# AREIX

*Everyone Worthy, Everyone Wealthy*
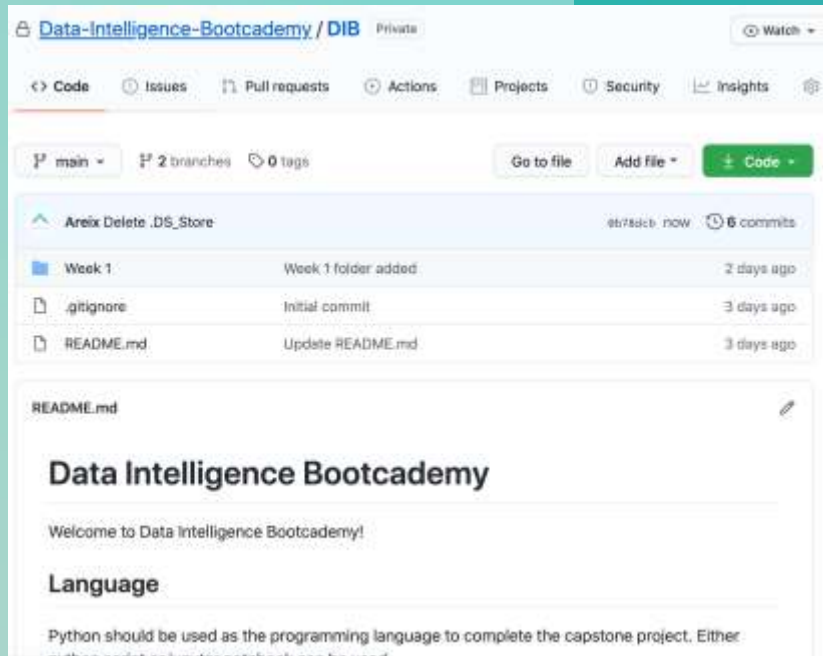
# Modelling in Practice

*Charon Guo*
*HEAD OF APPLICATION &*
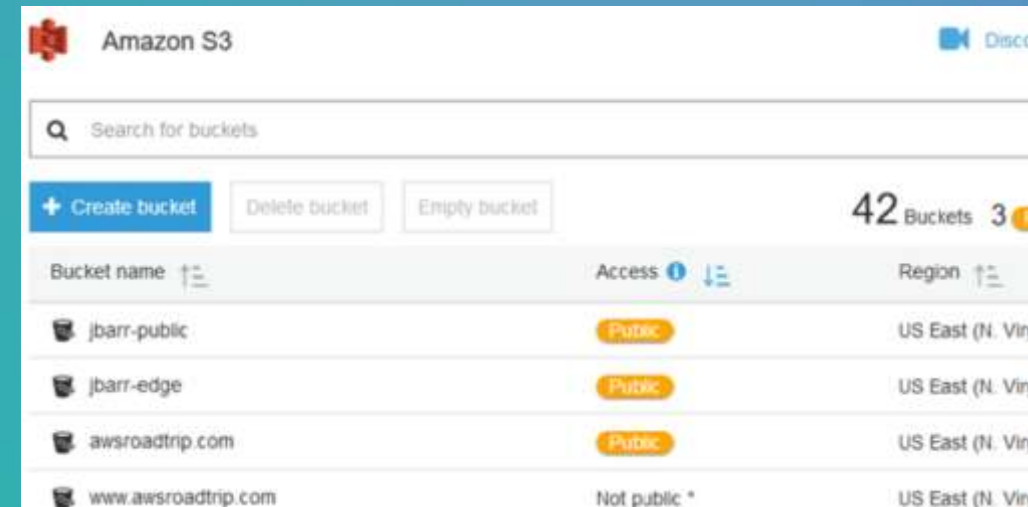*BACKEND DEVELOPMENT*

# Tools

**AREIX**

## Github



## AWS DynamoDB



## AWS S3

# What is Machine Learning?

A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T, as measured by P, improves with experience E.

--- Mitchell, T. (1997). Machine Learning, McGraw Hill

Machine learning is a technique of data science that helps computers **learn from existing data** in order to **forecast future behaviors, outcomes, and trends**.
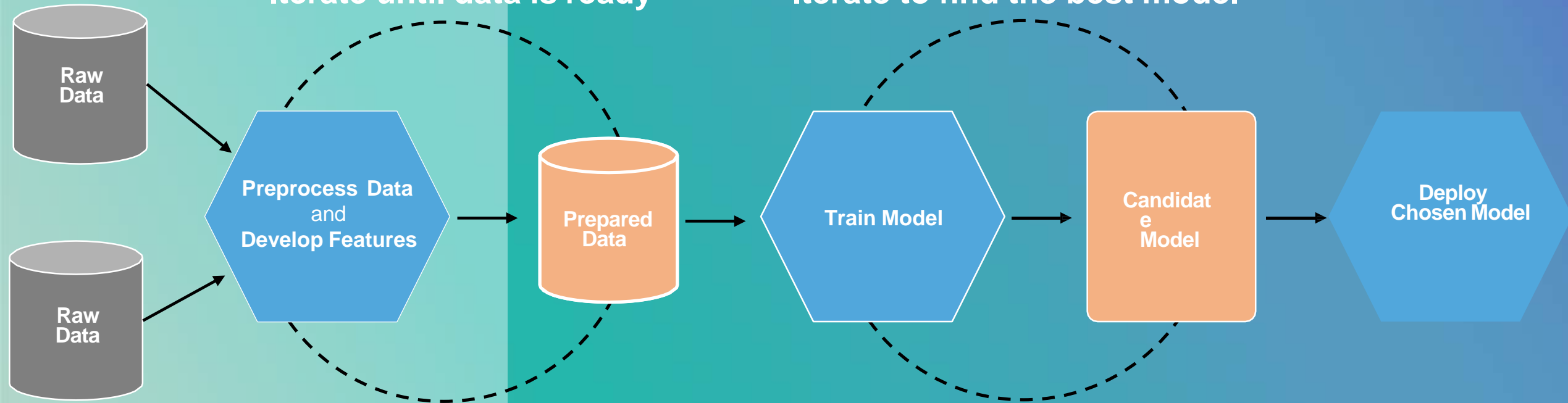
AREIX

# The process of Machine Learning
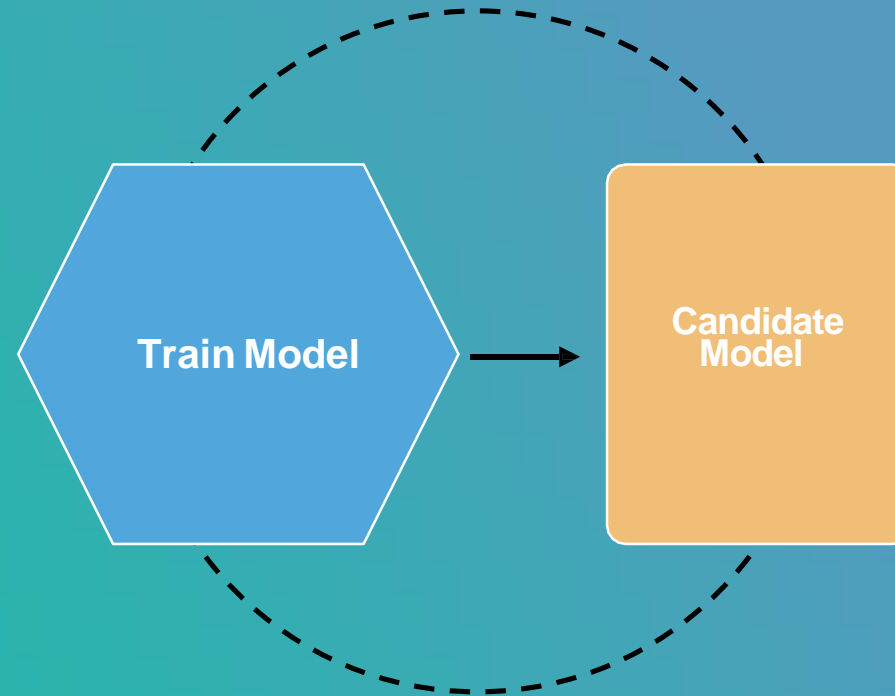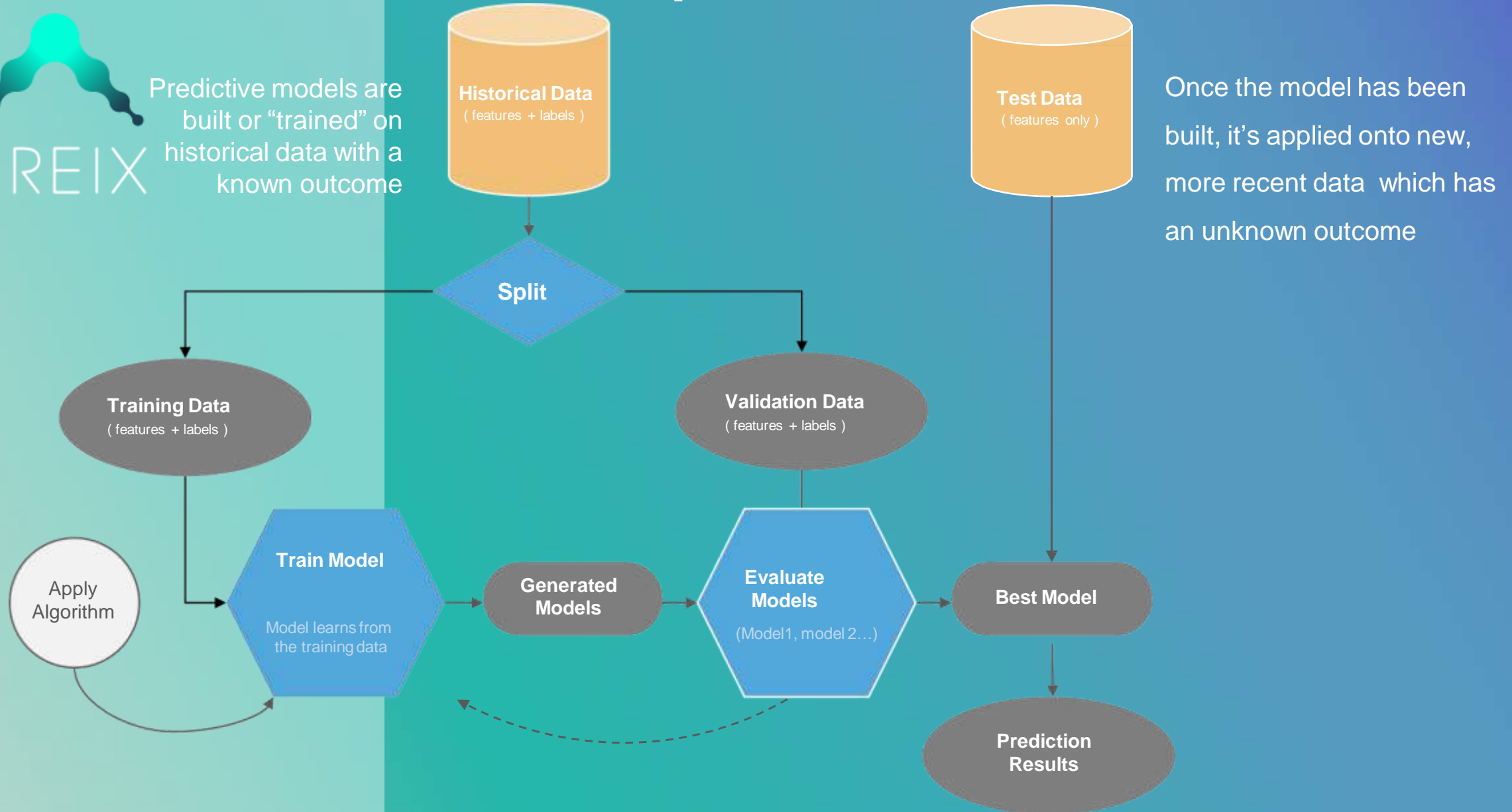
# Modelling

## Iterate to find the best model



**Train model, Evaluate Model & Optimization**
- Find the model that answers the question most accurately by comparing their success metrics
- Determine if your model is suitable for production

# How to develope model?



Predictive models are built or "trained" on historical data with a known outcome

**Historical Data**
( features + labels )

**Test Data**
( features only )

Once the model has been built, it's applied onto new, more recent data which has an unknown outcome

**Split**

**Training Data**
( features + labels )

**Validation Data**
( features + labels )

Apply Algorithm

**Train Model**
Model learns from the training data

**Generated Models**

**Evaluate Models**
(Model1, model 2…)

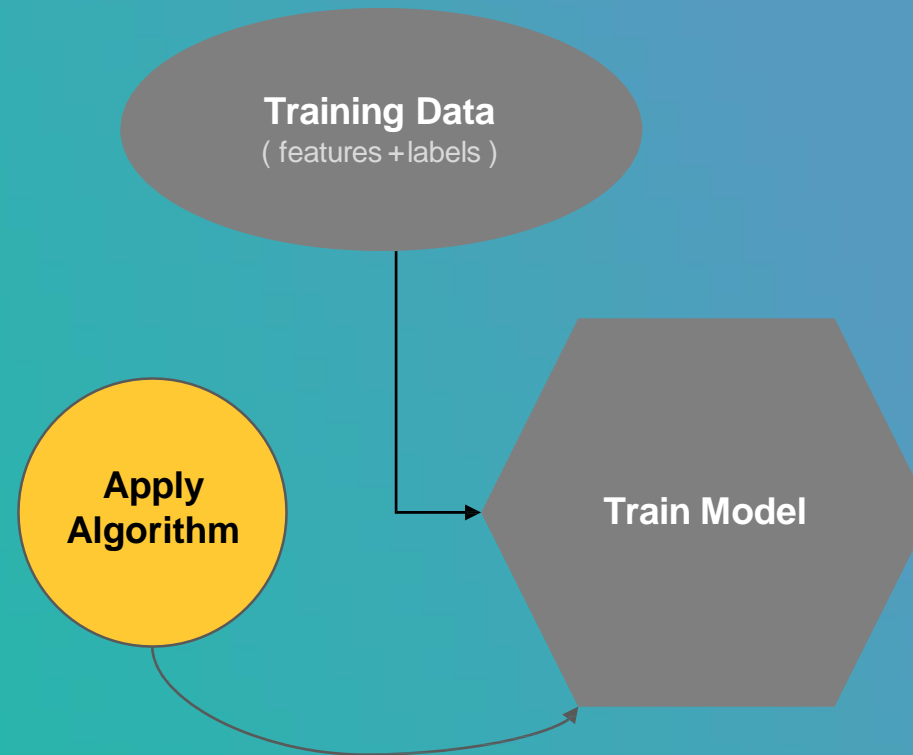**Best Model**

**Prediction Results**

# Train the model

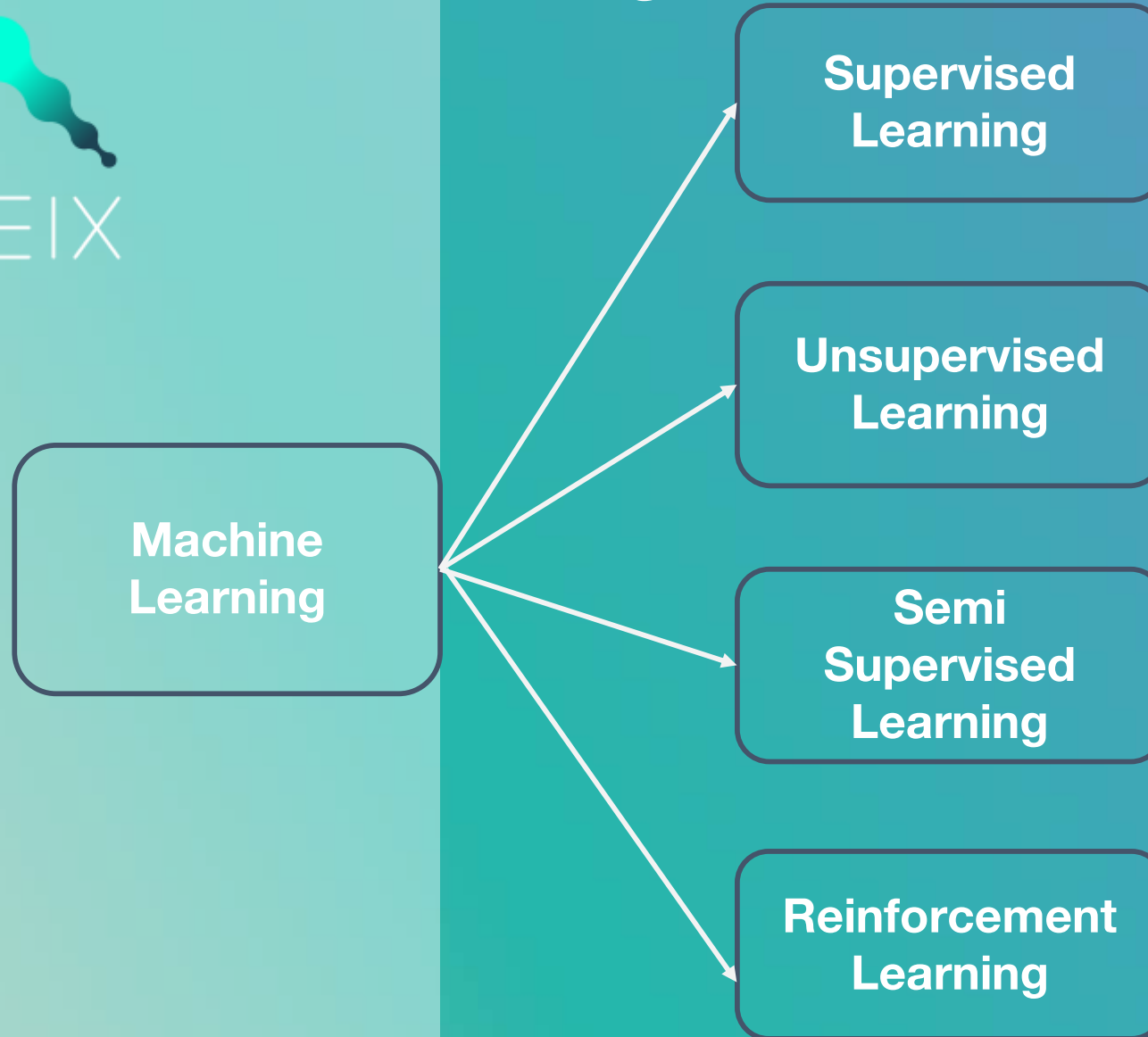## Train a model using algorithms and training data

Algorithm is a set of rules used to solve problems through data processing, math, or automated reasoning.

Machine learning algorithms use computational methods to "learn" information directly from data without relying on a predetermined equation as a model

The algorithm will learn from the training data

patterns that map the variables(features) to the

targets(labels), and it will output a model that
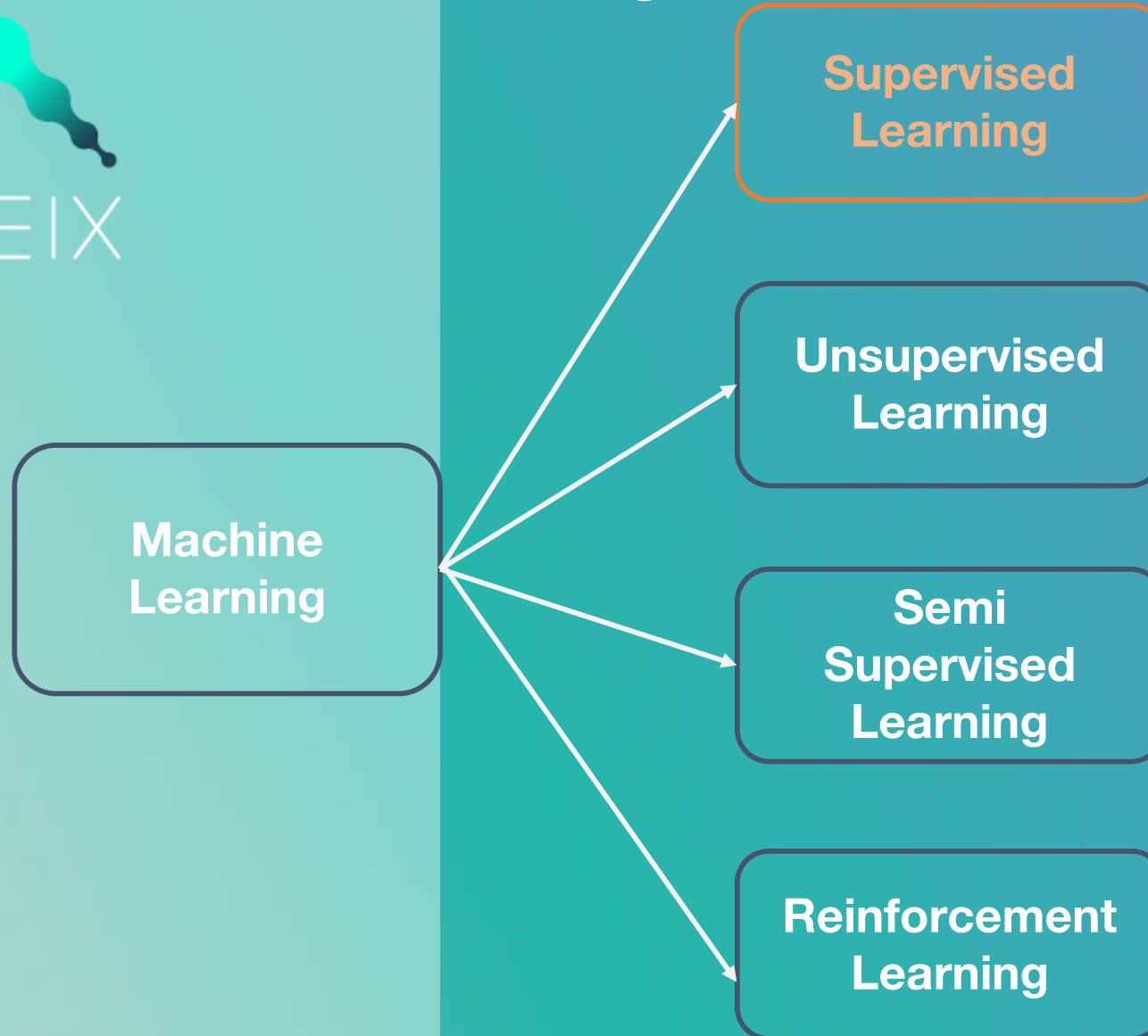
captures these relationships

**Training Data**
( features + labels )

**Apply Algorithm**

**Train Model**

# The categories of machine learning algorithms

**Machine Learning**

**Supervised Learning**

**Unsupervised Learning**

**Semi Supervised Learning**
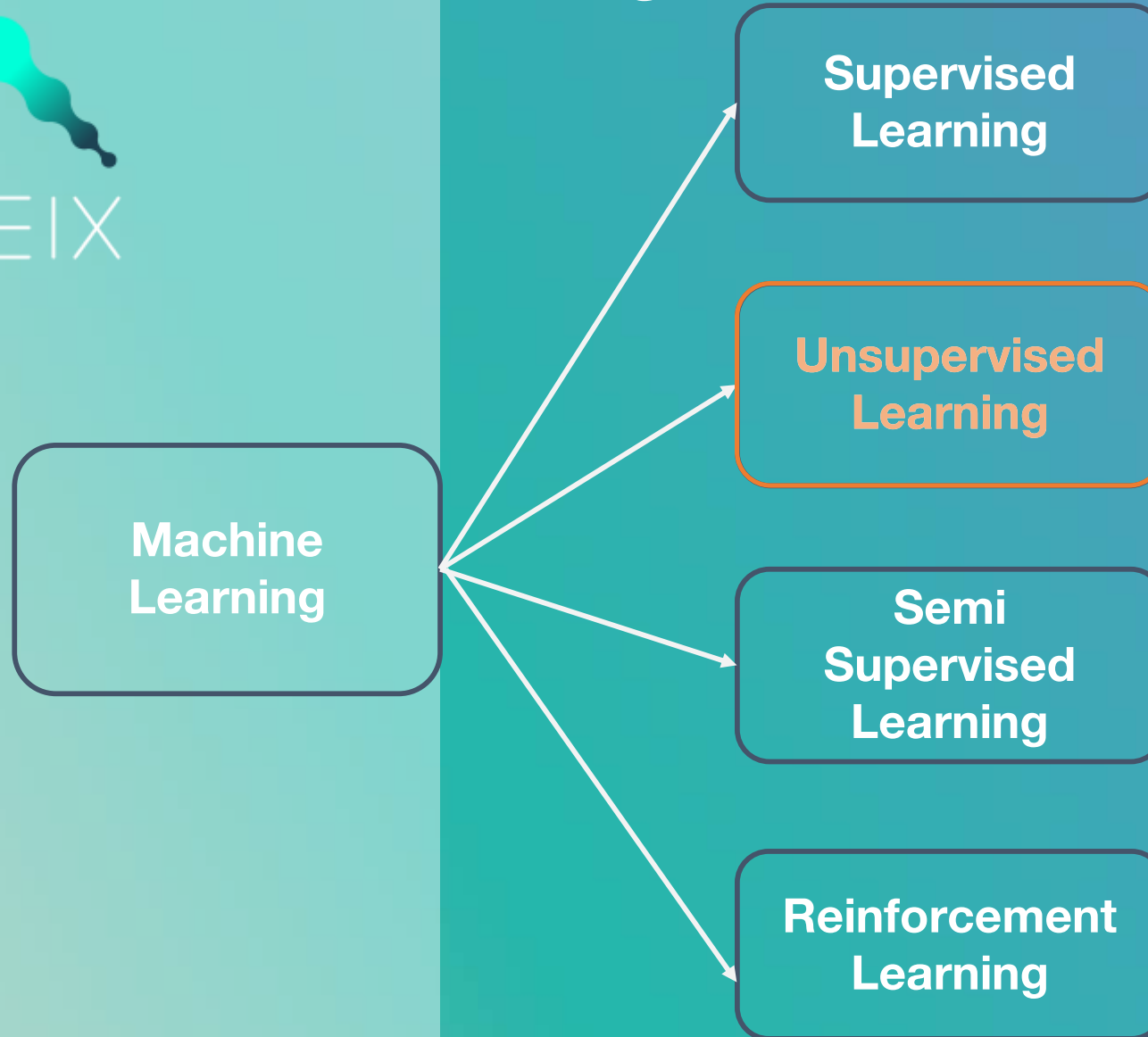
**Reinforcement Learning**

**Train the model with labeled data**

A supervised learning algorithm takes a known set of input data and known responses to the data (output) and trains a model to generate reasonable predictions

(Most machine learning is supervised)

AREIX

AREIX

**Machine Learning**

**Supervised Learning**

**Unsupervised Learning**

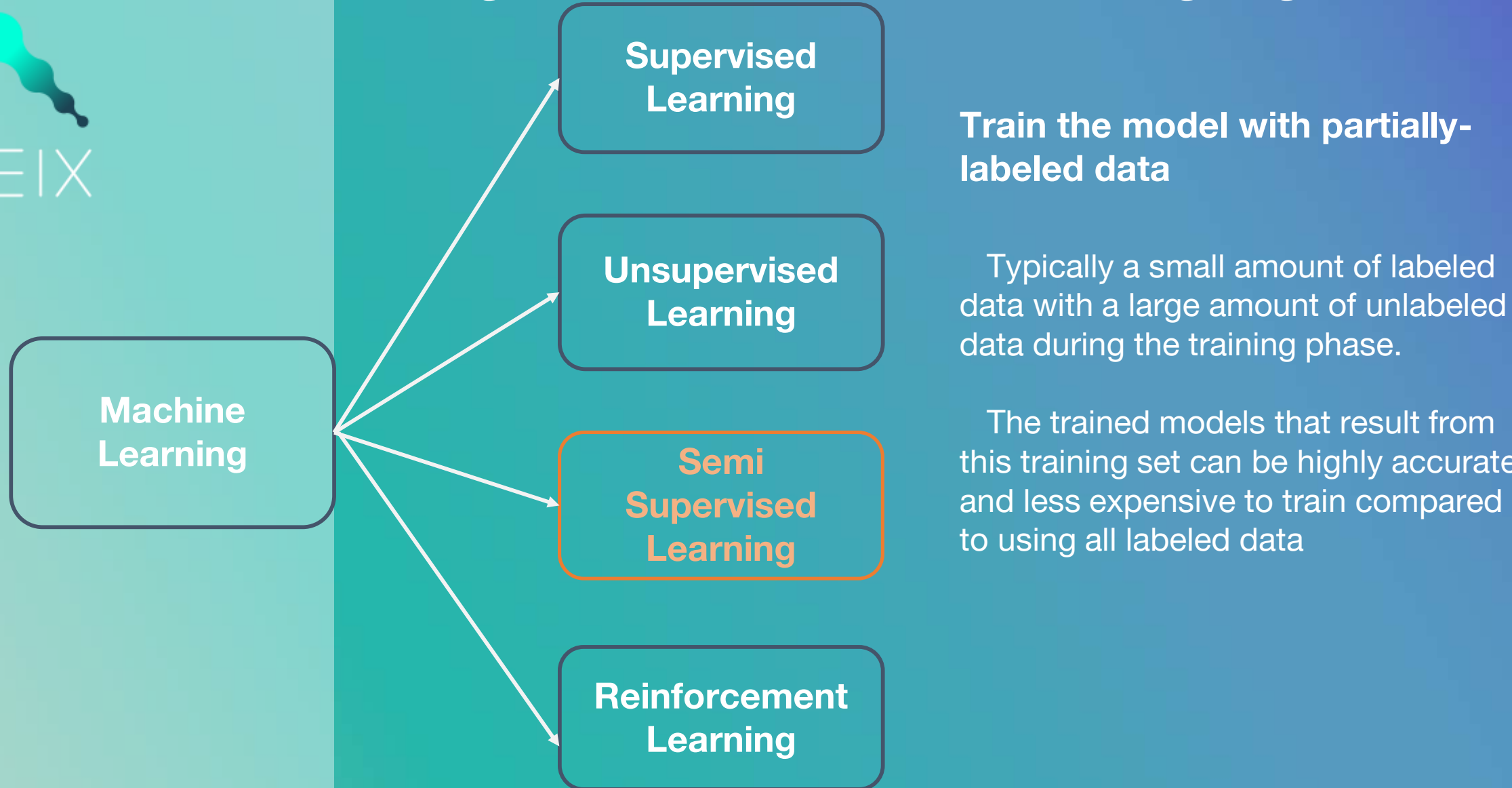**Semi Supervised Learning**

**Reinforcement Learning**

**Train the model with unlabeled data**

The goal is to find patterns or intrinsic structures in the data. It's also a good way to simplify data somehow (reduce dimensions, remove unnecessary variables or detect anomalies).

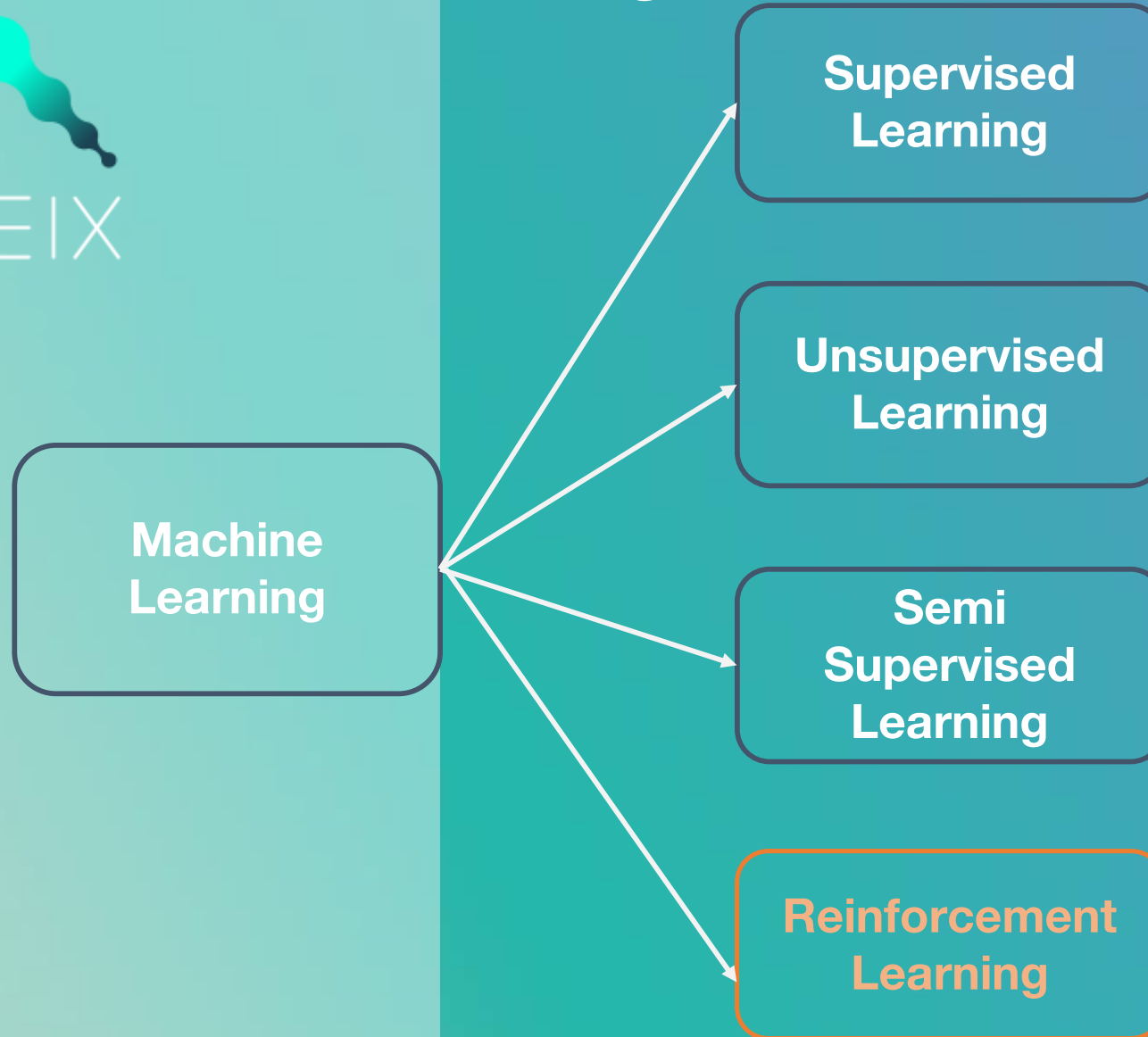Clustering is the most common unsupervised learning technique

AREIX

**Machine Learning**

**Supervised Learning**

**Unsupervised Learning**

**Semi Supervised Learning**

**Reinforcement Learning**

**Train the model with partially-labeled data**

Typically a small amount of labeled data with a large amount of unlabeled data during the training phase.

The trained models that result from this training set can be highly accurate and less expensive to train compared to using all labeled data

# The categories of machine learning algorithms

**Machine Learning**

**Supervised Learning**

**Unsupervised Learning**

**Semi Supervised Learning**

**Reinforcement Learning**

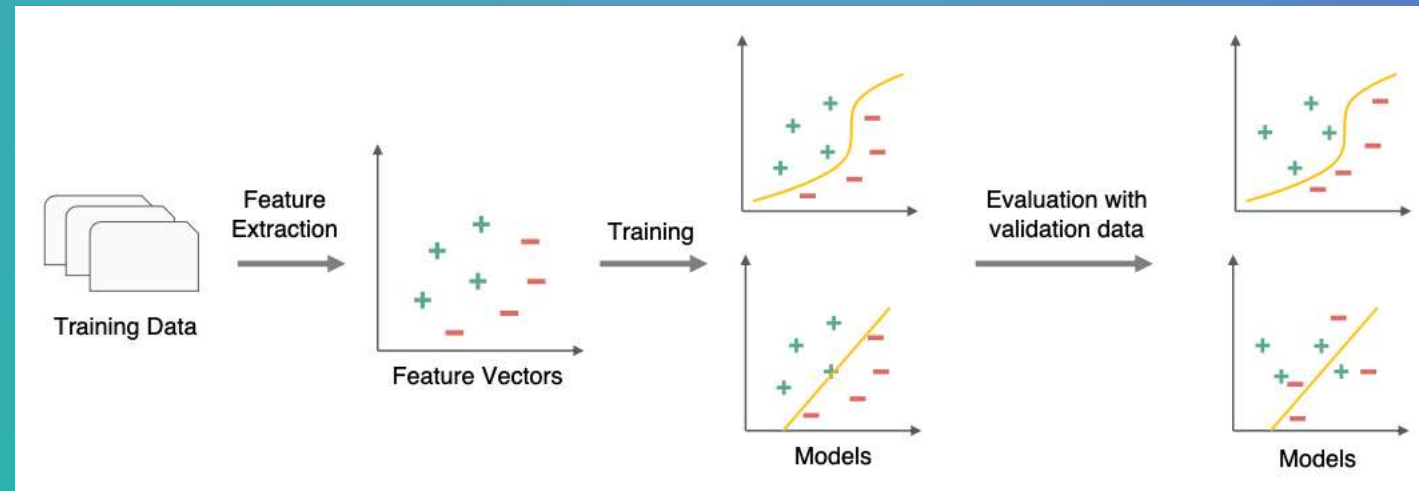**Train the model from a series of 'reward function'**

The model is learned from a series of actions by maximizing a "reward function". The reward function can either be maximized by penalizing "bad actions" and/or rewarding "good actions".
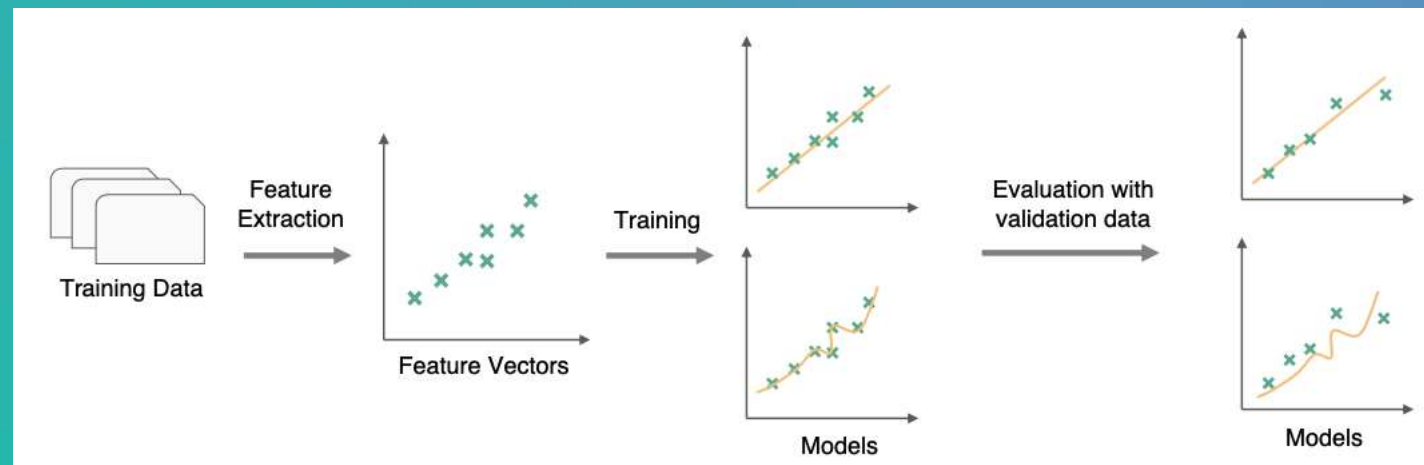
# Supervised Learning tasks

**Classification (Predict the category)**
Identifing to which cataegoty the objecct belongs to

**Regression (Predict the value)**
Predicting a cointinuous-valued attribute associated with an object

# Example for Regression problem

## Example: Predict house's price using linear regression

Suppose we have a dataset giving the living areas and prices of n houses from House Sales in HongKong. Given data like this, we can learn to predict the prices of other houses in Hong Kong

| latitude | longtitude | bedrooms | Living Area ( Feet$^2$ ) | Price ($) |
|---|---|---|---|---|
| -32.432 | 64.342 | 2 | 1180 | 221,900 |
| 34.543 | 43.532 | 3 | 2570 | 538,000 |
| 54.34 | 54.53 | 2 | 770 | 180,000 |
| -12.432 | 324.53 | 3 | 1960 | 604,000 |
| -43.432 | 5.345 | 2 | 1680 | 510,000 |
| 54.543 | 23.423 | 5 | 5420 | 1,225,000 |
| 56.32 | 53.525 | 4 | 1715 | 257,500 |
| -93.54 | 98.34 | 1 | 1060 | 291,850 |
| 45.65 | 54.89 | 2 | 1780 | 229,500 |
| 76.63 | 654.54 | 3 | 1890 | 323,000 |
| 25.654 | 543.63 | 1 | 3560 | 662,500 |
| 75.53 | 43.22 | 2 | 1160 | 468,000 |
| -54.00 | 43.3543 | 4 | 1430 | 310,000 |
| 543.2 | 65.654 | 3 | 1370 | 400,000 |
| 54.6 | 63.435 | 4 | 1810 | 530,000 |
| ... | ... | ... | ... | ... |

Dataset

How much for this house ?

living area = 4876 feet$^2$
bedrooms = 4
latitude = -34.244
logtitude = 31.42

A "mathy" approach where we weight each feature by how important it is then use a weighted sum to estimate housing prices:
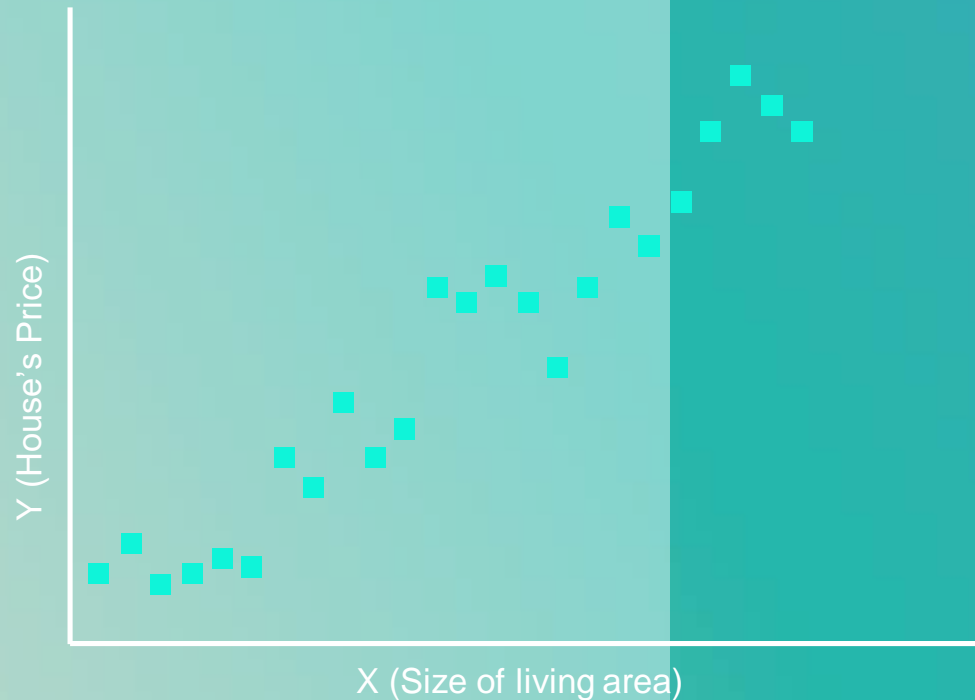
Housing price = w1*bedrooms + w2*size + w3*latitude + w4*longtitude + min_price

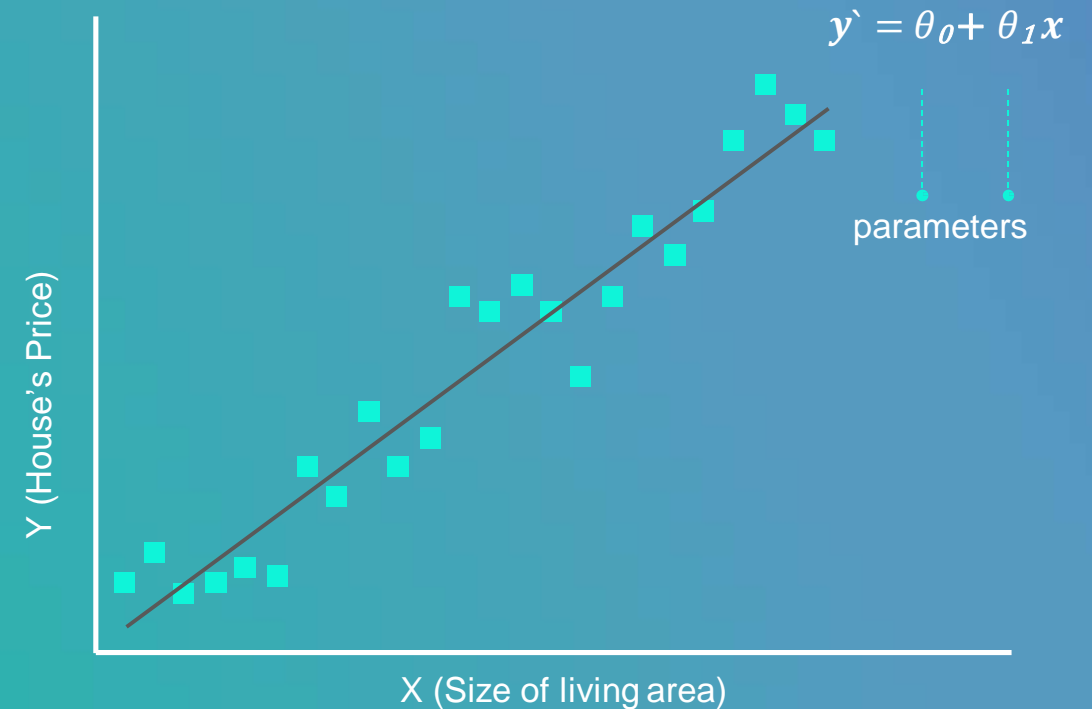# Example for Regression problem

**Use linear regression algorithm to approximate the relationship between $x$ and $y$**

Take linear regression as example, the algorithm is trying to find a best-fit line to represent the relationship between the input feature $x$ and target $y$



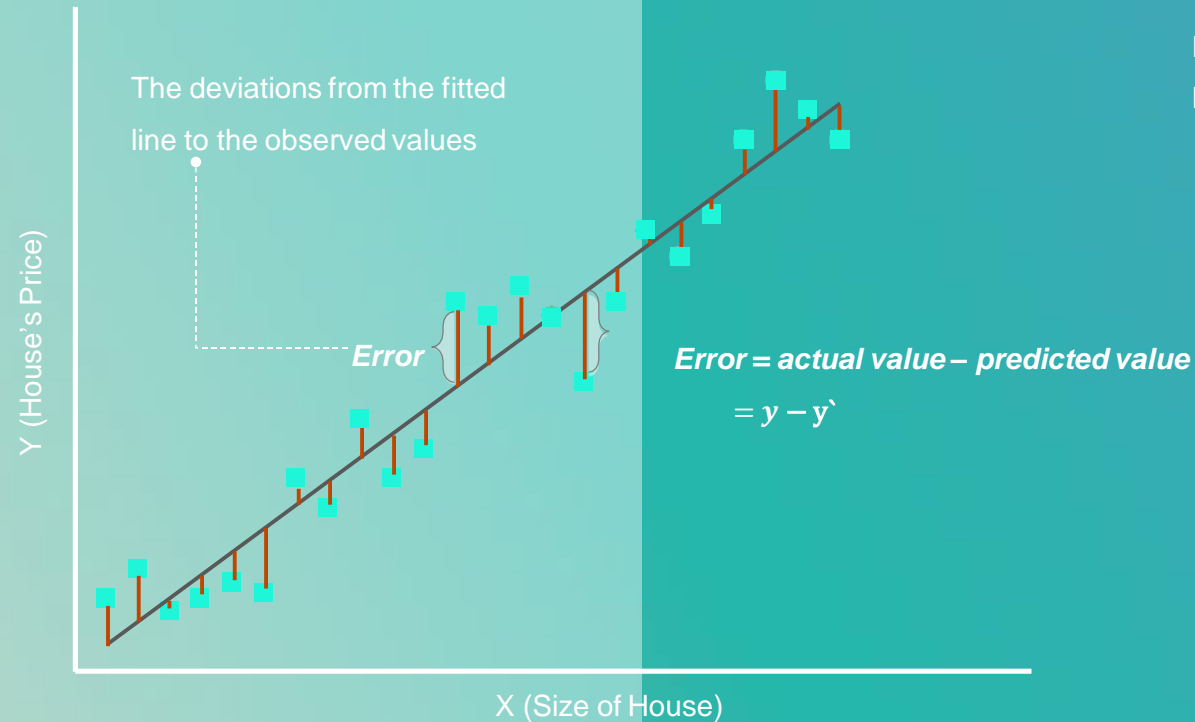Regression line (model) can be presented as :

$$y\grave{} = \theta_0 + \theta_1 x$$

parameters

Y (House's Price)

X (Size of living area)

# Train the model to minimize the loss/error

**The deviations indicate how bad the model's prediction was on the training examples**

**Loss** (i.e. error) is a number indicating how bad the model's prediction was on a single example. If the model's prediction is perfect, the loss is zero; otherwise, the loss is greater

Mean square error (MSE) is a commonly-used function to measure how large the loss is. It's called as **Loss function** or **Cost function**.

The deviations from the fitted line to the observed values

*Error*

*Error = actual value − predicted value*

$$= y - y\grave{}$$

Y (House's Price)

X (Size of House)

$$MSE = \frac{1}{n}\sum_{i=1}^{n}(y - y\grave{})^2$$

*y` is the prediction*
*y is the actual value*

Mean square error (MSE) is the average squared loss per example over the whole dataset.

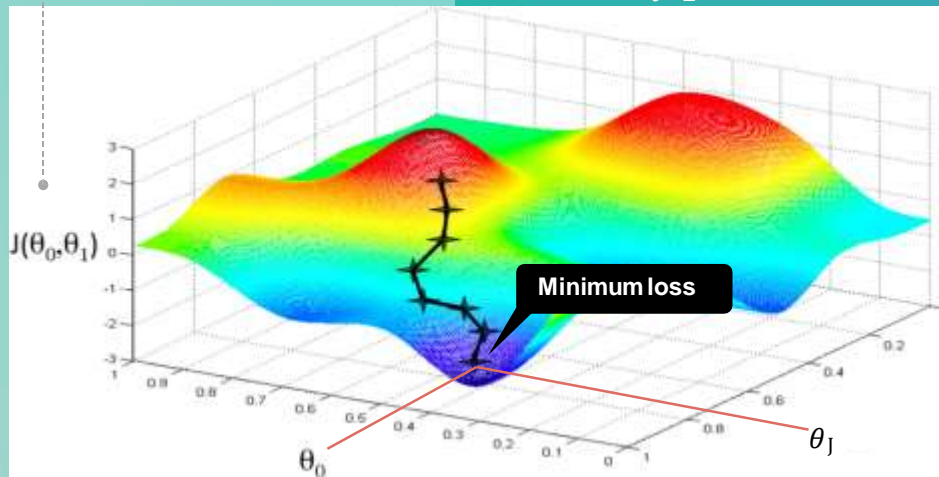The smaller the Mean square error, the better the fit of the line to the data.

AREIX

# How does the model find the "best" parameters ?

**Gradient Descent is one of the most common algorithms to find the good parameters**

A Machine Learning model is trained by starting with an initial guess for the parameters (e.g. weights and bias in neural network ) and iteratively adjusting those guesses until learning parameters with the lowest possible loss

**Loss function/cost function:**
(error)

$$MSE = \frac{1}{n}\sum_{i=1}^{n}(y - y`)^2$$

Usually iterate until overall loss stops changing or at least changes extremely slowly.

When that happens, we say that the model has **converged**.

$J(\theta_0,\theta_1)$

Minimum loss

$\theta_0$

$\theta_J$

**Parameter 1**

**Parameter 2**

$$w := w - \alpha\frac{\delta J(w,b)}{\delta w}$$

$$b := b - \alpha\frac{\delta J(w,b)}{\delta b}$$

**Other optimization algo:**

Gradient Descent (with momentum)

Mini-batch Gradient Descent

Stochastic Gradient Descent

Adam

Adagrad

RMSprop

....

With these 2 specific parameter value, the loss (i.e. MSE) is almost smallest.

# Example for classification problem

**AREIX**

**Classification Algorithm**
Identify what category new information belongs in

**Predict Categories**

**Predict between two categories**

**Binary-Class Classification**

It answers simple

two-choice questions, like

Yes-or-no, true-or-false

Example: Use CT scan to identify whether

has diabetes (True, False)

**Predict between several categories**
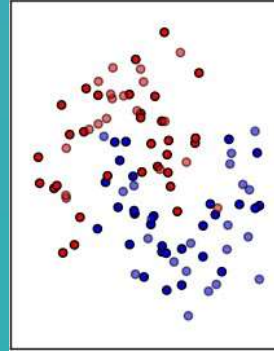
**Multi-Class Classification**

It answers complex

questions with multiple

possible answers

Example: Use CT scan to identify which

type of diabetes (Type1, Type2, Type3..)

# Example: how to classify the data points ?
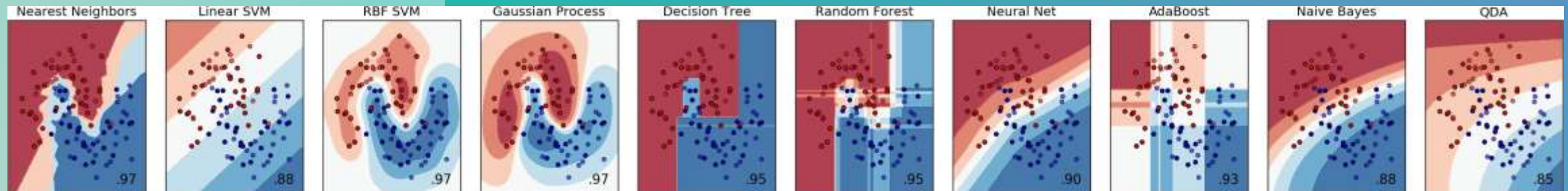
**Input data**



How to classify this dataset into 2 categories ?

"red" and "blue"
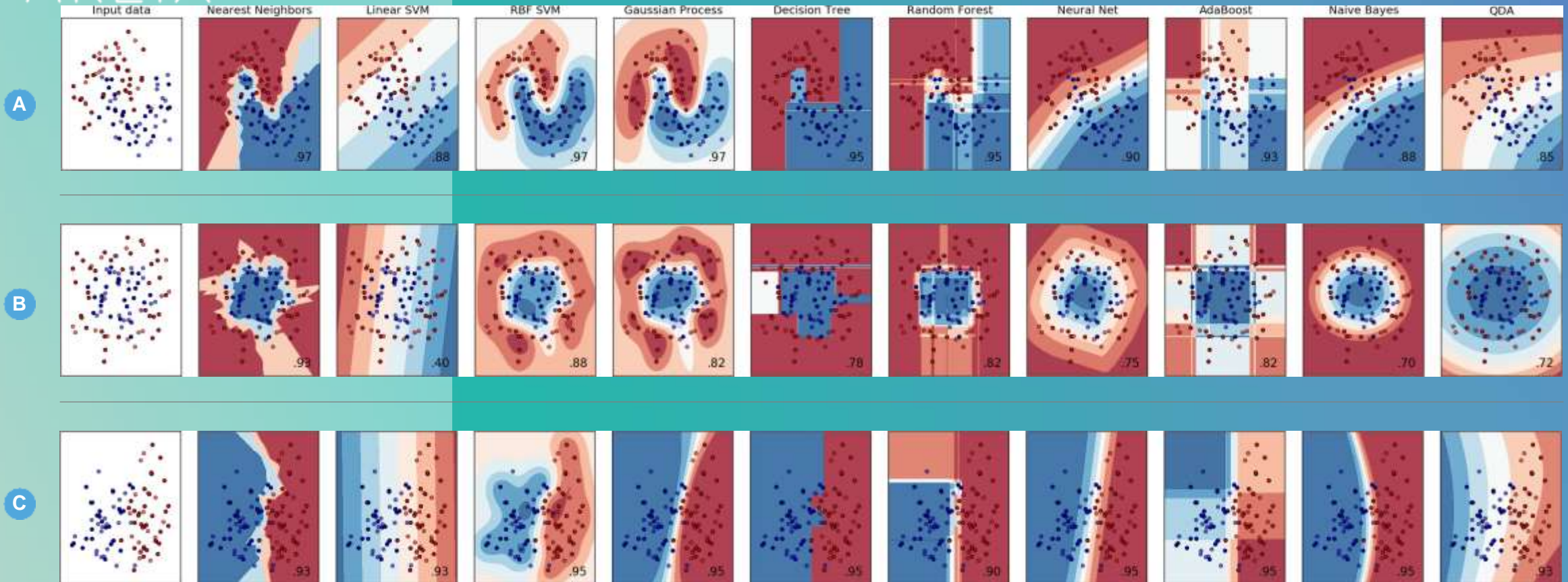
Apply different algorithms on the same data set



The plots show training points in solid colors and testing points semi-transparent.
The lower right shows the classification accuracy on the test set.

# Example: how to classify the data points ?

**Use different algorithms to classify the data set.**

In this example, there are 3 different data sets- A, B, C. You can see how they're classified with different algorithms.



The plots show training points in solid colors and testing points semi-transparent.
The lower right shows the classification accuracy on the test set.

# Example: Stock Prediction



- Data collection & construction
  - API & web scrape
  - Label: outperform S&P500 10%
- Feature engineering
  - Standarlization
  - Feature augmentation
- Model selection & training & evaluation & optimization
  - Model:
    - RandomForest (Bagging)
    - Lightgbm (Boosting)
  - Metrics: Accuracy, Annualized Return, Volatility...
  - HP tunning: RandomizedSearch & 10-fold Crocss Validation

# Cond.

**LightGBM**

**Random Forest**

```python
import lightgbm as lgb
from sklearn import metrics

param = {'num_leaves': 31,
         'min_data_in_leaf': 20,
         'max_depth': 15,
         'num_leaves': 20,
         'objective': 'binary',
         'learning_rate': 0.06,
         "boosting": "gbdt",
         'feature_fraction': 0.8,
         'subsample': 0.2,
#          "bagging_freq": 1,
#          "bagging_seed": 11,
#          'objective':'multiclass',
         "metric": 'None',
         "verbosity": -1}
```

```python
def predict_stocks(test_data):
    X_train, y_train = build_data_set()
    clf = RandomForestClassifier(n_estimators=100)
    clf.fit(X_train, y_train)

    test_data.dropna(axis=0, how="any", inplace=True)
    features = test_data.columns[6:]
    X_test = test_data[features].values
    z = test_data["Ticker"].values

    # Get the predicted tickers
    y_pred = clf.predict(X_test)
    if sum(y_pred) == 0:
        print("No stocks predicted!")
    else:
        invest_list = z[y_pred].tolist()
        print(
            f"{len(invest_list)} stocks predicted to outperform the S&P500 by more than {OUTPERFORMANCE}%:"
        )
        print(" ".join(invest_list))
        print(f"y_pred length: {len(y_pred)}")
        return invest_list


if __name__ == "__main__":
    print("Building dataset and predicting stocks...")
    predict_stocks()
```

```
Building dataset and predicting stocks...
28 stocks predicted to outperform the S&P500 by more than 10%:
CNX OI BAX SWK WGO MAC LH SNA LNC BIIB BWA GES GWW AIZ GNW VIAB DNR R BIG PBI BLK DLX GTN AMP X BBBY LM APD
y_pred length: 286
```

```python
st = time.time()
trn_data = lgb.Dataset(X_train, y_train)
num_round =1000
lgb_clf = lgb.train(param, trn_data, num_round, verbose_eval=300)
pred_y = lgb_clf.predict(X_test, num_iteration=lgb_clf.best_iteration
```

```python
    'column': features,
    'importance': lgb_clf.feature_importance(),
}).sort_values(by='importance', ascending=False)
```

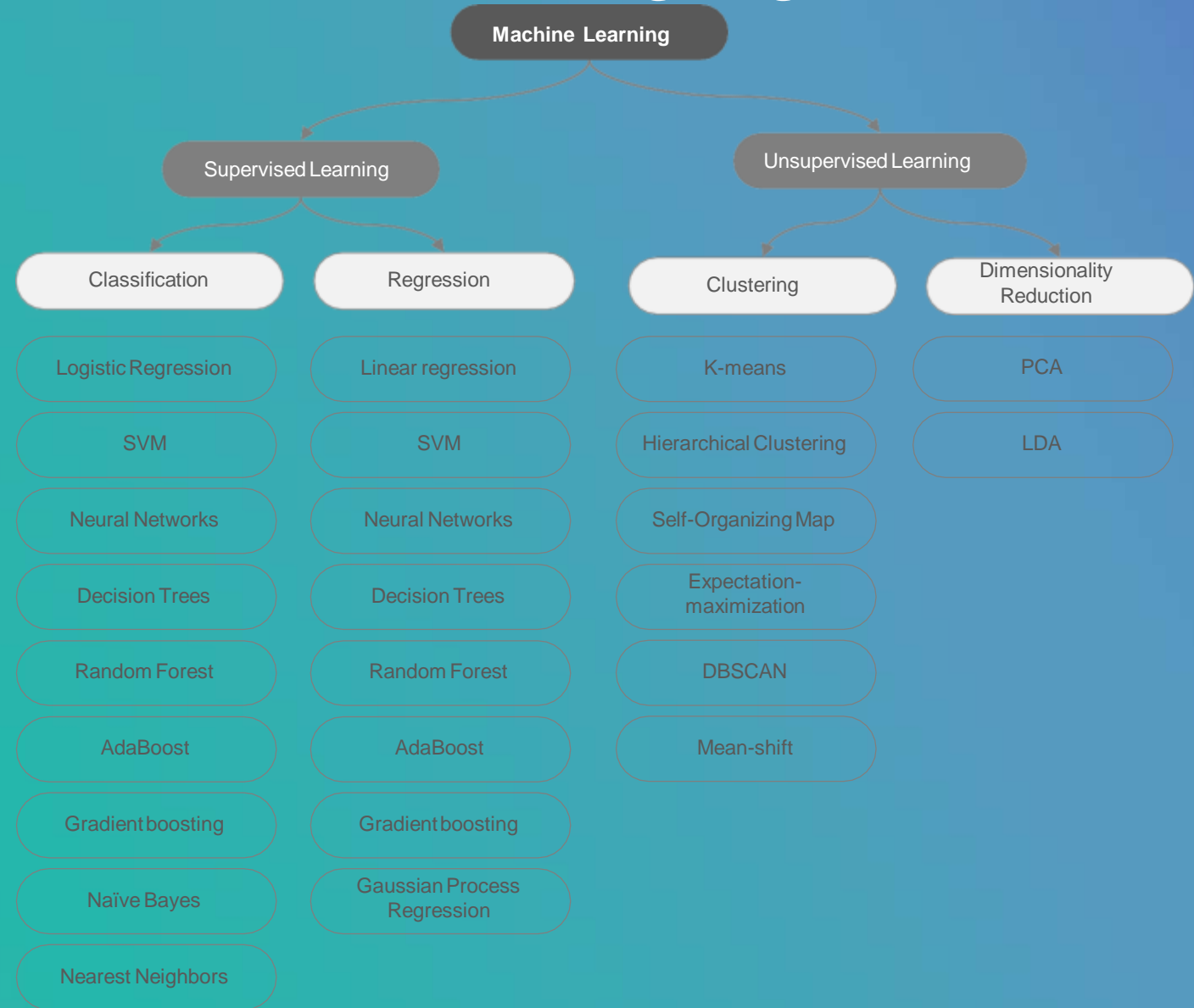|    | column | importance |
|----|--------|-----------|
| 4  | PEG Ratio | 741 |
| 29 | Beta | 705 |
| 3  | Forward P/E | 695 |
| 15 | Quarterly Revenue Growth | 659 |
| 24 | Total Debt/Equity | 627 |
| 38 | Short Ratio | 595 |
| 8  | Enterprise Value/EBITDA | 583 |
| 40 | Shares Short (prior month) | 581 |
| 32 | Avg Vol (3 month) | 571 |
| 20 | Quarterly Earnings Growth | 557 |
| 37 | Shares Short | 548 |
| 25 | Current Ratio | 547 |
| 30 | 50-Day Moving Average | 521 |
| 22 | Total Cash Per Share | 514 |

# Some of machine learning algorithms

AREIX

Q: How to select the right algorithm ?
A: The answer to the question varies depending on many factors, including:

• The size, quality, and nature of data.
• The available computational time.
• The urgency of the task.
• What you want to do with the data.

**Machine Learning**

Supervised Learning

Unsupervised Learning

| Classification | Regression | Clustering | Dimensionality Reduction |
|---|---|---|---|
| Logistic Regression | Linear regression | K-means | PCA |
| SVM | SVM | Hierarchical Clustering | LDA |
| Neural Networks | Neural Networks | Self-Organizing Map | |
| Decision Trees | Decision Trees | Expectation-maximization | |
| Random Forest | Random Forest | DBSCAN | |
| AdaBoost | AdaBoost | Mean-shift | |
| Gradient boosting | Gradient boosting | | |
| Naïve Bayes | Gaussian Process Regression | | |
| Nearest Neighbors | | | |

# Model Evaulation

**Use different metrics to measure the performance of the model**

By using Metrics and scoring to quantify the quality of predictions

- **For Classfication**
  - Accuracy
  - Precision
  - Recall
  - F1
  - ROC_AUC
  - Jaccard Similarity
  - ....

- **For Regression**
  - Max error
  - Mean square error
  - $R^2$ score
  - ....

- **For Clustering**
  - Mutual Information
  - ....

Precision = TP / (TP + FP)
Recal = TP / (TP + FN)
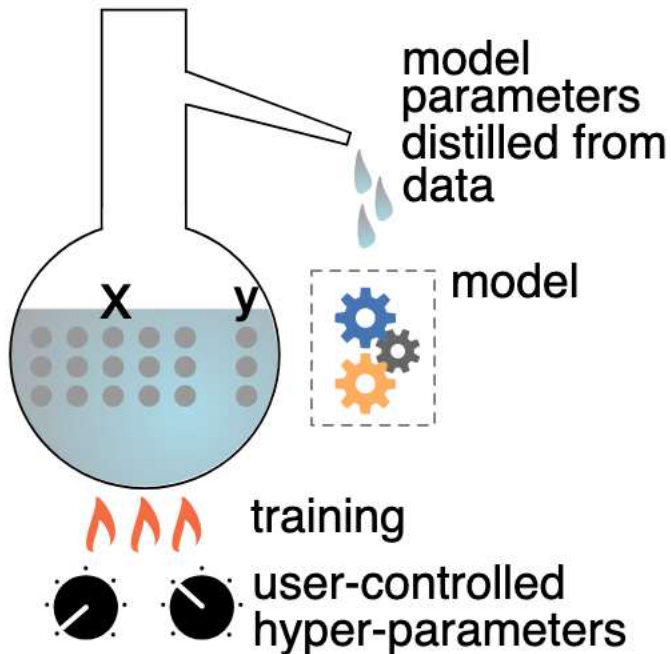F1 = 2 * (precision * recall) / (precision + recall)
ROC_AUC : TPR VS FPR

https://scikit-learn.org/stable/modules/model_evaluation.html

# Parameters and Hyper-parameters



AREIX



- X is the feature vectors
- Y is the target variable

**Iteratively tuning the Hyperparameter so that the model can learn the "best" Parameters from data**

The hyper-parameters are specified by the developer/data scientist while parameters are computed from the data via the algorithms.

- Model's parameters are the variables that your chosen machine learning technique uses to adjust to your data. They are internal to the model. They are estimated or learned from data. They are often not set manually by the practitioner.

- Hyperparameters control how a machine learning algorithm fits the model to the data. Hyper-parameters are specified by the programmer, not computed from the training data, and are often used to tune a model to improve accuracy for a particular data set

The examples of hyper-parameters :

- o Number of layers/units, learning-rate, dropout rate weight d ecay, activation function... in Neural network

- o Number of trees, max depth... in Random Forest

# Model Optimization

**Model Optimization aka hypterparameter tunning is one of the key step to optimize the performance of the model**

**N-fold Cross Validation**

- Use n-1 of the folds as training data
- Validate on the remaining data
- Averagre the values computed in the loop

**Grid Search**

- Gaussian Process based
- Easy to try, but some crucial drawbacks

**Random Search**

- Often leads better result than grid search

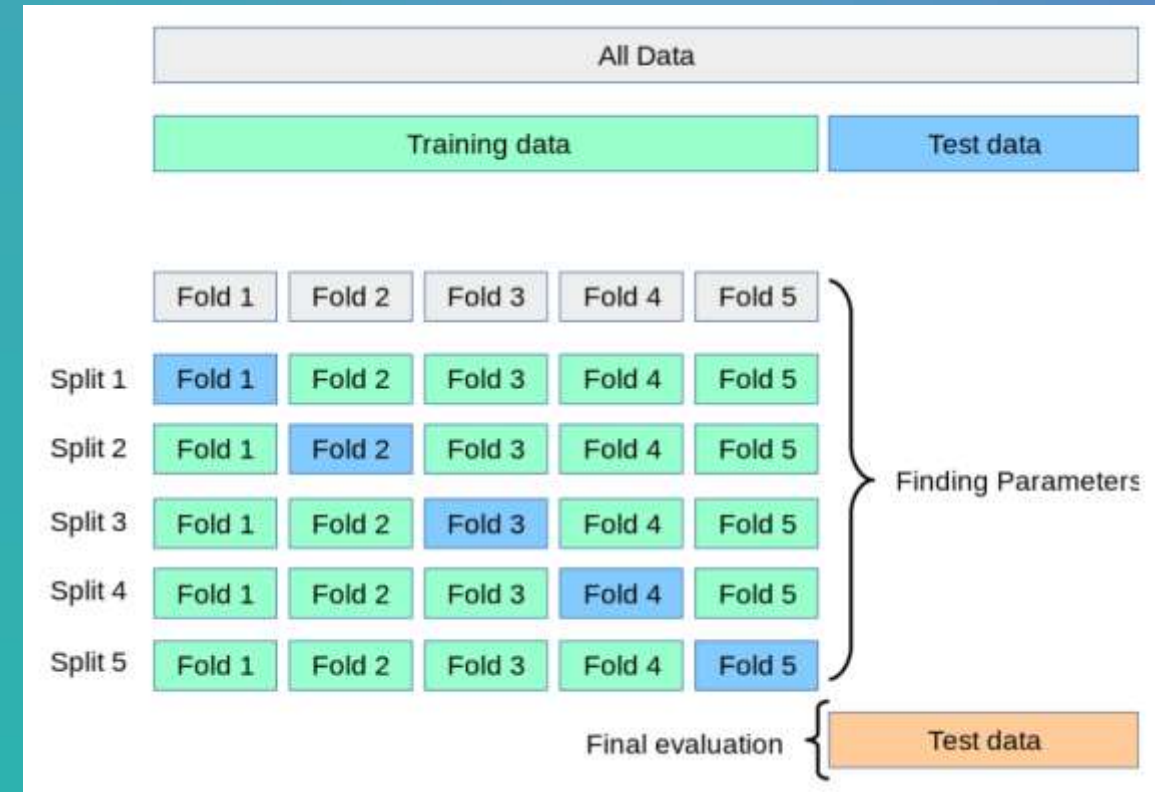**Bayesian Optimazation**

- Random Forest based
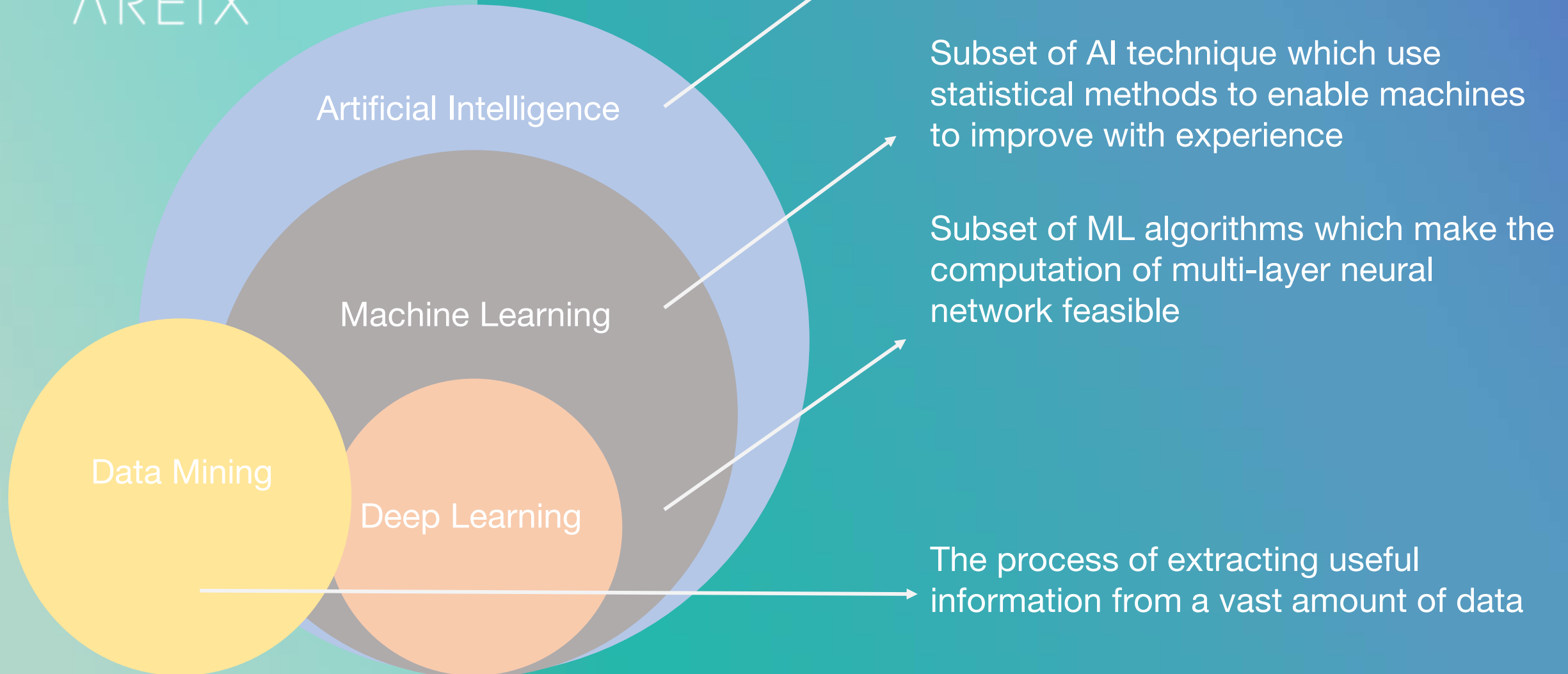- Deep Neural Net based
- Tree Parzan Estimators based

**Recommend Lib: Sklearn, Skopt, Hyperopt**



https://scikit-optimize.github.io/stable/

# Wrap up

- Machine Learning uses historical data to make predictions
- ML process
- Modelling involved training, evaulation & optimization
- Loss function: Measure the error
- Parameters: Computed from training data (automatically trained)
- Hyper-parameters: Human-defined and need (manually fine-tuned)
- Classfication & Regression
- Model evaluation by metrics and scoring
- Tunning HP techniques


- ML VS DM
  - Both are good at pattern recognition and learning from data
  - But serves different purpose


- ML VS DL
  - Structured data & unstructured data
  - Less feature engineering