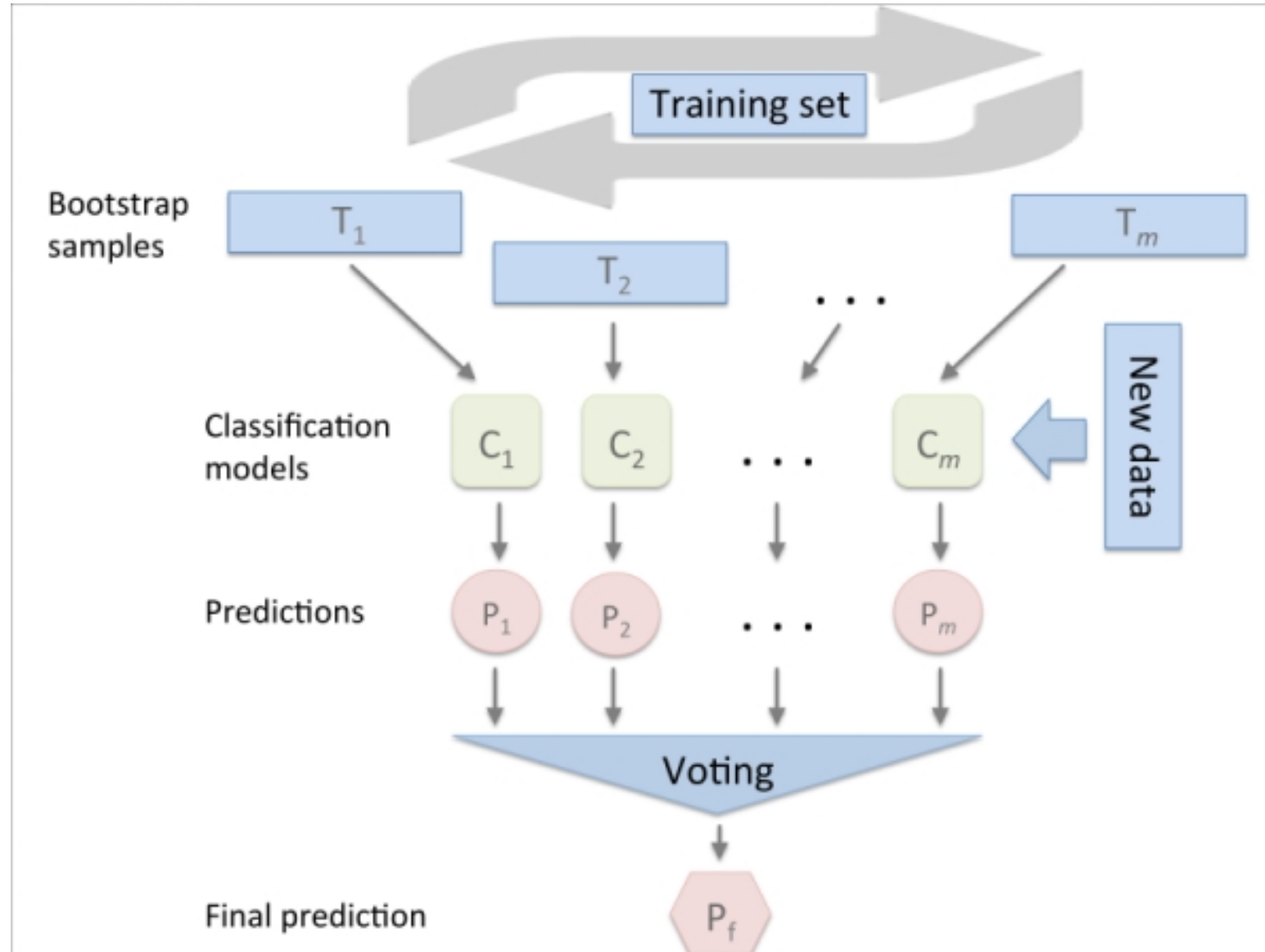
## Microscopic view of data

```
1 from sklearn.model_selection import RandomizedSearch
2
3 clf = RandomForestClassifier()
4
5 depth_range = range(1, 250, 5)
6 n_range = range(1, 1000, 10)
7
8 param_dist = dict(n_estimators = n_range, max_depth
9
10 rand = RandomizedSearchCV(clf, param_dist, cv = 5,
11                          n_iter = 50, random_sta
12
13 rand.fit(X, y)
14
15 print(rand.best_score_)
16 rand.best_params_
```

0.641929499072

{'max\_depth': 221, 'n\_estimators': 261}

# Bagging



# Bagging

Macroscopic view of data  
(random forest bagging)

```
1 from sklearn.ensemble import RandomForestClassifier
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
3 clf = RandomForestClassifier(max_depth=2, random_state=0)
4 bag = BaggingClassifier(base_estimator=clf, n_estimators=100)
5 bag.fit(X_train, y_train)
6 y_pred = bag.predict(X_test)
7 accuracy_score(y_pred, y_test)
```

1.65925925925925921

Microscopic view of data  
(logistic regression bagging)

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
2 logreg = LogisticRegression()
3 bag = BaggingClassifier(base_estimator=logreg, n_estimators=100)
4 bag.fit(X_train, y_train)
5 y_pred = bag.predict(X_test)
6 accuracy_score(y_pred, y_test)
```

0.688888888888888888

Bagging does not really improve accuracy score  
Gradient boosting and Adaboost also did not improve accuracy score

# Conclusion

- Can not seem to beat ~65% accuracy with any models
- DLBCL subtype must be explained by more than just copy number or we need more samples/accurate data
  - DNA mutations
  - Epigenetics

<b>Normal DNA</b>	TAT	CAT	CCT	AAG	GTA	
	└┐	└┐	└┐	└┐	└┐	
<b>Protein</b>	Tyr	His	Pro	Lys	Val	
<b>Substitution</b>	TAT	CAT	CGT	AAG	GTA	
	└┐	└┐	└┐	└┐	└┐	
<b>Protein</b>	Tyr	His	Arg	Lys	Val	
<b>Insertion</b>	TAT	CAT	CGC	TAA	GGT	A
	└┐	└┐	└┐	└┐	└┐	
<b>Protein</b>	Tyr	His	Arg	Stop	Gly	
<b>Deletion</b>	TAT	C_TC	CTA	AGG	TA	
	└┐	└┐	└┐	└┐	└┐	
<b>Protein</b>	Tyr	Leu	Leu	Arg	...	



# Future Directions

- Look at more than accuracy score for model evaluation
- Include data on other types of DNA mutations
- Try other types of classification models