

# RESAMPLING TECHNIQUES

**Professor Ernesto Lee**

# EVALUATE ML ALGORITHMS

# BETTER WAYS TO TRAIN DATA

- ^ Train and Test Sets.
- ^ k-fold Cross-Validation.
- ^ Leave One Out Cross-Validation.
- ^ Repeated Random Test-Train Splits.



# SPLIT INTO TRAIN AND TEST SETS

```
from pandas import read_csv

from sklearn.model_selection import train_test_split

from sklearn.linear_model import LogisticRegression

names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']

dataframe = read_csv('pima-indians-diabetes.data.csv', names=names)

array = dataframe.values

X = array[:,0:8]

Y = array[:,8]

X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.33,
random_state=7)

model = LogisticRegression(solver='liblinear')

model.fit(X_train, Y_train)

result = model.score(X_test, Y_test)

print("Accuracy: %.3f%%" % (result*100.0))
```

# EXERCISE

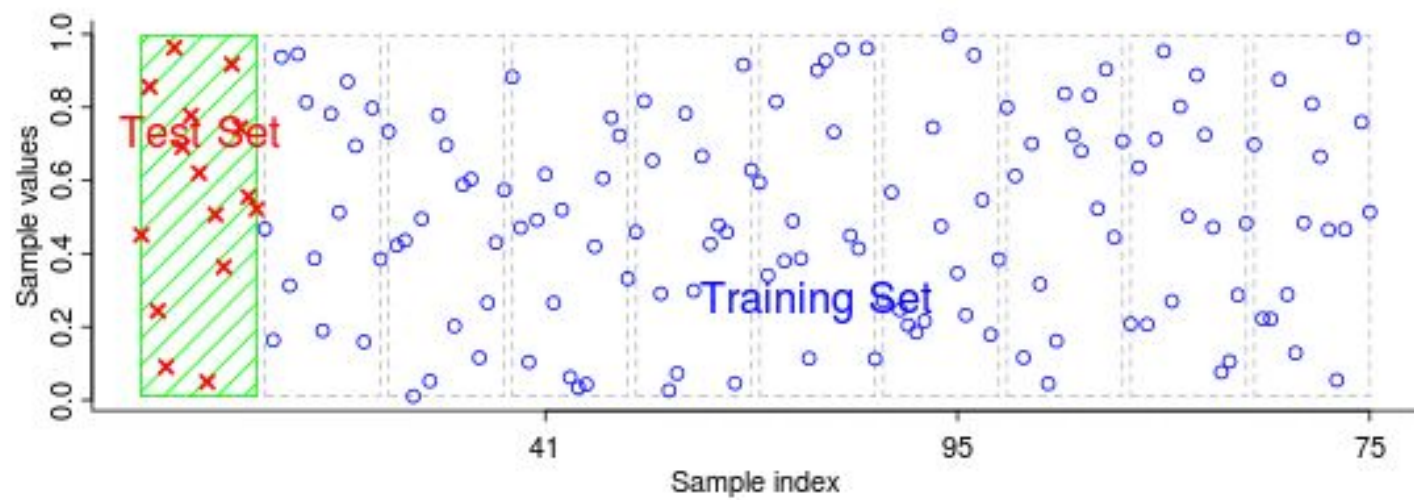
<https://archive-beta.ics.uci.edu/ml/datasets/statlog+german+credit+data>

Download and split the data into a train and test set.

# K-FOLD CROSS VALIDATION

```
# Evaluate using Cross Validation
from pandas import read_csv
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score
from sklearn.linear_model import LogisticRegression

filename = 'pima-indians-diabetes.data.csv'
names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
dataframe = read_csv(filename, names=names)
array = dataframe.values
X = array[:,0:8]
Y = array[:,8]
kfold = KFold(n_splits=10, random_state=7)
model = LogisticRegression(solver='liblinear')
results = cross_val_score(model, X, Y, cv=kfold)
print("Accuracy: %.3f%% (%.3f%%)" % (results.mean()*100.0, results.std()*100.0))
```




# EXERCISE

Use K-Fold Cross Validation to test and train (using Log Regression) on this dataset:

<https://archive-beta.ics.uci.edu/ml/datasets/hepatitis>

**$n = 8$**

	Test		Train
---	------	---	-------

Model 1





# LEAVE ONE OUT CROSS VALIDATION

```
# Evaluate using Leave One Out Cross Validation
from pandas import read_csv
from sklearn.model_selection import LeaveOneOut
from sklearn.model_selection import cross_val_score
from sklearn.linear_model import LogisticRegression
filename = 'pima-indians-diabetes.data.csv'
names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
dataframe = read_csv(filename, names=names)
array = dataframe.values
X = array[:,0:8]
Y = array[:,8]
loocv = LeaveOneOut()
model = LogisticRegression(solver='liblinear')
results = cross_val_score(model, X, Y, cv=loocv)
print("Accuracy: %.3f%% (%.3f%%)" % (results.mean()*100.0, results.std()*100.0))
```

# REPEATED RANDOM TEST-TRAIN SPLITS

# Evaluate using Shuffle Split Cross Validation

```
from pandas import read_csv
```

```
from sklearn.model_selection import ShuffleSplit
```

```
from sklearn.model_selection import cross_val_score
```

```
from sklearn.linear_model import LogisticRegression
```

```
filename = 'pima-indians-diabetes.data.csv'
```

```
names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
```

```
dataframe = read_csv(filename, names=names)
```

```
array = dataframe.values
```

```
X = array[:,0:8]
```

```
Y = array[:,8]
```

```
n_splits = 10
```

```
test_size = 0.33
```

```
seed = 7
```

```
kfold = ShuffleSplit(n_splits=n_splits, test_size=test_size, random_state=seed)
```

```
model = LogisticRegression(solver='liblinear')
```

```
results = cross_val_score(model, X, Y, cv=kfold)
```

```
print("Accuracy: %.3f%% (%.3f%%)" % (results.mean()*100.0, results.std()*100.0))
```

# WHAT TECHNIQUES TO USE WHEN

- K-fold cross validation is the gold standard
- Use train/test for speed with slow algorithms
- Consider the others only when trying to balance variance, performance and speed

# SUMMARY

- ^ Train and Test Sets.
- ^ Cross-Validation.
- ^ Leave One Out Cross-Validation.
- ^ Repeated Random Test-Train Splits.