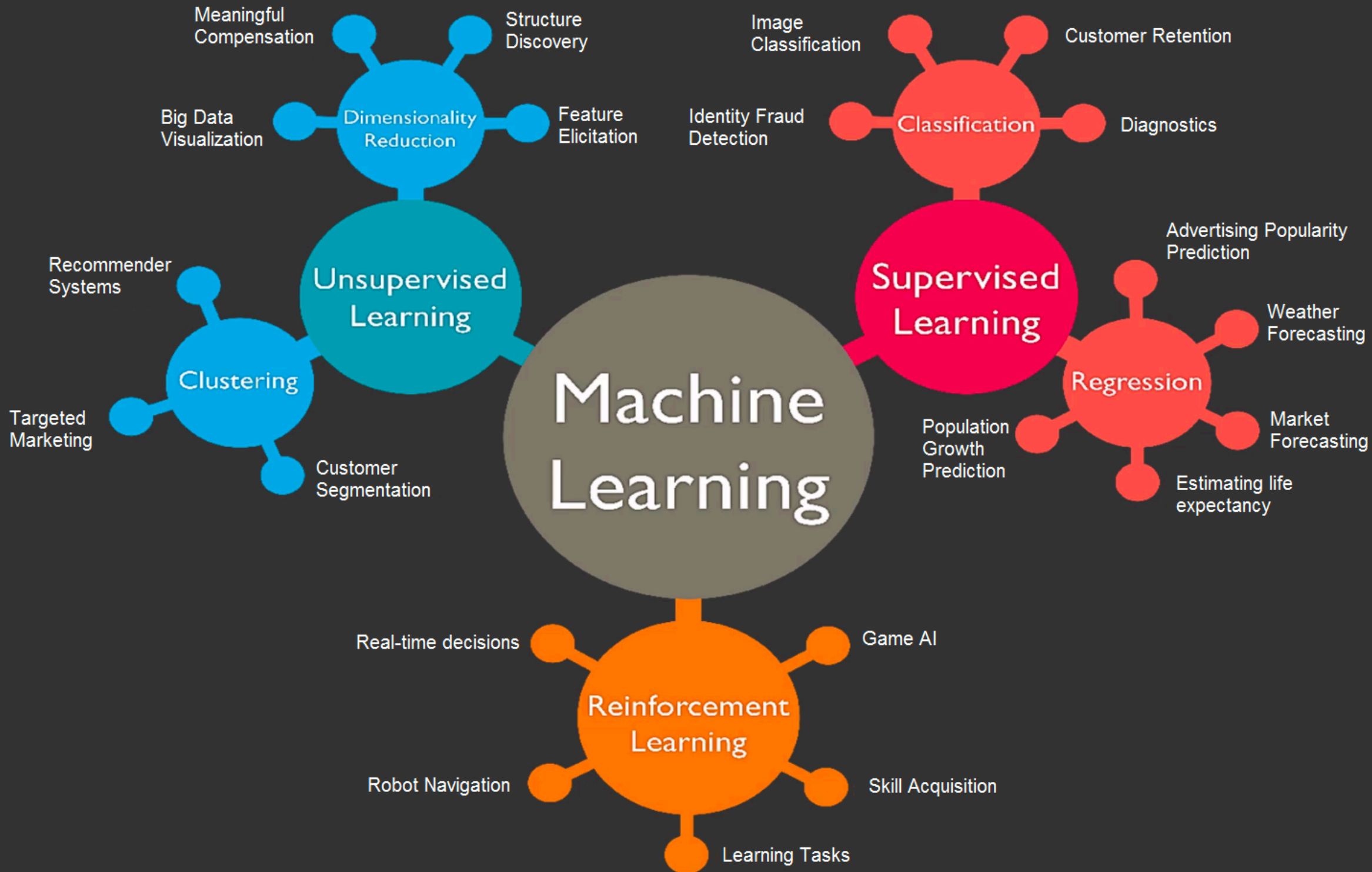# INTRODUCTION TO LOGISTIC REGRESSION

# INTRO TO CLASSIFICATION

# WHERE ARE WE?

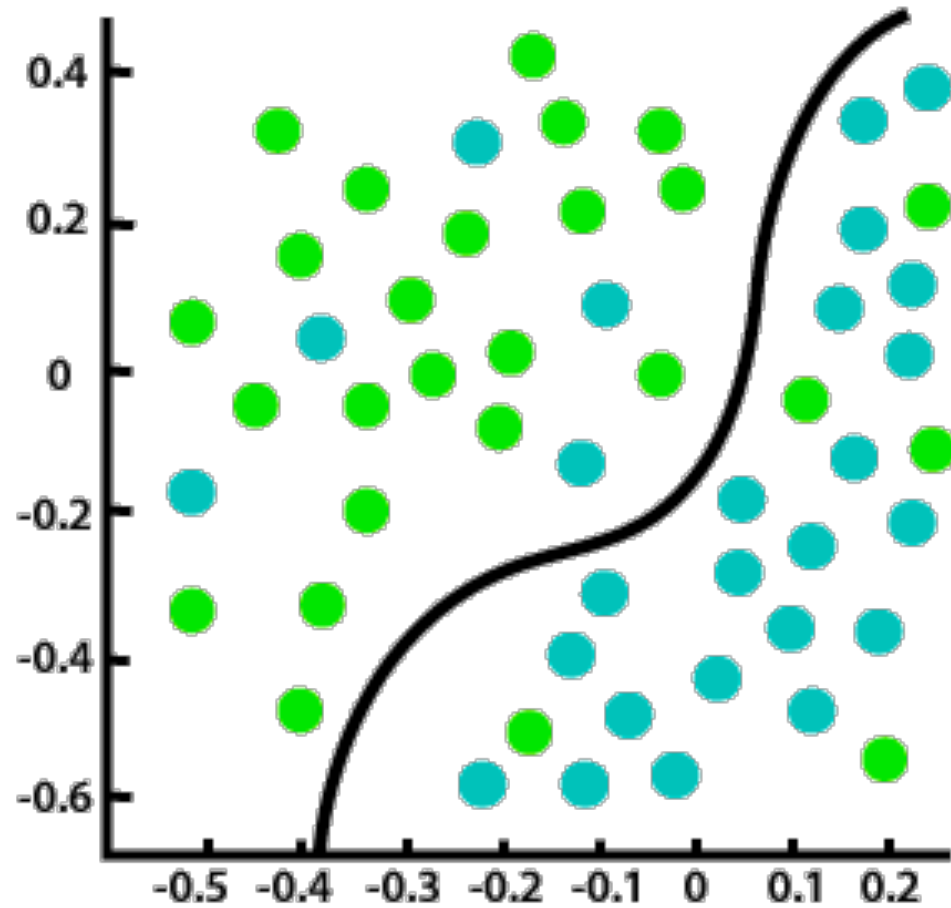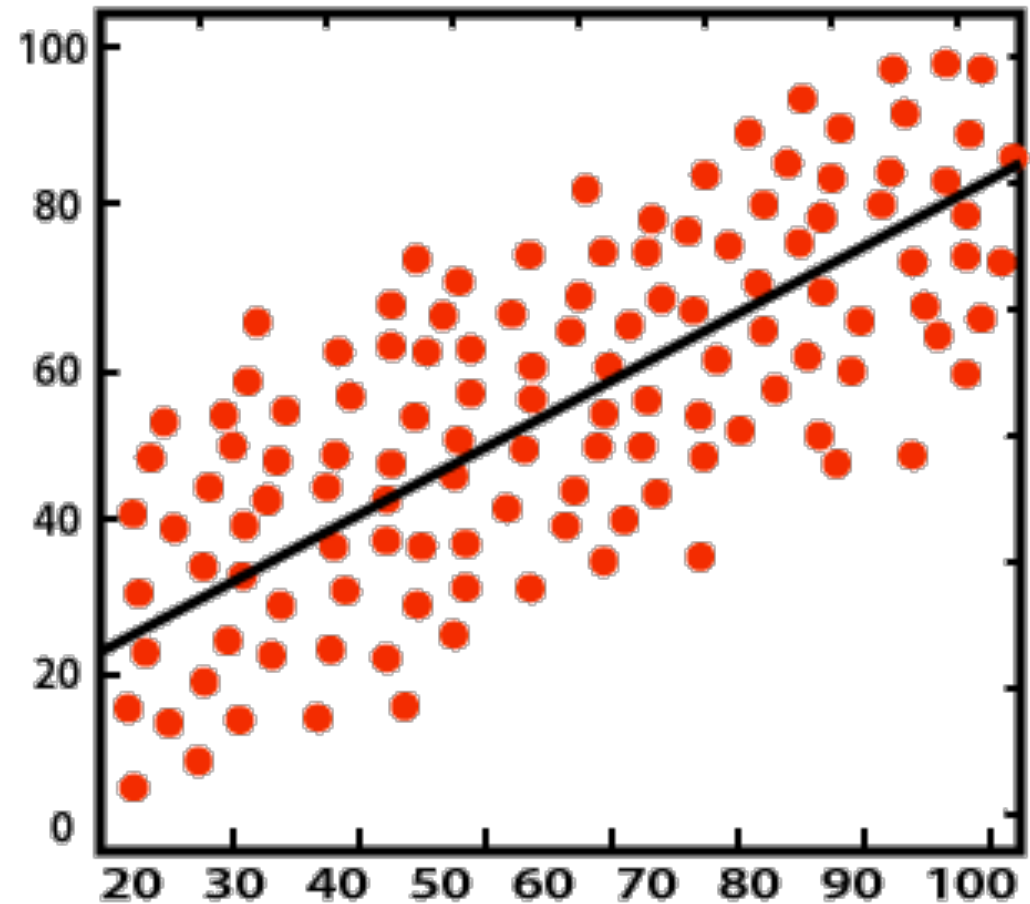|  | Supervised Learning | Unsupervised Learning |
|---|---|---|
| **Discrete** | classification or categorization | clustering |
| **Continuous** | regression | dimensionality reduction |

# CLASSIFICATION VS REGRESSION



Classification

Regression

# COMMON ALGORITHMS

## Regression

- Simple Linear Regression
- Multiple Linear Regression
- Polynomial Linear Regression
- Support Vector Regression
- Decision Tree Regression
- Random Forest Regression
- Neural Network

## Classification

- Logistic Regression
- K Nearest Neighbors
- Support Vector Machine
- Naïve Bayes
- Decision Tree Classification
- Random Forest Classification
- Neural Network

# WHAT IS THE OUTCOME VARIABLE?

# CLASSIFICATION vs REGRESSION

# CLASSIFICATION

# CLASSIFICATION

# LOGISTIC REGRESSION

# LOGISTIC REGRESSION

▸ Logistic regression is a linear approach to solving a classification problem. It will use a linear regression *style* approach to predict the class of an item, but retain the interpretability of linear regression model.



Classification          Regression

# LINEAR REGRESSION can't model a binary outcome

▸ We need a way to *transform* our regression model so that its range changes from [-∞, ∞] to [0, 1].

# LINEAR REGRESSION can't model a binary outcome

# LOGISTIC REGRESSION

‣ To do this, we'll use a log-based transformation called the **logit function**.

‣ It will limit our range to [0,1] and create the right shape for our regression line to match the categorical outcome variable.

$$\frac{1}{1+e^{-x}}$$

# EQUATIONS

▸ Linear regression equation:

$$y = \beta_1 X + \beta_0$$

▸ Logistic regression equation:

$$p = P(y \mid X) = \frac{1}{1 + e^{-\beta_1 X + \beta_0}}$$

$$\text{logit}(p) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

Multinomial Logistic Regression

# INTERPRETING COEFFICIENTS

# COEFFICIENTS

▸ Coefficients are expressed as logodds

▸ Positive coefficients increase the log odds of the response (and thus increase the probability)

▸ Negative coefficients decrease the log odds of the response (and thus decrease the probability).

| | probability | odds | logodds |
|---|---|---|---|
| 0 | 0.10 | 0.111111 | -2.197225 |
| 1 | 0.20 | 0.250000 | -1.386294 |
| 2 | 0.25 | 0.333333 | -1.098612 |
| 3 | 0.50 | 1.000000 | 0.000000 |
| 4 | 0.60 | 1.500000 | 0.405465 |
| 5 | 0.80 | 4.000000 | 1.386294 |
| 6 | 0.90 | 9.000000 | 2.197225 |

# LOGISTIC REGRESSION COEFFICIENTS

▸The intercept is the log of the odds when all predictors are zero

▸Coefficients are the log of the odds of each predictor

## Binary Logit: Churn

| | Estimate | Standard Error | z | p |
|---|---|---|---|---|
| (Intercept) | -1.41 | 0.16 | -8.73 | < .001 |
| Senior Citizen: Yes | 0.41 | 0.11 | 3.60 | < .001 |
| Tenure | -0.03 | 0.00 | -11.38 | < .001 |
| Internet Service: DSL | 0.92 | 0.21 | 4.39 | < .001 |
| Internet Service: Fiber optic | 1.82 | 0.32 | 5.66 | < .001 |
| Contract: One year | -0.88 | 0.14 | -6.25 | < .001 |
| Contract: Two year | -1.68 | 0.24 | -7.02 | < .001 |
| Monthly Charges | 0.00 | 0.00 | 1.11 | .266 |

*n = 3,522 cases used in estimation (Training sample); R-squared: 0.1898; Correct predictions: 79.05%; McFadden's rho-squared: 0.2564; AIC: 3,065.1; multiple comparisons correction: None*

# COEFFICIENTS EXAMPLE: TITANIC SURVIVAL

▸ These coefficients are "log odds"

▸ log odds = 0 indicates even probability

▸ log odds < 0 indicates less likely to occur

▸ log odds > 0 indicates more likely to occur

|  | Log Odds | odds_ratios |
|---|---|---|
| Pclass | -0.796128 | 0.451072 |
| Sex_male | -0.637771 | 0.528469 |
| Sex_female | 0.637771 | 1.892258 |
| Age | -0.441080 | 0.643341 |
| SibSp | -0.324210 | 0.723098 |
| Parch | -0.109567 | 0.896222 |
| Fare | 0.165687 | 1.180204 |
| Embarked_S | -0.094984 | 0.909388 |
| Embarked_C | 0.093482 | 1.097991 |
| Embarked_Q | 0.022137 | 1.022384 |

# INTERPRET THE COEFFICIENTS

Changing the $\beta_0$ value shifts the curve horizontally, whereas changing the $\beta_1$ value changes the slope of the curve.