

"From bits to information"

Introduction to  
Splunk and  
Machine Learning

# Outline

- Introduction to Machine Learning.
- Classification Metrics.
- Numeric and Category Prediction.

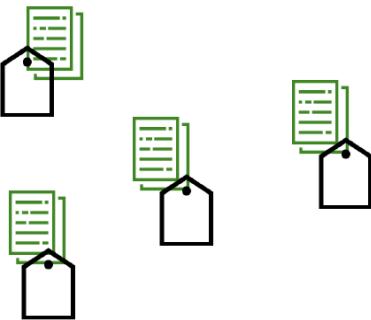


cyber  
& data

# Supervised Learning

## Supervised learning (Classification)

Training set (data with labels)



New instance



Classification

Support Vector Machine (SVM).  
Logistic Regression.

## Supervised learning (Regression)

Training set



Features

Linear Regression.

Decision Trees and Random Forests.

k-Nearest Neighbors (Cluster)

Neural Networks



Value

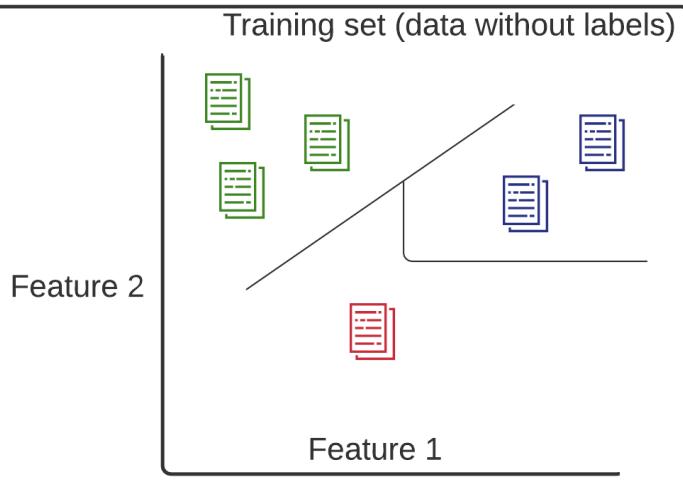
New instance

Attribute: Speed

Feature: Speed + Value

# Unsupervised and Semi-Supervised Learning

## Unsupervised learning



## Semi-supervised

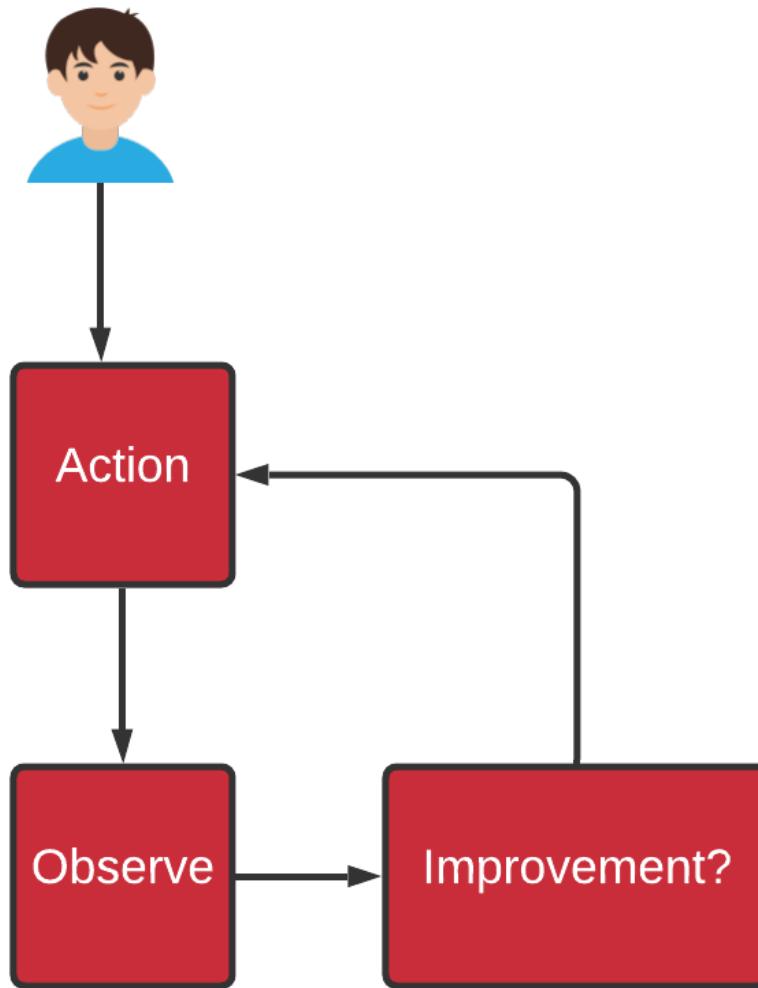


Clustering: k-Means  
Dimensional reduction: Principal Component Analysis (PCA), Kernel PCA.

## Anomaly detection



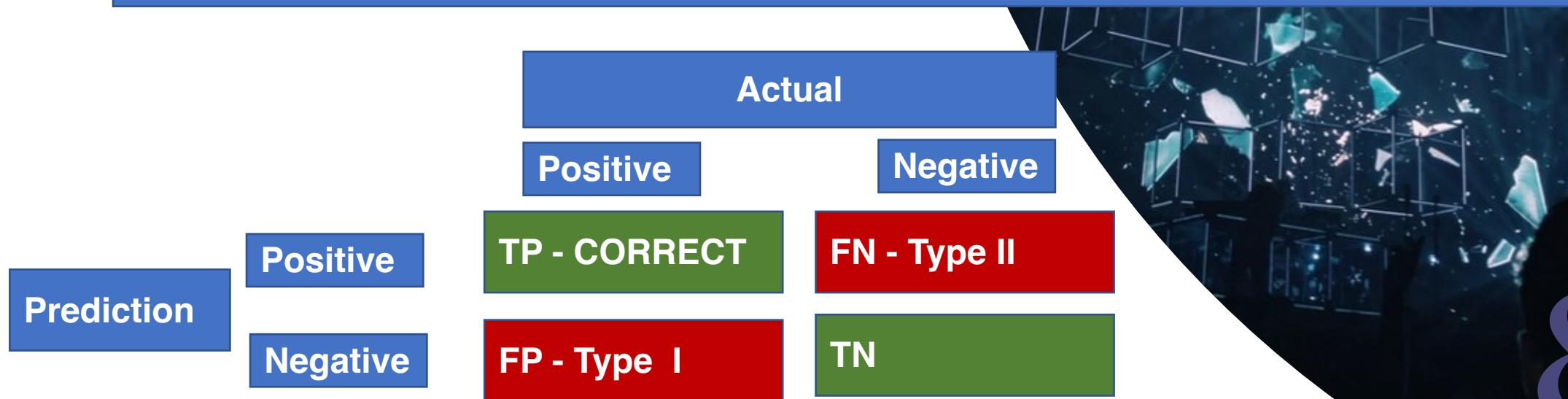
# Re-enforcement learning



cyber  
&  
data

# Classification

- **True-positive (TP)**. That something has been successfully classified as the thing that we want it classified as.
- **False-positive (FP)**. That we have classified something that is incorrect. This would be defined as a **Type I error**. An example might be where a system classes an alert as a hack and where a user enters an incorrect password a number of times, but, on investigation, it is found that the valid user had just forgotten their password.
- **True-negative (TN)**. That we rejected something, and it is not a match.
- **False-negative (FN)**. That we have dismissed something, but, in fact, it is True. This is defined as a miss and is a **Type II error**. With this, a hacker might try a number of passwords for a user, and but where the system does not create an alert for the intrusion.



# Metrics

N=100	Classified as positive	Classified as Negative	Total
Actually positive	10 (TP)	3 (FN)	13
Actually negative	5 (FP)	82 (TN)	87
Total	15	85	100

$$Precision = \frac{TP}{TP + FP} = \frac{10}{10 + 5} = 0.67$$

$$Accuracy = \frac{TP + TN}{total} = \frac{92}{100} = 0.92$$

$$Recall = \frac{FN + TP}{Total} = \frac{10 + 3}{100} = 0.13$$

$$F1 = \frac{precision \times recall}{precision + recall}$$

Precision ↗

Recall ↗

Accuracy ↗

F1 ↗

**0.83**

**0.82**

**0.82**

**0.83**

Classification Results (Confusion Matrix) ↗

Predicted actual ↗	Predicted no ↗	Predicted yes ↗
no	4818 (80.7%)	1152 (19.3%)
yes	1486 (16.4%)	7580 (83.6%)

# Metrics

N=100	Classified as positive	Classified as Negative	Total
Actually positive	10 (TP)	3 (FN)	13
Actually negative	5 (FP)	82 (TN)	87
Total	15	85	100

$$Sensitivity = \frac{TP}{TP + FN} = \frac{10}{13} = 0.77$$

$$Prevalence = \frac{FN + TP}{Total} = \frac{10 + 3}{100} = 0.13$$

$$False\ Positive\ Rate = \frac{FN}{TN + FN} = \frac{5}{82 + 5} = 0.06$$

$$Misclassification\ Rate = \frac{FP + FN}{total} = \frac{18}{100} = 0.18$$

$$Specificity = \frac{TN}{TN + FP} = \frac{82}{82 + 5} = 0.94$$

Precision ↗

Recall ↗

Accuracy ↗

F1 ↗

**0.83**

**0.82**

**0.82**

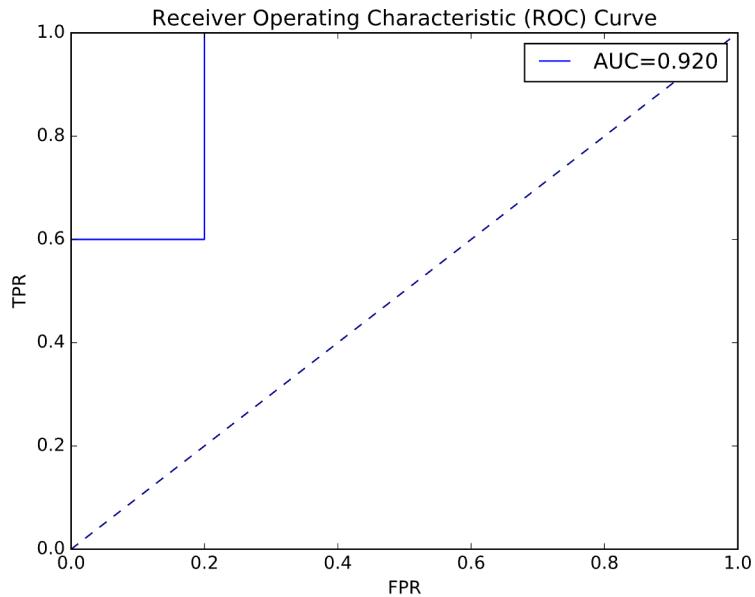
**0.83**

Classification Results (Confusion Matrix) ↗

Predicted actual ↗		Predicted no ↗	Predicted yes ↗
no	yes	4818 (80.7%)	1152 (19.3%)
yes	no	1486 (16.4%)	7580 (83.6%)

# Metrics

```
FPR: [0.  0.  0.  0.2 0.2 1. ]
TPR: [0.  0.2 0.6 0.6 1.  1. ]
Thresholds: [61  60  50  42  39  16]
```



Coding

```
from sklearn import metrics
from sklearn.metrics import roc_auc_score
import matplotlib.pyplot as plt

def show_roc(FPR, TPR, AUC):
    plt.plot(FPR, TPR, color='blue', label='ROC')
    plt.plot([0, 1], [0, 1], color='darkblue', linestyle='--')
    plt.xlabel('FPR')
    plt.ylabel('TPR')
    plt.title('Receiver Operating Characteristic (ROC) Curve')
    plt.legend(["AUC=%3f" % AUC])

plt.show()

y = ['Eve', 'Eve', 'Eve', 'Eve', 'Eve', 'Bob', 'Bob', 'Bob', 'Bob', 'Bob']
scores = [20, 25, 16, 42, 22, 50, 41, 60, 54, 39]

positive_label = 'Bob'

fpr, tpr, thresholds = metrics.roc_curve(y, scores, pos_label=
    positive_label)

auc=metrics.auc(fpr, tpr)

print "FPR:",fpr
print "TPR:",tpr

print "Thresholds:",thresholds

show_roc(fpr, tpr,auc)
```

# Metrics

Confusion matrix:

```
[[87  0  0  0  1  0  0  0  0  0]
 [ 0 88  1  0  0  0  0  0  1  1]
 [ 0  0 85  1  0  0  0  0  0  0]
 [ 0  0  0 79  0  3  0  4  5  0]
 [ 0  0  0  0 88  0  0  0  0  4]
 [ 0  0  0  0  0 88  1  0  0  2]
 [ 0  1  0  0  0  0 90  0  0  0]
 [ 0  0  0  0  1  0 88  0  0  0]
 [ 0  0  0  0  0  0  0 88  0  0]
 [ 0  0  0  1  0  1  0  0  0 90]]
```

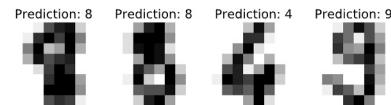
```
Classification report for classifier SVC(C=1.0, break_ties=False, cache_size=20
0, class_weight=None, coef0=0.0,
decision_function_shape='ovr', degree=3, gamma=0.001, kernel='rbf',
max_iter=-1, probability=False, random_state=None, shrinking=True,
tol=0.001, verbose=False):
precision    recall    f1-score   support
```

	precision	recall	f1-score	support
0	1.00	0.99	0.99	88
1	0.99	0.97	0.98	91
2	0.99	0.99	0.99	86
3	0.98	0.87	0.92	91
4	0.99	0.96	0.97	92
5	0.95	0.97	0.96	91
6	0.99	0.99	0.99	91
7	0.96	0.99	0.97	89
8	0.94	1.00	0.97	88
9	0.93	0.98	0.95	92

Tr: 0 Tr: 1 Tr: 2 Tr: 3 Tr: 4 Tr: 5 Tr: 6 Tr: 7 Tr: 8 Tr: 9

0 1 2 3 4 5 6 7 8 9

Prediction: 8 Prediction: 8 Prediction: 4 Prediction: 9



Coding

```
from sklearn import datasets, svm, metrics

digits = datasets.load_digits()

images_and_labels = list(zip(digits.images, digits.target))
for index, (image, label) in enumerate(images_and_labels[:10]):
    plt.subplot(2, 10, (index + 1))
    plt.axis('off')
    plt.imshow(image, cmap=plt.cm.gray_r, interpolation='nearest')
    plt.title('Tr: %i' % label)

# To apply a classifier on this data, we need to flatten the image, to
# turn the data in a (samples, feature) matrix:
n_samples = len(digits.images)
data = digits.images.reshape((n_samples, -1))

# Create a classifier: a support vector classifier
classifier = svm.SVC(gamma=ga)

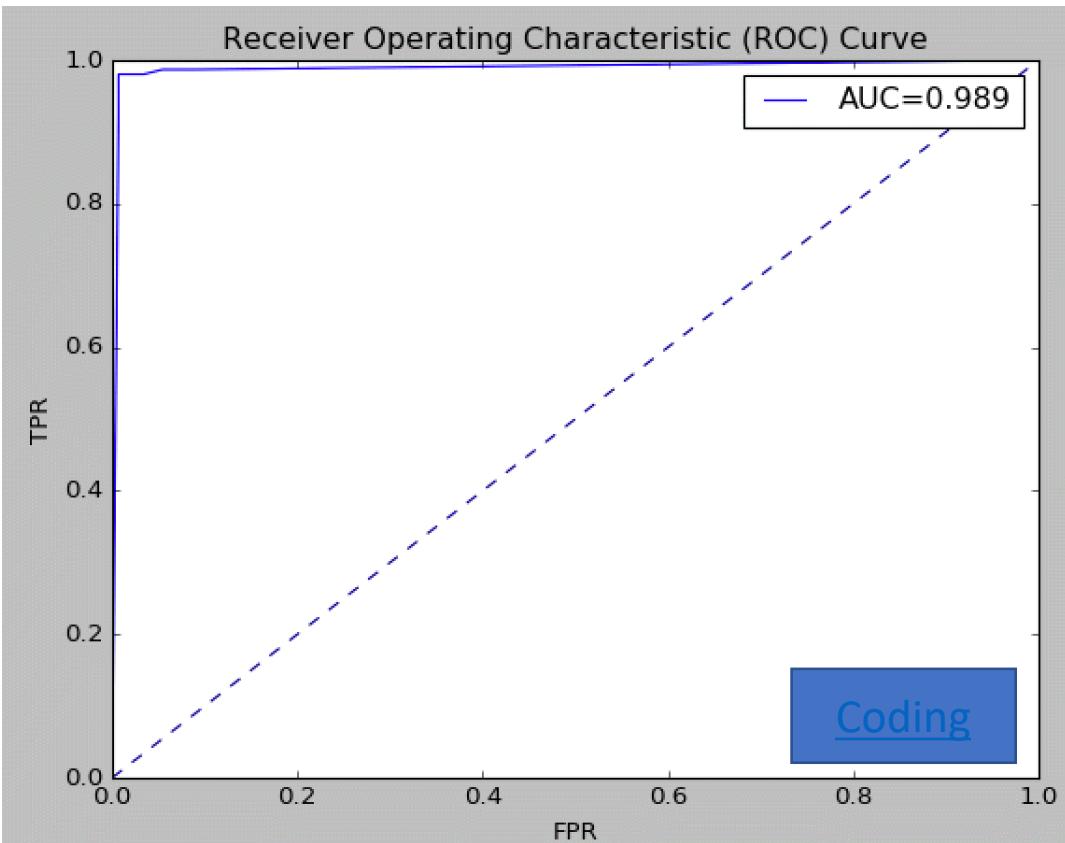
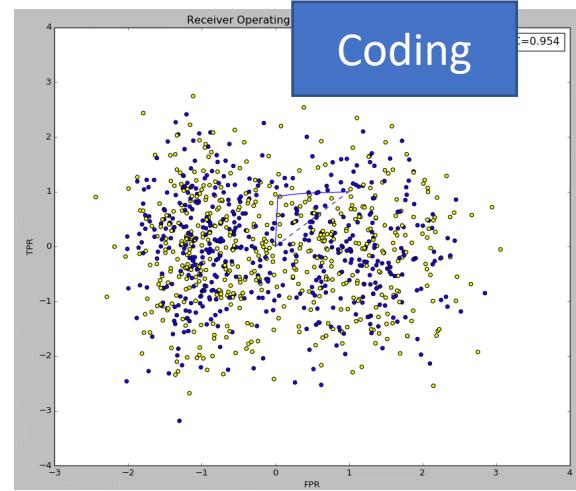
# We learn the digits on the first half of the digits
classifier.fit(data[:n_samples // 2], digits.target[:n_samples // 2])

# Now predict the value of the digit on the second half:
expected = digits.target[n_samples // 2:]
predicted = classifier.predict(data[n_samples // 2:])

print("Classification report for classifier %s:\n%s\n"
      % (classifier, metrics.classification_report(expected, predicted)))
print("Confusion matrix:\n%s" % metrics.confusion_matrix(expected, predicted))

images_and_predictions = list(zip(digits.images[n_samples // 2:], predicted))
for index, (image, prediction) in enumerate(images_and_predictions[:4]):
    plt.subplot(2, 4, index + 5)
    plt.axis('off')
    plt.imshow(image, cmap=plt.cm.gray_r, interpolation='nearest')
    plt.title('Prediction: %i' % prediction)
```

# Classification Metrics



```
from sklearn.datasets import make_classification
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import roc_auc_score
from sklearn.metrics import roc_curve
import matplotlib.pyplot as plt

def show_roc(FPR, TPR, AUC):
    plt.plot(FPR, TPR, color='blue', label='ROC')
    plt.plot([0, 1], [0, 1], color='darkblue', linestyle='--')
    plt.xlabel('FPR')
    plt.ylabel('TPR')
    plt.title('Receiver Operating Characteristic (ROC) Curve')
    plt.legend(["AUC=%3f" % AUC])
    plt.show()

data_vals, class_label = make_classification(n_samples=1000, n_features=2, n_redundant=0,
n_informative=1, n_clusters_per_class=1)
X_train, X_test, y_train, y_test = train_test_split(data_vals, class_label, test_size=0.3,
random_state=1)

model = RandomForestClassifier()

model.fit(X_train, y_train)
print ("Model score: ",model.score(X_test, y_test))

probs = model.predict_proba(X_test)

# probabilities of getting a 1
probs = probs[:, 1]

auc = roc_auc_score(y_test, probs)

FPR, TPR, thresholds = roc_curve(y_test, probs)

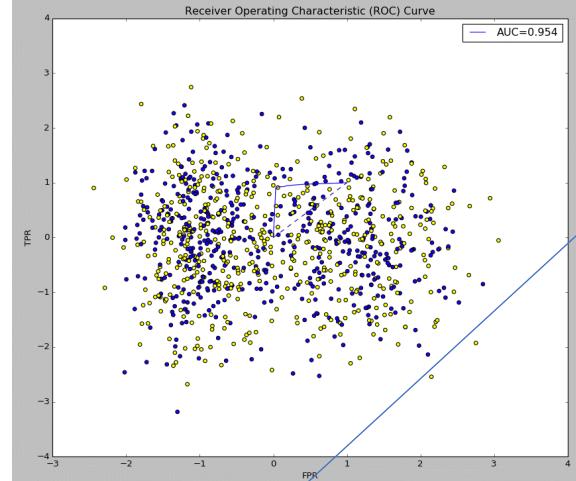
print ("Thresholds: ",thresholds)
print ("FPR: ",FPR)
print ("TPR: ",TPR)

show_roc(FPR, TPR,auc)

colors = ("blue", "yellow")

plt.scatter(data_vals[:, 0], data_vals[:, 1], marker='o', c=colors,s=25, edgecolor='k')
```

# Classification Metrics



metrics.accuracy\_score(). Accuracy classification score.  
metrics.auc() - Area Under the Curve (AUC).  
metrics.confusion\_matrix(). Compute confusion matrix.  
metrics.f1\_score(). Compute the F1 score.  
metrics.precision\_recall\_fscore\_support(). Compute precision, recall, F-measure and support.  
metrics.precision\_score(). Compute the precision.  
metrics.recall\_score(). Compute the recall  
metrics.roc\_auc\_score(). AUC scores.  
metrics.roc\_curve(). Receiver operating characteristic (ROC)

```
from sklearn.datasets import make_classification
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import roc_auc_score
from sklearn.metrics import roc_curve
import matplotlib.pyplot as plt

def show_roc(FPR, TPR, AUC):
    plt.plot(FPR, TPR, color='blue', label='ROC')
    plt.plot([0, 1], [0, 1], color='darkblue', linestyle='--')
    plt.xlabel('FPR')
    plt.ylabel('TPR')
    plt.title('Receiver Operating Characteristic (ROC) Curve')
    plt.legend(["AUC=%3f" % AUC])
    plt.show()

data_vals, class_label = make_classification(n_samples=1000, n_features=2, n_redundant=0,
n_informative=1, n_clusters_per_class=1)
X_train, X_test, y_train, y_test = train_test_split(data_vals, class_label, test_size=0.3,
random_state=1)

model = RandomForestClassifier()

model.fit(X_train, y_train)
print ("Model score: ",model.score(X_test, y_test))

probs = model.predict_proba(X_test)

# probabilities of getting a 1
probs = probs[:, 1]

auc = roc_auc_score(y_test, probs)

FPR, TPR, thresholds = roc_curve(y_test, probs)

print ("Thresholds: ",thresholds)
print ("FPR: ",FPR)
print ("TPR: ",TPR)

show_roc(FPR, TPR,auc)

colors = ("blue", "yellow")

plt.scatter(data_vals[:, 0], data_vals[:, 1], marker='o', c=colors,s=25, edgecolor='k')
```

# Regression Metrics

$R^2$  is a regression score function for numeric predictions:

+ve: The best score is 1.0.

-ve: A model with a negative is worse than arbitrary.

0.0: we can predict the output without the inputs.

The **RMSE** takes the difference between the actual and predicted values and sums them, and then takes the square root of this.

[R<sup>2</sup> Statistic](#)

[Root Mean Squared Error \(RMSE\)](#)

**0.3631**

**13.10**

Coding



```
from sklearn.metrics import r2_score, mean_squared_error, max_error,  
mean_squared_log_error  
bob_login = [48, 12, 7, 11, 43, 44]  
bob_predicted = [41, 14, 9, 15, 40, 41]  
  
print ("R^2 score: ",r2_score(bob_login, bob_predicted))  
print ("RMSE score: ",mean_squared_error(bob_login, bob_predicted))  
print ("Mean squared error: ",mean_squared_log_error(bob_login, bob_predicted))  
  
print ("Max error: ",max_error(bob_login, bob_predicted))
```

# Example (Predict Numeric)

Coding

```
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import r2_score,mean_squared_error

# Features
x1= "blood_pressure"
x2= "age"
# Prediction
x3 = "BMI"

fdata="diabetes.csv"

print ("Training data:\t\t",x1,",",x2)
print ("Training against:\t",x3)
print ("Data set:\t\t",fdata)

ver=pd.read_csv(fdata)

dataset=ver[[x1,x2]]
train=ver[x3]
print (dataset)

x_train, x_test, y_train, y_test= train_test_split(dataset,train,
test_size=0.3, random_state=1)

model= RandomForestRegressor()

model.fit(x_train,y_train)

y_predictions =model.predict(x_test)

accuracy = r2_score(y_test, y_predictions)
mse = mean_squared_error(y_test, y_predictions)

print ("R^2=",accuracy)
print ("MSE=", mse)
```

cyber  
& data

# Example (Category Prediction)

Coding

```
# Features
x1= "blood_pressure"
x2= "age"
# Prediction
x3 = "response"

fdata="diabetes.csv"

print ("Training data:\t\t",x1,",",x2)
print ("Training against:\t",x3)
print ("Data set:\t\t",fdata)

ver=pd.read_csv(fdata)

dataset=ver[[x1,x2]]
train=ver[x3]
print (dataset)

x_train, x_test, y_train, y_test= train_test_split(dataset,train,test_size=0.3,
random_state=1)

model= KMeans(n_clusters=2, random_state=0)

model.fit(x_train,y_train)

y_predictions =model.predict(x_test)

conf=confusion_matrix(y_test,y_predictions)
print (conf)

fpr, tpr, thresholds = roc_curve(y_test,y_predictions)
auc=auc(fpr, tpr)
print ("FPR:",fpr)
print ("TPR:",tpr)
print ("Thresholds:",thresholds)
print ("AUC: ",auc)
```



"From bits to information"

Introduction to  
Splunk and  
Machine Learning

# Outline

---

- Introduction to Machine Learning.
- Introduction to model fitting.



cyber  
& data

# Machine Learning

The two main approaches used within the detection of threats is within:

- **Signature detection**, where we match against well-known patterns of malicious behavior.
- **Use anomaly detection**, where we define a normal behavior pattern.

Within a decision engine, we often use the concept of correct guesses (true) and incorrect ones (false). So, a true positive is where we determine that an event was correctly detected, while a false positive is where a true event was not detected (and thus missed by the system).

Within IDS (Intrusion Detection Systems) there is often a balance to be struck when tuning the systems so that users do not get swamped by too many false alerts (false positives), or from too many fake alerts.



# Machine Learning

- **Information sources.** This involves defining the sources of information that would be required to capture the right information.
- **Data capturing tools.** This involves creating the software agents required to the required data.
- **Data pre-processing.** This involves processing the data into a format which is ready for the analysis part.
- **Feature extraction.** This involves defines the key features that would be required to the analysis engine.
- **Analysis engine.** This involves the creation of an analysis engine which takes the features and creates scoring to evaluate risks.
- **Decision engine.** This takes the scoring systems from the analysis stage and makes a reasoned decision on the level of risk involved.



cyber  
& data

# Splunk and Machine Learning

- fit. Fit a model
- apply. Apply a model run by the fit command.
- summary. Show summary of model.
- listmodels. List the models.
- deletemodel. Delete a model.
- score. Show scores for tests.

```
| inputlookup iris.csv  
| fit GaussianNB petal_length from * into myModel  
| apply myModel as new_petal link.
```

```
| summary myModel link
```

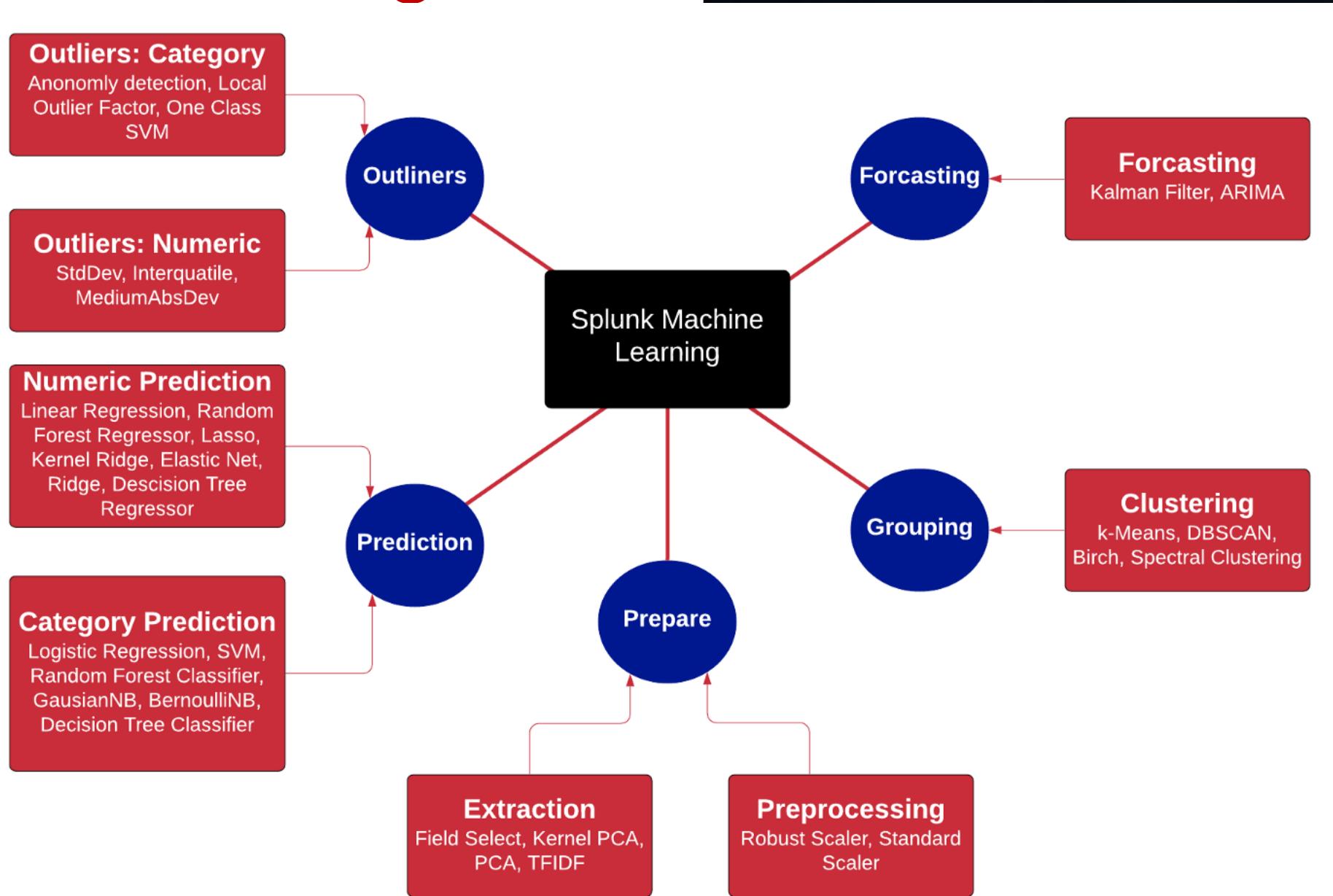
```
| listmodels link
```

A screenshot of the Splunk interface showing search results for "listmodels". The results table displays a single entry:

name	type	options	owner	app
_exp_089ef93fa472494795c64d332410bbb	RandomForestRegressor	(args: ["Class", "Amount", "Time", "V1", "V10", "V11", "V12", "V13", "V14", "V15", "V16", "V17", "V18", "V19", "V2", "V20", "V21", "V22", "V23", "V24", "V25", "V26", "V27", "V28", "V3", "V4", "V5", "V6", "V7", "V8", "V9"], "target_variable": "Class", "feature_variables": ["Amount", "Time", "V1", "V10", "V11", "V12", "V13", "V14", "V15", "V16", "V17", "V18", "V19", "V2", "V20", "V21", "V22", "V23", "V24", "V25", "V26", "V27", "V28", "V3", "V4", "V5", "V6", "V7", "V8", "V9"], "model_name": "exp_089ef93fa472494795c64d332410bbb", "algos_name": "RandomForestRegressor", "mispLimits": {"handle_new_cat": "default", "max_distinct_cat_values": "100", "max_distinct_cat_values_for_classifiers": "100", "max_distinct_cat_values_for_scoring": "100", "max_fit_time": "600", "max_inputs": "100000", "max_memory_usage_mb": "1000", "max_model_size_mb": "15", "max_score_time": "600", "streaming_apply": "false", "use_sampling": "true"}, "kfold_cv": null)	administrator	Splunk

The interface includes a search bar with "listmodels", a statistics tab showing 51 results, and a detailed view of the selected model's configuration.

# Machine Learning



cyber  
&  
data

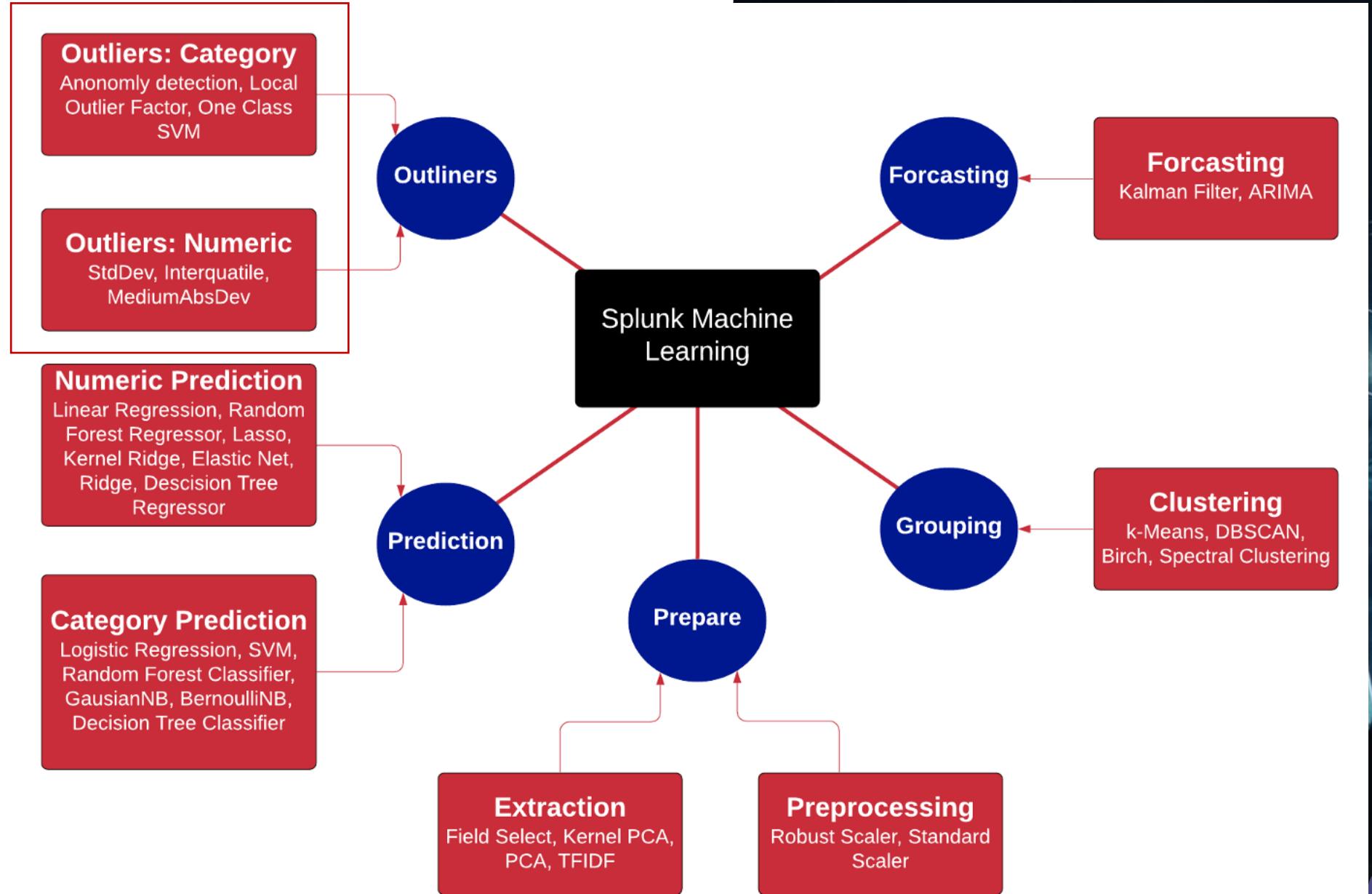
# cyber & data

---

“From bits to information”

fit command

# Outliners



# fit- Anomaly Detection

```
| inputlookup iris.csv  
| fit LocalOutlierFactor petal_length petal_width n_neighbors=10  
algorithm=kd_tree metric=minkowski p=1 contamination=0.14  
leaf_size=10 Link.
```

- *class sklearn.neighbors.LocalOutlierFactor(n\_neighbors=20, \*,  
algorithm='auto', leaf\_size=30, metric='minkowski', p=2,  
metric\_params=None, contamination='auto', novelty=False,  
n\_jobs=None* [Link](#). -> Local Outlier Factor.

```
| inputlookup iris.csv  
| fit OneClassSVM * kernel="poly" nu=0.5 coef0=0.5 gamma=0.5  
tol=1 degree=3 shrinking=f into TESTMODEL_OneClassSVM Link.
```

- *class sklearn.svm.OneClassSVM(\*, kernel='rbf', degree=3,  
gamma='scale', coef0=0.0, tol=0.001, nu=0.5, shrinking=True,  
cache\_size=200, verbose=False, max\_iter=-1)* [Link](#)

```
| inputlookup call_center.csv  
| fit DensityFunction count by "source" into mymodel Link.
```

The dashboard interface includes a header with code snippets, a search bar, and navigation buttons. The main area contains three data tables:

- Top Table:** Statistics (150) view. Columns: sepal\_length, sepal\_width, petal\_length, petal\_width, species, isOutlier, anomaly\_score. Data rows (partial):

sepal_length	sepal_width	petal_length	petal_width	species	isOutlier	anomaly_score
4.3	3.0	1.1	0.1	Iris Setosa	1.0	-2.58
5.8	4.0	1.2	0.2	Iris Setosa	1.0	-1.63
5.7	4.4	1.5	0.4	Iris Setosa	1.0	-1.41
5.4	3.9	1.3	0.4	Iris Setosa	1.0	-1.51
5.7	3.8	1.7	0.3	Iris Setosa	1.0	-1.6
5.4	3.4	1.7	0.2	Iris Setosa	1.0	-1.41
5.1	3.7	1.5	0.4	Iris Setosa	1.0	-1.41
4.6	3.6	1.0	0.2	Iris Setosa	1.0	-2.66
5.1	3.3	1.7	0.5	Iris Setosa	1.0	-1.52
4.8	3.4	1.9	0.2	Iris Setosa	1.0	-2.58

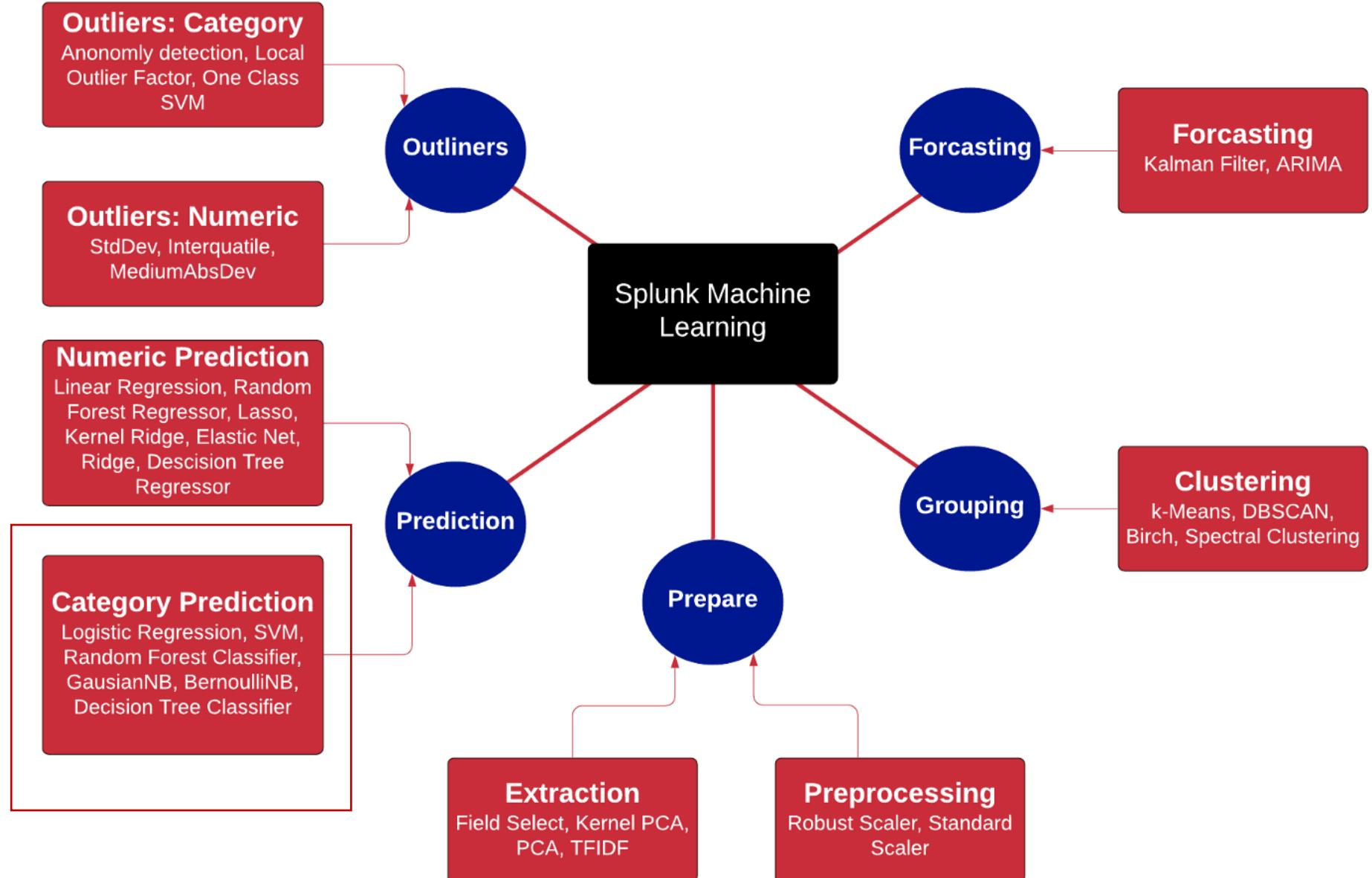
- Middle Table:** Statistics (150) view. Columns: sepal\_length, sepal\_width, petal\_length, petal\_width, species, isNormal. Data rows (partial):

sepal_length	sepal_width	petal_length	petal_width	species	isNormal
5.1	3.5	1.4	0.2	Iris Setosa	-1
4.9	3.0	1.4	0.2	Iris Setosa	-1
4.7	3.2	1.3	0.2	Iris Setosa	-1
4.6	3.1	1.5	0.2	Iris Setosa	-1
5.0	3.6	1.4	0.2	Iris Setosa	-1
5.4	3.9	1.7	0.4	Iris Setosa	-1
4.6	3.4	1.4	0.3	Iris Setosa	-1
5.0	3.4	1.5	0.2	Iris Setosa	-1
4.4	2.9	1.4	0.2	Iris Setosa	-1
4.9	3.1	1.5	0.1	Iris Setosa	-1
5.4	3.7	1.5	0.2	Iris Setosa	-1

- Bottom Table:** Statistics (150) view. Columns: source, \_time, count, IsOutlier(count), BoundaryRanges. Data rows (partial):

source	_time	count	IsOutlier(count)	BoundaryRanges
si_active_agents	2017-09-01 00:00	1	0.0	-Infinity:1.0:0.005 1.0:Infinity:0.005
si_active_agents	2017-09-01 01:00	1	0.0	-Infinity:1.0:0.005 1.0:Infinity:0.005
si_active_agents	2017-09-01 02:00	1	0.0	-Infinity:1.0:0.005 1.0:Infinity:0.005
si_active_agents	2017-09-01 03:00	1	0.0	-Infinity:1.0:0.005 1.0:Infinity:0.005
si_active_agents	2017-09-01 04:00	1	0.0	-Infinity:1.0:0.005 1.0:Infinity:0.005

# Prediction (Cat)



# fit- Prediction

```
| inputlookup iris.csv  
| fit AutoPrediction random_state=42 petal_length from *  
max_features=0.1 into auto_classify_model test_split_ratio=0.3  
random_state=42 Link.
```

```
| inputlookup iris.csv  
| fit BernoulliNB petal_length from * into TESTMODEL_BernoulliNB  
alpha=0.5 binarize=0 fit_prior=f Link.
```

- *class sklearn.naive\_bayes.BernoulliNB(\*, alpha=1.0, binarize=0.0, fit\_prior=True, class\_prior=None)* [Link](#).

```
| inputlookup iris.csv  
| fit DecisionTreeClassifier petal_length from * into sla_MOD Link.
```

```
| inputlookup iris.csv  
| fit GaussianNB petal_length from * into MOD Link.
```

```
| inputlookup iris.csv  
| fit GradientBoostingClassifier petal_length from * into MOD link
```

sepal_length	sepal_width	petal_length	petal_width	species	predicted(petal_length)
5.1	3.5	1.4	0.2	Iris Setosa	1.4200000000000002
4.9	3.0	1.4	0.2	Iris Setosa	1.73
4.7	3.2	1.3	0.2	Iris Setosa	1.3400000000000003
4.6	3.1	1.5	0.2	Iris Setosa	1.4700000000000002
5.0	3.6	1.4	0.2	Iris Setosa	1.4200000000000002
5.4	3.9	1.7	0.4	Iris Setosa	1.6600000000000001
4.6	3.4	1.4	0.3	Iris Setosa	1.44
5.0	3.4	1.5	0.2	Iris Setosa	1.4700000000000002



sepal_length	sepal_width	petal_length	petal_width	species	predicted(petal_length)
5.1	3.5	1.4	0.2	Iris Setosa	1.2
4.9	3.0	1.4	0.2	Iris Setosa	1.2
4.7	3.2	1.3	0.2	Iris Setosa	1.2
4.6	3.1	1.5	0.2	Iris Setosa	1.2
5.0	3.6	1.4	0.2	Iris Setosa	1.2
5.4	3.9	1.7	0.4	Iris Setosa	1.7
4.6	3.4	1.4	0.3	Iris Setosa	1.4
5.0	3.4	1.5	0.2	Iris Setosa	1.2

cyber  
&  
data

# fit- Prediction

| inputlookup iris.csv

| fit **LogisticRegression** petal\_length from \* into MOD [Link](#).

| inputlookup iris.csv

| fit **MLPClassifier** petal\_length from \* into MOD [Link](#).

| inputlookup iris.csv

| fit **RandomForestClassifier** petal\_length from \* into MOD [Link](#).

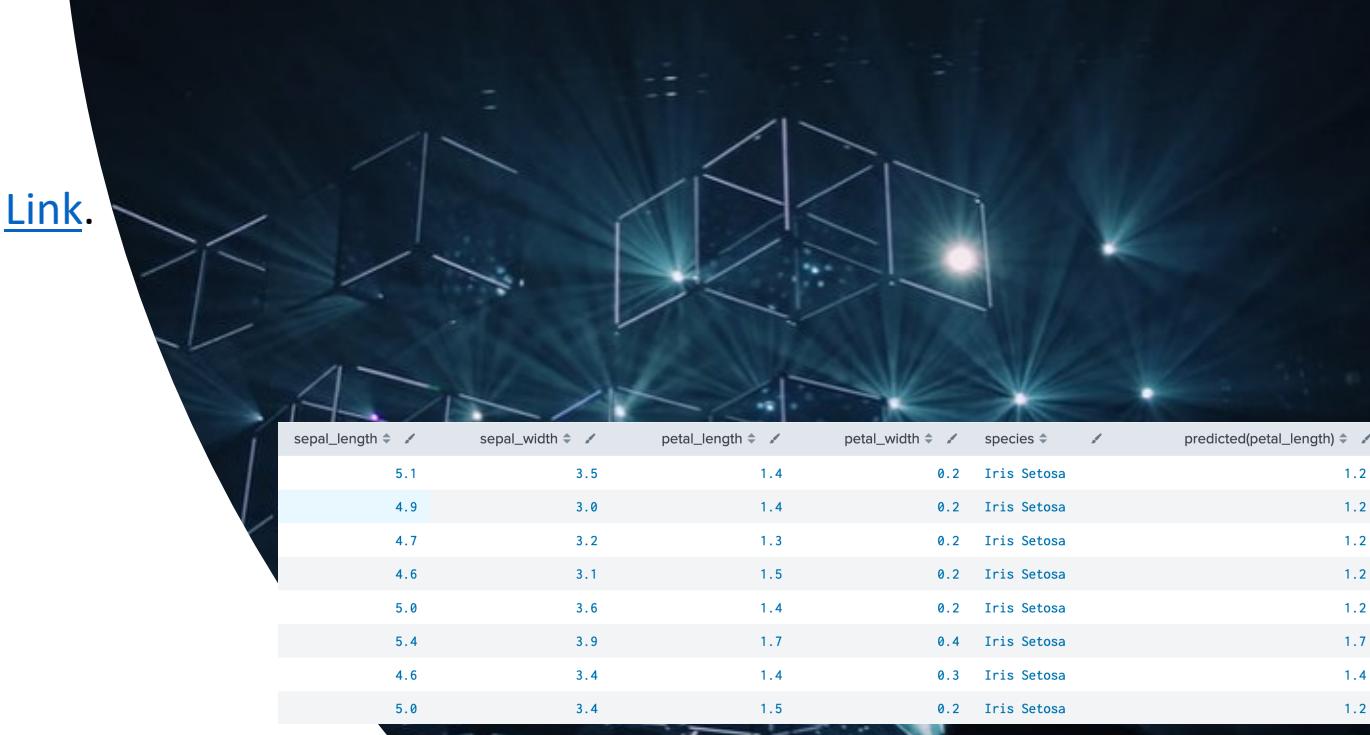
| inputlookup iris.csv

| fit **SGDClassifier** petal\_length from \* into MOD [Link](#).

| inputlookup iris.csv

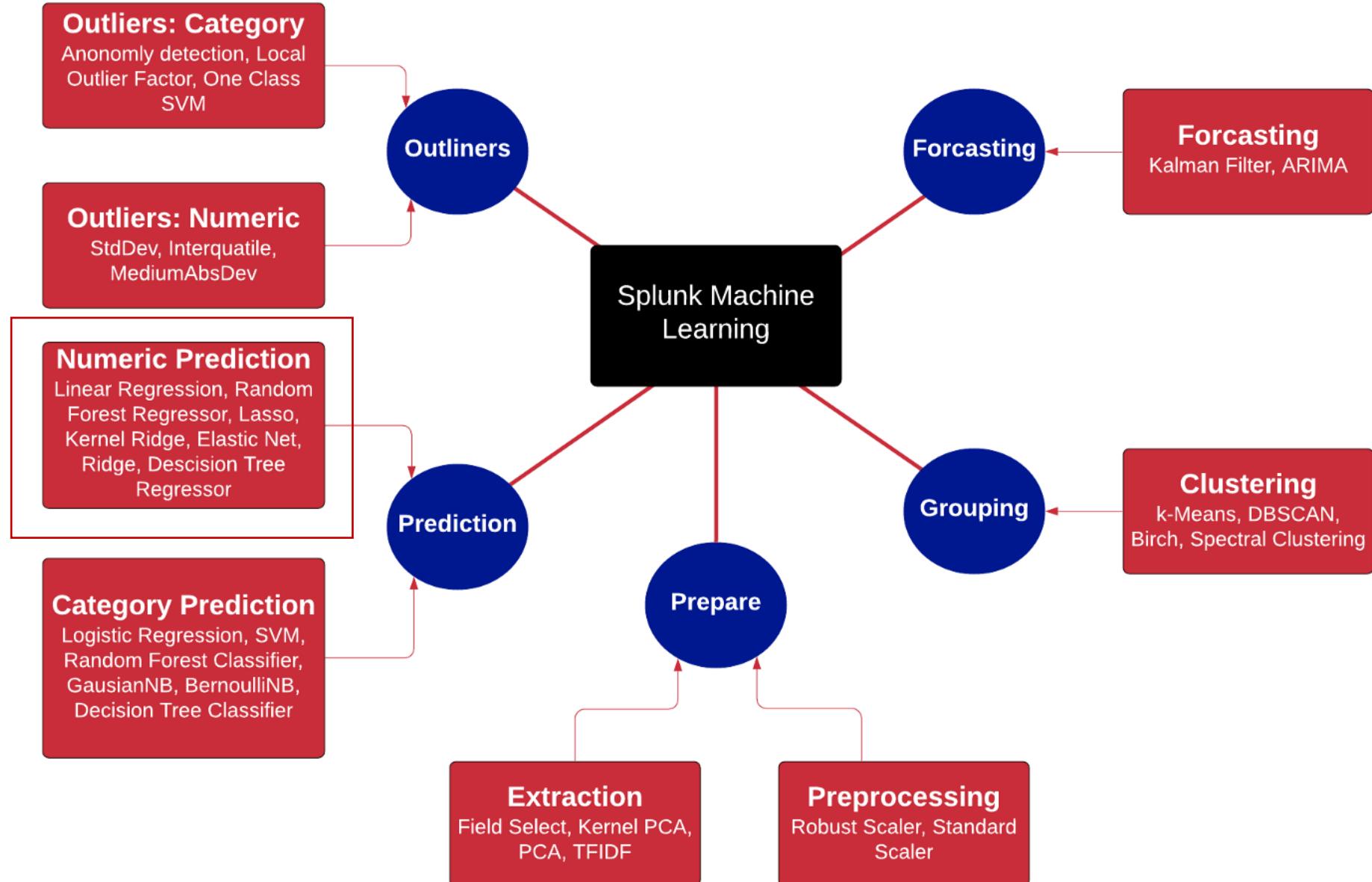
| fit **SVM** petal\_length from \* into MOD. [Link](#).

sepal_length	sepal_width	petal_length	petal_width	species	predicted(petal_length)
5.1	3.5	1.4	0.2	Iris Setosa	1.4200000000000002
4.9	3.0	1.4	0.2	Iris Setosa	1.73
4.7	3.2	1.3	0.2	Iris Setosa	1.3400000000000003
4.6	3.1	1.5	0.2	Iris Setosa	1.4700000000000002
5.0	3.6	1.4	0.2	Iris Setosa	1.4200000000000002
5.4	3.9	1.7	0.4	Iris Setosa	1.6600000000000001
4.6	3.4	1.4	0.3	Iris Setosa	1.44
5.0	3.4	1.5	0.2	Iris Setosa	1.4700000000000002



cyber  
&  
data

# Prediction (Numeric)



# fit- Prediction

```
| inputlookup track_day_missing.csv  
| fit AutoPrediction batteryVoltage target_type=numeric  
test_split_ratio=0.7 from * into PM Link.
```

AutoPrediction can support category or numeric, and then calls **RandomForestRegressor**

```
| inputlookup track_day_missing.csv  
| fit DecisionTreeRegressor batteryVoltage from * into PM Link.
```

```
| inputlookup track_day_missing.csv  
| fit ElasticNet batteryVoltage from * into EN. Link.
```

ElasticNet is a linear regression model and is a generalised form of Lasso and Ridge.

```
| inputlookup track_day_missing.csv  
| fit GradientBoostingRegressor batteryVoltage from * into GB Link.
```

```
| inputlookup track_day_missing.csv  
| fit KernelRidge batteryVoltage from * into KR Link.
```

vehicleType	batteryVoltage	engineCoolantTemperature	engineSpeed	lateralGForce	longitudeGForce	speed	verticalGForce	predicted(batteryVoltage)
2015 Porsche GT3		93.0	6060	1.11	0.5	69	-2.0	13.90133565917362
2013 Audi RS5	13.937000000000001	94.0	4957	0.56	0.7	56	0.95	13.964810194805192
2015 Porsche GT3	13.827	93.0	6163	0.71	0.26	70	-2.0	13.893069210760919
2011 Ford Mustang GT500	14.035	87.0	2846	0.81	-0.71	47	-2.0	14.038379523809521
2015 Porsche GT3	13.827	93.0	6542	0.49	-0.18	76	-2.0	13.965466269841269
2014 Chevrolet Corvette	14.624	105.0	4425	0.32	0.05	100	-0.17	14.70799623015873
2015 Porsche GT3	13.827	93.0	7763	0.24	-0.18	91	-2.0	14.12003358974359



cyber  
&  
data

# fit- Prediction

| inputlookup track\_day\_missing.csv

| fit **Lasso** batteryVoltage from \* into LA [Link](#).

| inputlookup track\_day\_missing.csv

| fit **LinearRegression** batteryVoltage from \* into LR [Link](#).

| inputlookup track\_day\_missing.csv

| fit **RandomForestRegressor** batteryVoltage

min\_samples\_split=30000 from \* into RF [Link](#).

| inputlookup track\_day\_missing.csv

| fit **Ridge** batteryVoltage from \* into RD [Link](#).

| inputlookup track\_day\_missing.csv

| fit **SGDRegressor** batteryVoltage from \* into SG [Link](#).

| inputlookup app\_usage.csv

| fit **SystemIdentification** Expenses from HR1 HR2 ERP dynamics=3-2-2-3 layers=64-64-64 [Link](#).

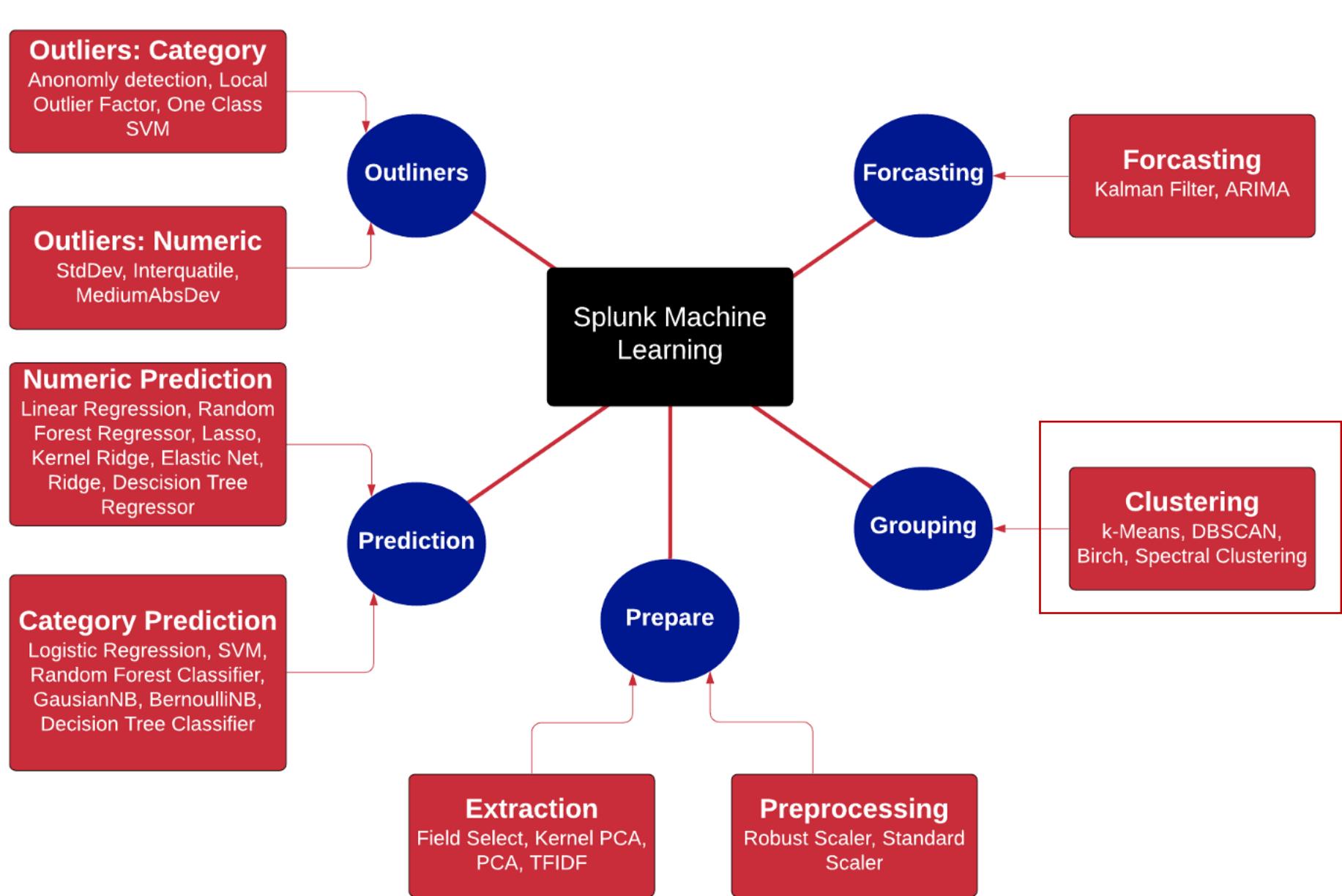
vehicleType	batteryVoltage	engineCoolantTemperature	engineSpeed	lateralGForce	longitudeGForce	speed	verticalGForce	predicted(batteryVoltage)
2015 Porsche GT3		93.0	6060	1.11	0.5	69	-2.0	14.047398311500638
2013 Audi RS5	13.937000000000001	94.0	4957	0.56	0.7	56	0.95	14.050742114165605
2015 Porsche GT3	13.827	93.0	6163	0.71	0.26	70	-2.0	14.047891131013976
2011 Ford Mustang GT500	14.035	87.0	2846	0.81	-0.71	47	-2.0	13.841659798088022
2015 Porsche GT3	13.827	93.0	6542	0.49	-0.18	76	-2.0	14.050547588476304
2014 Chevrolet Corvette	14.624	105.0	4425	0.32	0.05	100	-0.17	14.723515281942252
2015 Porsche GT3	13.827	93.0	7763	0.24	-0.18	91	-2.0	14.052966400552615
2015 Porsche GT3	13.827	93.0	6365	0.38	-0.2	95	-2.0	14.06058143488513
2015 Porsche GT3	13.827	93.0	6713	0.04	0.13	100	-2.0	14.06037488719077



Uses multi-layer neural network. In this case five layers, with three hidden layers of 64/64/64. Logs are 3 for Expenses, 2 for H1, and so on.

cyber  
&  
data

# Cluster



# fit- Cluster

- | inputlookup iris.csv
- | fit **Birch** petal\_length k=3 partial\_fit=true into MOD [Link](#)
  
- | inputlookup iris.csv
- | fit **DBSCAN** petal\_length min\_samples=4 [link](#)
  
- | inputlookup iris.csv
- | fit **GMeans** petal\_length random\_state=42 into MOD3 [link](#).  
[based on k-means]
  
- | inputlookup iris.csv
- | fit **KMeans** petal\_length k=3 into MOD4 [link](#).
  
- | inputlookup iris.csv
- | fit **SpectralClustering** petal\_length k=3 [link](#).
  
- | inputlookup iris.csv
- | fit **XMeans** petal\_length [link](#).

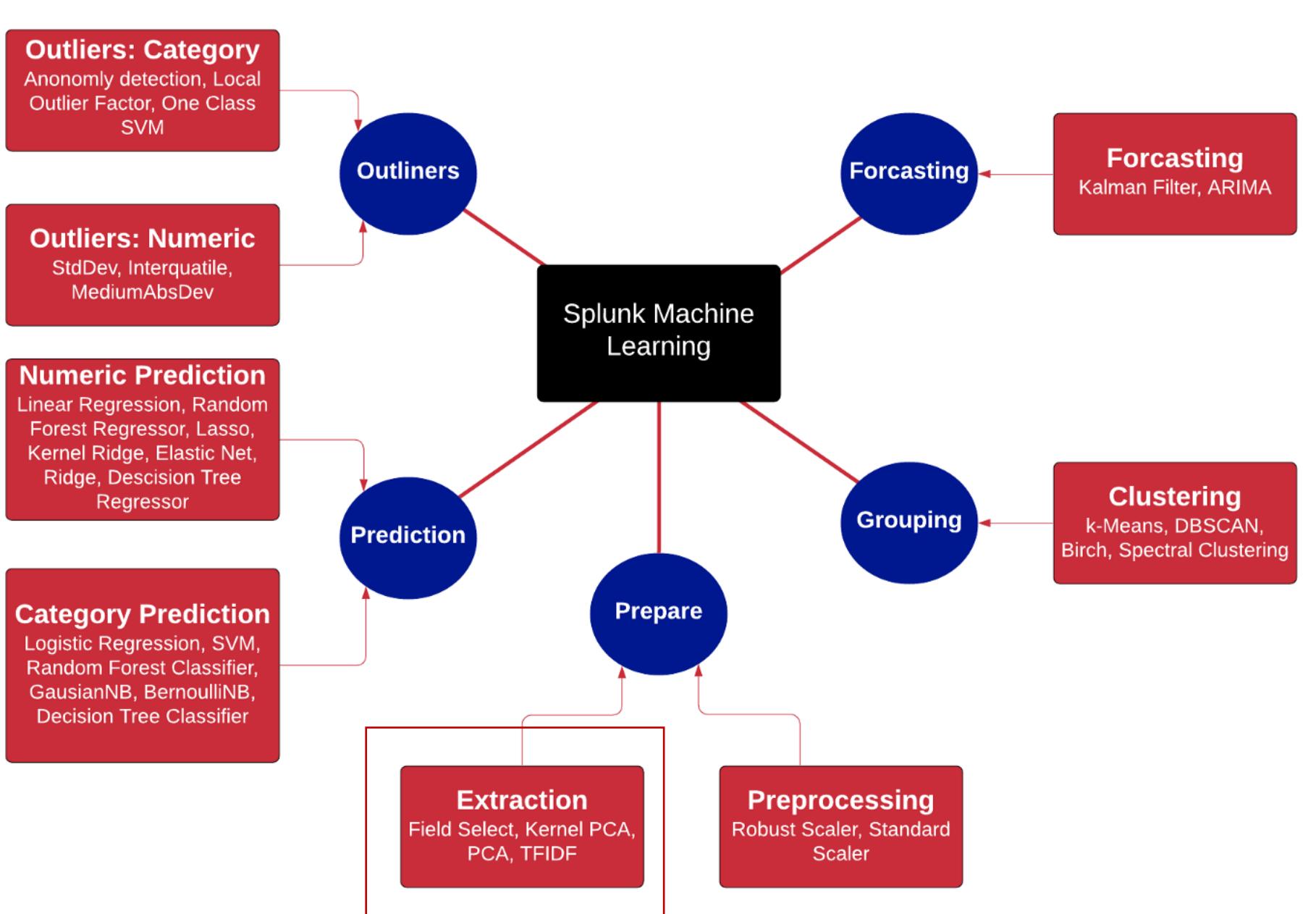
sepal_length	sepal_width	petal_length	petal_width	species	cluster
7.0	3.2	4.7	1.4	Iris Versicolor	0
6.4	3.2	4.5	1.5	Iris Versicolor	0
6.9	3.1	4.9	1.5	Iris Versicolor	0
5.5	2.3	4.0	1.3	Iris Versicolor	0
6.5	2.8	4.6	1.5	Iris Versicolor	0
5.7	2.8	4.5	1.3	Iris Versicolor	0
6.3	3.3	4.7	1.6	Iris Versicolor	0
4.9	2.4	3.3	1.0	Iris Versicolor	0



sepal_length	sepal_width	petal_length	petal_width	species	cluster	cluster_distance
5.1	3.5	1.4	0.2	Iris Setosa	1	0.00409600000000121
4.9	3.0	1.4	0.2	Iris Setosa	1	0.00409600000000121
4.7	3.2	1.3	0.2	Iris Setosa	1	0.02689600000000267
4.6	3.1	1.5	0.2	Iris Setosa	1	0.001295999999999383
5.0	3.6	1.4	0.2	Iris Setosa	1	0.00409600000000121
5.4	3.9	1.7	0.4	Iris Setosa	1	0.05569599999999957
4.6	3.4	1.4	0.3	Iris Setosa	1	0.00409600000000121

er  
data

# Feature Extraction



# fit- Feature Extraction

| inputlookup track\_day.csv

| fit **FieldSelector** batteryVoltage from engineCoolantTemperature, engineSpeed, lateralGForce,longitudeGForce, speed **type=numeric** [Link](#)

| inputlookup track\_day.csv

| fit **FieldSelector** vehicleType from engineCoolantTemperature, engineSpeed, lateralGForce,longitudeGForce, speed **type=categorical** [link](#)

| inputlookup passwords.csv

| fit **HashingVectorizer** Passwords ngram\_range=1-2 k=10 [Link](#).

| inputlookup track\_day.csv

| fit **ICA** batteryVoltage, engineSpeed n\_components=2 as IC [Link](#)

| inputlookup track\_day.csv

| fit **KernelPCA** batteryVoltage, engineSpeed k=3 gamma=0.001 as PCA

[Link](#)



vehicleType	batteryVoltage	engineCoolantTemperature	engineSpeed	lateralGForce	longitudeGForce	speed	verticalGForce
2015 Porsche GT3	13.785	93.0	6060	1.11	0.5	69	-2.0
2013 Audi RS5	13.93700000000001		94.0	4957	0.56	0.7	56
2015 Porsche GT3	13.827	93.0	6163	0.71	0.26	70	-2.0
2011 Ford Mustang GT500	14.035	87.0	2846	0.81	-0.71	47	-2.0

vehicleType	batteryVoltage	engineCoolantTemperature	engineSpeed	lateralGForce	longitudeGForce	speed	verticalGForce	fs_lateralGForce
2015 Porsche GT3	13.785	93.0	6060	1.11	0.5	69	-2.0	1.11
2013 Audi RS5	13.93700000000001	94.0	4957	0.56	0.7	56	0.95	0.56
2015 Porsche GT3	13.827	93.0	6163	0.71	0.26	70	-2.0	0.71
2011 Ford Mustang GT500	14.035	87.0	2846	0.81	-0.71	47	-2.0	0.81

Passwords	Passwords_hashed_0	Passwords_hashed_1	Passwords_hashed_2	Passwords_hashed_3	Passwords_hashed_4	Passwords_h
qwerty123	-0.0	0.9999999999999998	0.0	-0.0	-0.0	-0.0
qwerty5	-1.4719616800160393e-17	0.0	0.015088875493768687	0.4224794055919135	0.7253279855866746	-0.1059026953
cisco1	1.0000000000000002	2.2187595626058292e-16	1.0686969971588345e-17	1.2583791481697097e-17	4.44133817001074e-18	-8.1242533116
cisco43	-9.941516562259971e-35	-3.760597526072945e-19	-0.18988975014269943	0.000903986945728833	0.10829789914398458	0.9752121
liverpool	-3.189250306701419e-17	1.9421073798334343e-17	0.48208015086067485	0.34827807416123513	-0.586295299832911	0.17695195
liverp001	1.4719616800160396e-17	-1.583776923683663e-17	0.7867056762806923	0.13091943941483192	0.29857321546377086	0.10184005
cisco1	1.0	2.2187595626064292e-16	7.027787639449764e-18	-1.6066250926232365e-17	-1.0737492244786327e-17	3.56693015115
cisc01	2.698596413362739e-17	-9.76739240727416e-18	-0.3352717042653099	0.8264821946661363	-0.17112198797126701	-0.06159265
qwerty123	-4.43751912521258e-16	1.0	0.0	-0.0	0.0	0.0

vehicleType	batteryVoltage	engineCoolantTemperature	engineSpeed	lateralGForce	longitudeGForce	speed	verticalGForce
2015 Porsche GT3	13.785	93.0	6060	1.11	0.5	69	-2.0
2013 Audi RS5	13.93700000000001	94.0	4957	0.56	0.7	56	-2.0
2015 Porsche GT3	13.827	93.0	6163	0.71	0.26	70	-2.0
2011 Ford Mustang GT500	14.035	87.0	2846	0.81	-0.71	47	-2.0
2015 Porsche GT3	13.716	93.0	5798	0.21	-0.4	68	-2.0

# fit- Feature Extraction

| inputlookup track\_day.csv

| fit **NPR** vehicleType from engineSpeed as npr01 [Link](#)

Normalized Perlich Ratio (NPR) algorithm categorical field values into numeric field entries.

| inputlookup track\_day.csv

| fit **PCA** engineCoolantTemperature, engineSpeed, lateralGForce,speed k=3 as pca01 [Link](#)

Principal Component Analysis (PCA) reduce the number of fields in the data.

| inputlookup track\_day.csv

| fit **TFIDF** vehicleType ngram\_range=1-2 max\_df=0.6 min\_df=0.2 stop\_words=english as tf01 [Link](#)

TFIDF converts raw text data into a matrix.

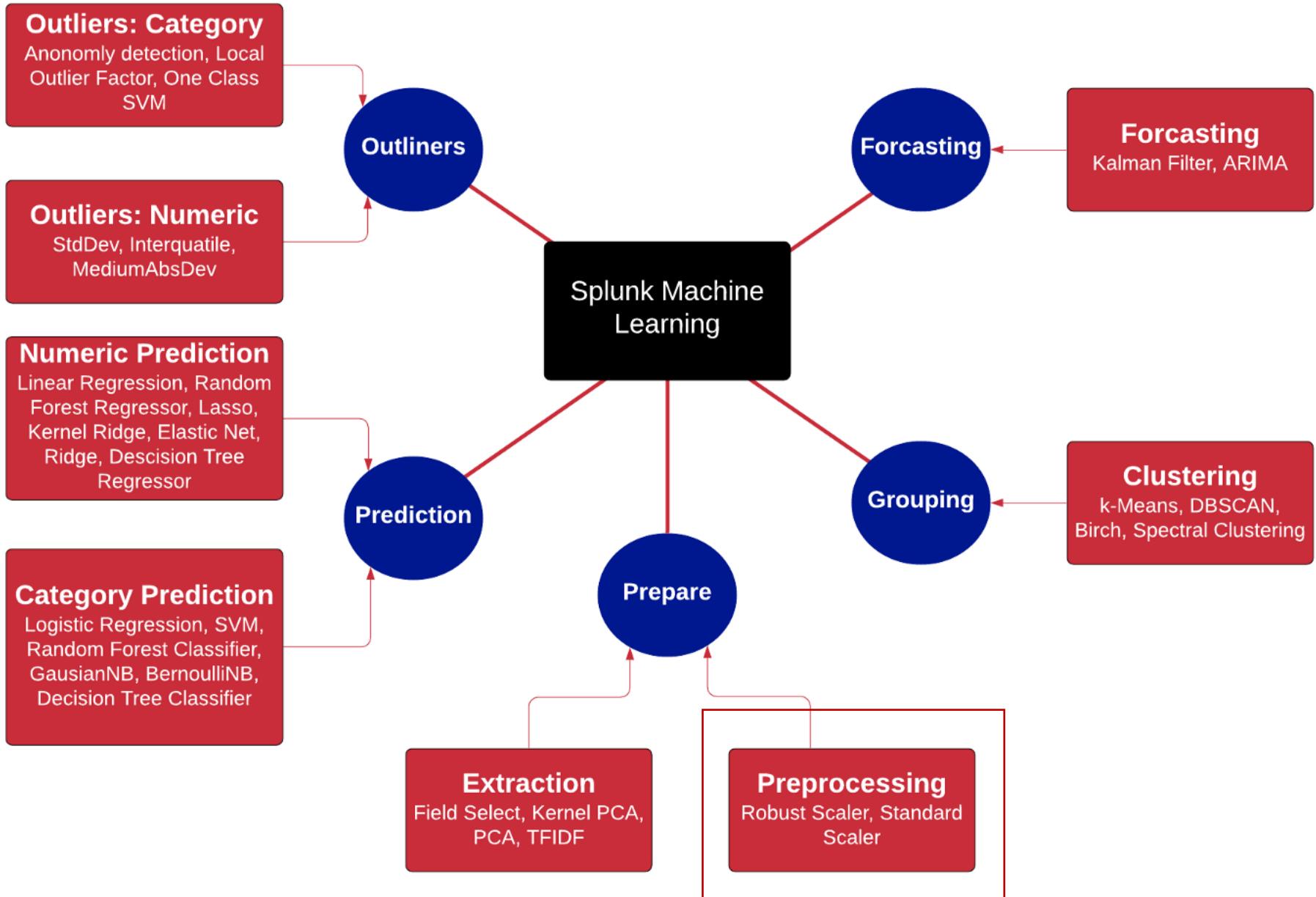
batteryVoltage	engineCoolantTemperature	engineSpeed	lateralGForce	longitudeGForce	speed	verticalGForce	NPR_engineSpeed_2008	NPR_engineSpeed_2011
13.785	93.0	6060	1.11	0.5	69	-2.0	0.20341933065483636	0.0
13.937000000000001	94.0	4957	0.56	0.7	56	0.95	0.4750875723470781	0.0
13.827	93.0	6163	0.71	0.26	70	-2.0	0.15948805818592193	0.0
14.035	87.0	2846	0.81	-0.71	47	-2.0	0.0	0.5529840905437634
13.827	93.0	6542	0.49	-0.18	76	-2.0	0.0	0.0

vehicleType	batteryVoltage	engineCoolantTemperature	engineSpeed	lateralGForce	longitudeGForce	speed	verticalGForce	pca01_1	pca01_2
2015 Porsche GT3	13.785	93.0	6060	1.11	0.5	69	-2.0	1973.31538399881	-3.8336505338179956
2013 Audi RS5	13.937000000000001	94.0	4957	0.56	0.7	56	0.95	870.2387264963066	-3.305800082392769
2015 Porsche GT3	13.827	93.0	6163	0.71	0.26	70	-2.0	2076.3197896193524	-3.893886305786853
2011 Ford Mustang GT500	14.035	87.0	2846	0.81	-0.71	47	-2.0	-1240.6992929015553	0.3620517152381604
2015 Porsche GT3	13.827	93.0	6542	0.49	-0.18	76	-2.0	2455.3643752526996	-4.943621239329986



cyber  
&  
data

# Preprocessing



cyber  
&  
data

# fit- Preprocessing

| inputlookup track\_day\_missing.csv

| fit Imputer batteryVoltage [Link](#)

vehicleType	batteryVoltage	engineCoolantTemperature	engineSpeed	lateralGForce	longitudeGForce	speed	verticalGForce	Imputed_batteryVoltage
2015 Porsche GT3		93.0	6060	1.11	0.5	69	-2.0	14.113871297425948
2013 Audi RS5	13.937000000000001	94.0	4957	0.56	0.7	56	0.95	13.937000000000001
2015 Porsche GT3	13.827	93.0	6163	0.71	0.26	70	-2.0	13.827
2011 Ford Mustang GT500	14.035	87.0	2846	0.81	-0.71	47	-2.0	14.035
2015 Porsche GT3	13.827	93.0	6542	0.49	-0.18	76	-2.0	13.827
2014 Chevrolet Corvette	14.624	105.0	4425	0.32	0.05	100	-0.17	14.624
2015 Porsche GT3	13.827	93.0	7763	0.24	-0.18	91	-2.0	13.827

Imputer fills in blanks for medium, or specific value.

| inputlookup track\_day\_missing.csv

| fit RobustScaler \* [Link](#)

vehicleType	batteryVoltage	engineCoolantTemperature	engineSpeed	lateralGForce	longitudeGForce	speed	verticalGForce	RS_batteryVoltage
2015 Porsche GT3		93.0	6060	1.11	0.5	69	-2.0	
2013 Audi RS5	13.937000000000001	94.0	4957	0.56	0.7	56	0.95	-0.1406926406926397
2015 Porsche GT3	13.827	93.0	6163	0.71	0.26	70	-2.0	-0.37878787878788056
2011 Ford Mustang GT500	14.035	87.0	2846	0.81	-0.71	47	-2.0	0.07142857142857033
2015 Porsche GT3	13.827	93.0	6542	0.49	-0.18	76	-2.0	-0.37878787878788056

RobustScaler scales to median and interquartile range to 0 and 1. Avoids dominance of fields with large values.

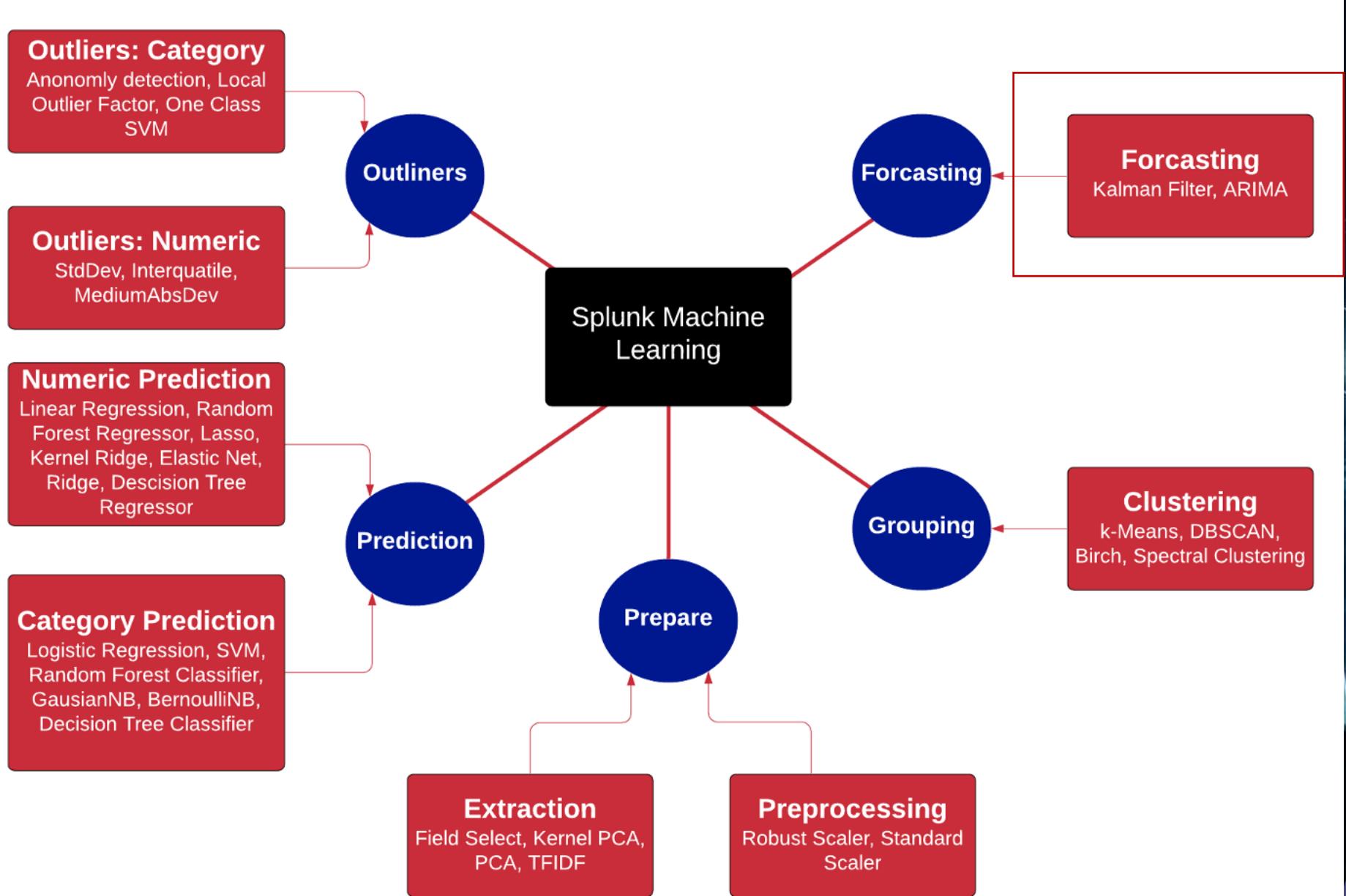
| inputlookup track\_day\_missing.csv

| fit StandardScaler \* [Link](#)

vehicleType	batteryVoltage	engineCoolantTemperature	engineSpeed	lateralGForce	longitudeGForce	speed	verticalGForce	SS_batteryVoltage	SS_engineCoolantTemperature
2015 Porsche GT3		93.0	6060	1.11	0.5	69	-2.0		
2013 Audi RS5	13.937000000000001	94.0	4957	0.56	0.7	56	0.95	-0.4981803519277019	0.324647341076498
2015 Porsche GT3	13.827	93.0	6163	0.71	0.26	70	-2.0	-0.803465525894621	0.27818024136906
2011 Ford Mustang GT500	14.035	87.0	2846	0.81	-0.71	47	-2.0	-0.22619901512081622	-0.00062235687556058
2015 Porsche	13.827	93.0	6542	0.49	-0.18	76	-2.0	-0.803465525894621	0.27818024136906

StandardScalar scales to median and interquartile range to 0 and 1. Avoids dominance of fields with large values.

# Preprocessing



cyber  
&  
data

# Forecasting

```
| inputlookup app_usage.csv
```

```
| fit StateSpaceForecast CRM ERP Expenses holdback=12  
into SF Link.
```

RemoteAccess	Webmail	_time	lower95(predicted(CRM))	lower95(predicted(ERP))	lower95(predicted(Expenses))	predicted(CRM)	predicted(ERP)
897	760	2015-07-28				787.6566719248075	245.06352282287813
938	904	2015-07-27				760.523128268509	243.80372342031262
859	736	2015-07-21				756.5699826494289	240.70924134242486
853	636	2015-08-04				748.0078713245281	231.7229700844895
828	804	2015-07-29				744.0975823377198	241.92535659112744
902	725	2015-08-11				736.8298191126142	226.07432574518296
953	680	2015-08-19				733.5698420631398	219.37842803297875

StateSpaceForecast based on Kalman filters.

```
| inputlookup logins.csv
```

```
| fit ARIMA _time logins holdback=0 conf_interval=95  
order=0-0-0 forecast_k=5 as AR link.
```



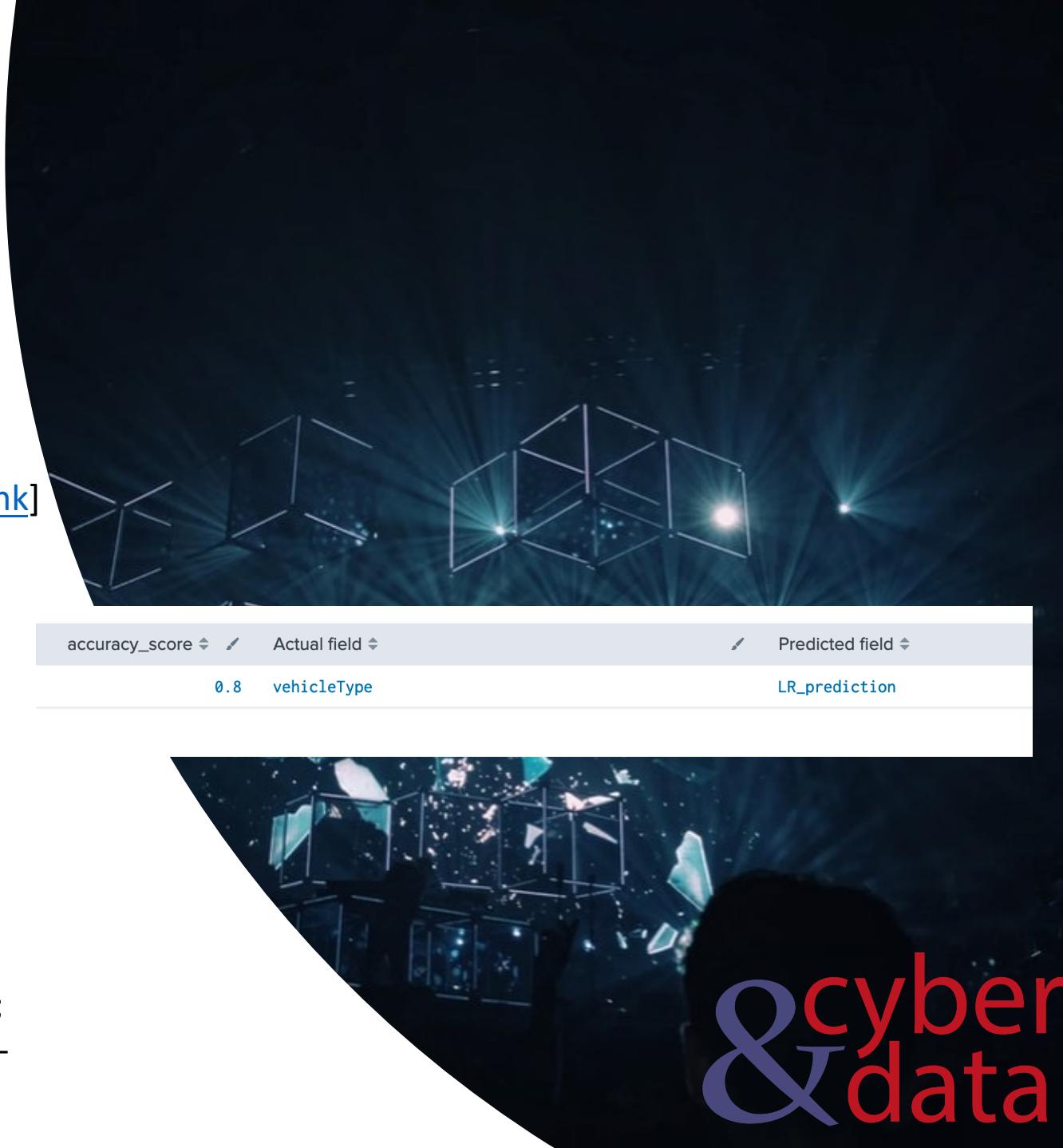
cyber  
&  
data

# Scoring

- **score.** Show scores for tests.

```
| inputlookup track_day.csv  
| head 100  
| fit LogisticRegression vehicleType from batteryVoltage  
engineCoolantTemperature engineSpeed into LR_model  
| apply LR_model as LR_prediction  
| score accuracy_score vehicleType against LR_prediction [Link]
```

- **Classification:** Accuracy, Confusion matrix, F1-score, Precision, Precision-Recall-F1-Support, Recall, ROC-AUC-score, ROC-curve.
- **Clustering:** Silhouette score.
- **Pairwise distances scoring:** Pairwise distances score.
- **Regression scoring:** Explained variance score; Mean absolute error score; Mean squared error; and R2 score.
- **Statistical functions** (statsfunctions): Describe; Moment, Pearson; Spearman; Tmean; Trim, and Tvar.
- **Statistical testing** (statstest): Analysis of Variance (Anova); Augmented Dickey-Fuller (Adfuller); Energy distance; One-way ANOVA; and T-test (1 sample).





"From bits to information"

Introduction to  
Splunk and  
Machine Learning