# Introduction to Machine Learning

Melodie Rush

SAS Customer Success – Principal Data Scientist

Connect with me:

LinkedIn: https://www.linkedin.com/in/melodierush

Twitter: @Melodie_Rush

# Agenda

💡 What is Machine Learning?

💡 Machine Learning Terminology

💡 Intro to ML Modeling Algorithms

💡 Machine Learning in SAS Viya

Machine Learning

HEAVILY HYPED SELF-DRIVING CAR

ONLINE RECOMMENDATION OFFERS

CHAT BOTS

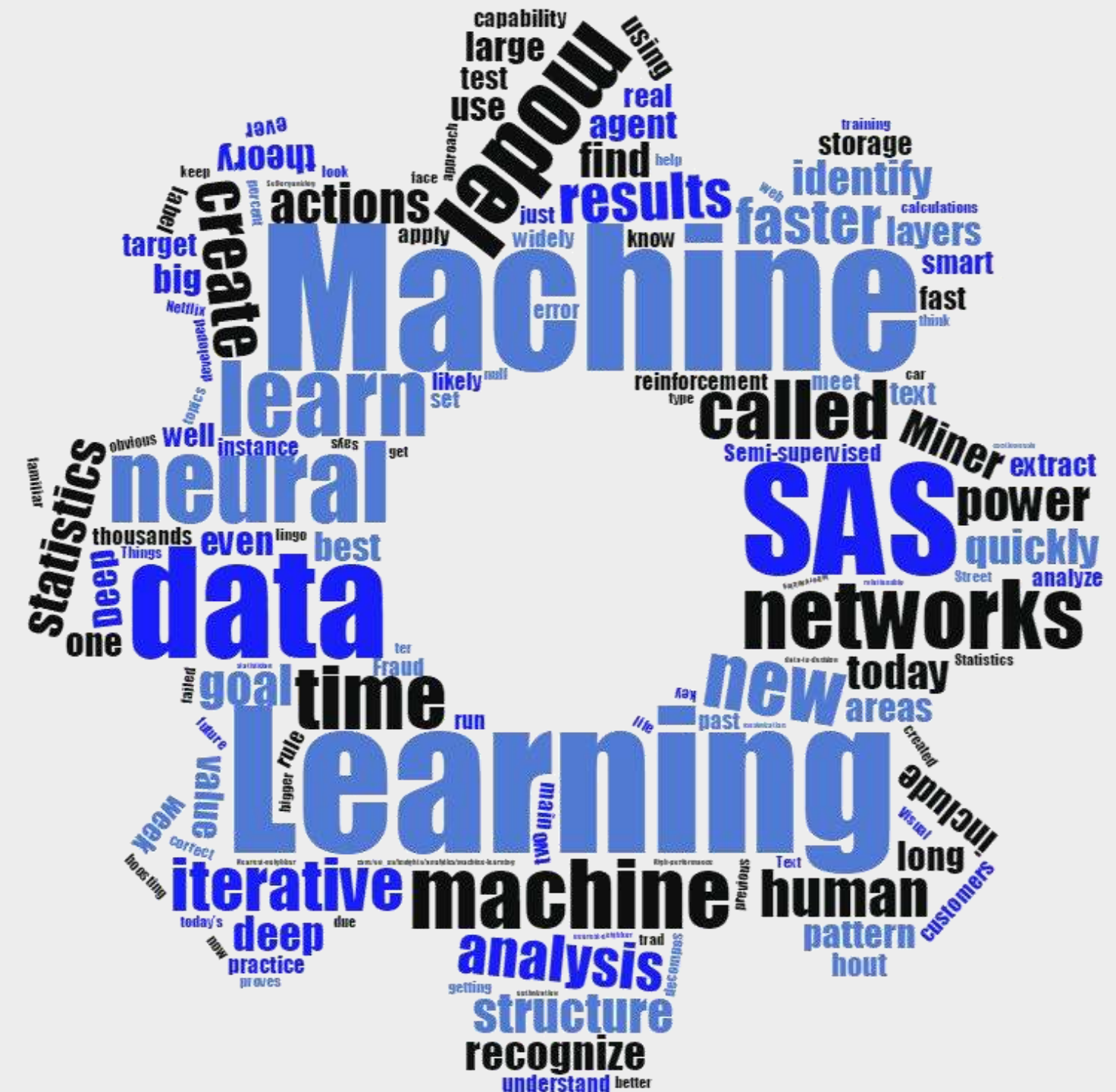FRAUD DETECTION

TARGETED ADS

§sas

# What is Machine Learning?
## Definition

- Automatic
- Adaptive

*Using iterative processes, machine learning builds models that **automatically adapt** with little or no human intervention.*
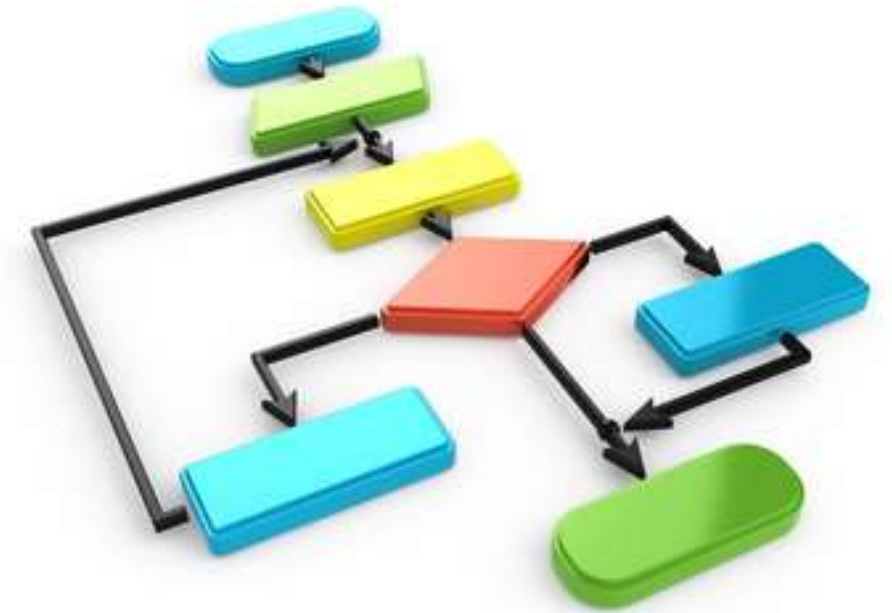
§sas

# Why is it so important now?



Data

Computing Power

Algorithms

§sas

What is
Artificial
Intelligence?

# AI

Science of designing computer systems to support and accelerate human decisions and actions.

§sas

# What is Machine Learning?

ML is a branch of artificial intelligence that **automates** the building of systems that learn from data, identify patterns, and predict future results – with **minimal human intervention**.

§sas

# AI or ML?

## What's the Difference

**AI systems perform tasks that typically require human-level intelligence**

- Understanding Language
- Recognizing images and patterns
- Making Decisions
- Learning from the past

**Machine Learning uses data & algorithms to learn and make decisions**

- ML may be part of the brains in an AI system, or it may be used in a stand-alone usage
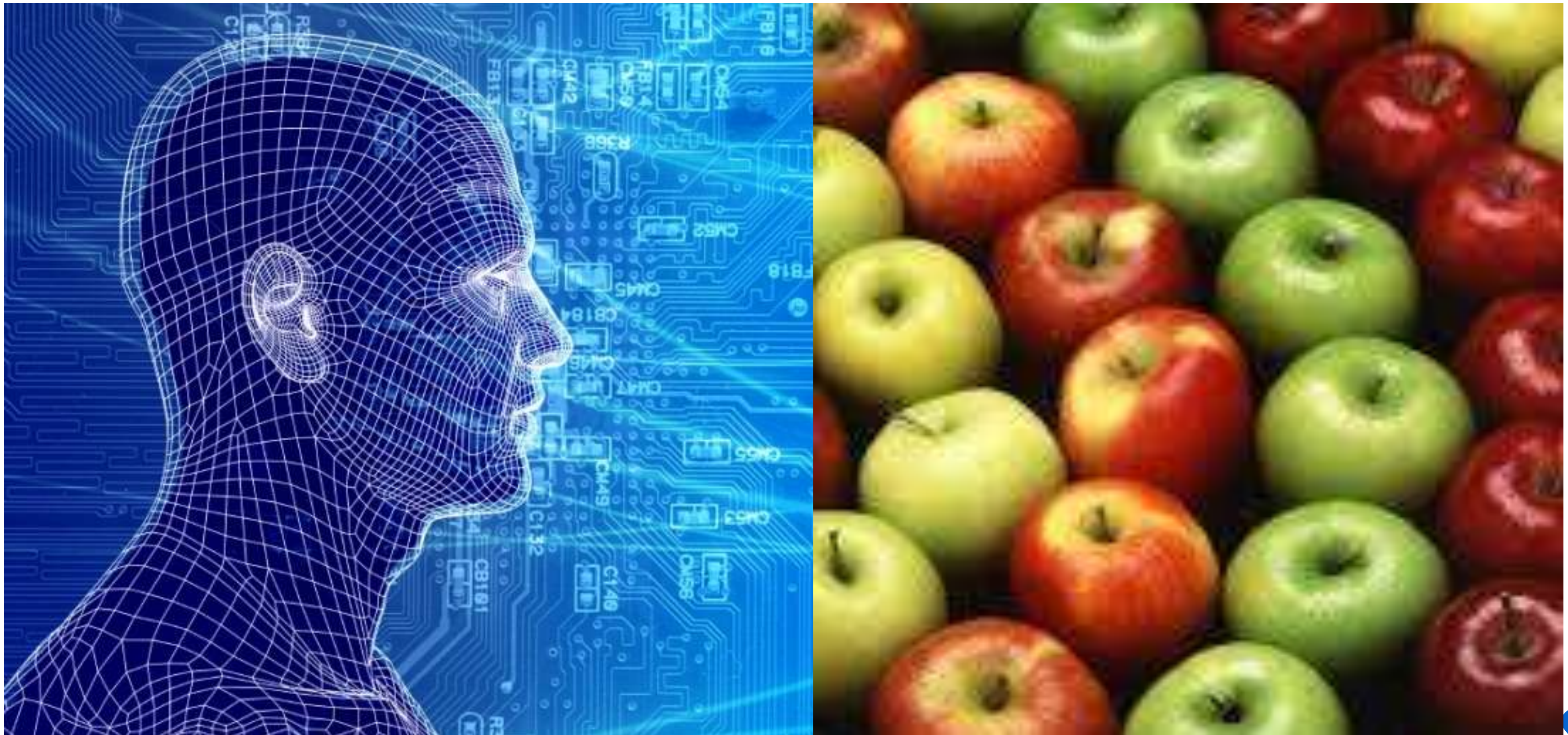- Generally, we think of predictive modeling

§sas

# Terminology

In Machine Learning

# Terminology

Machine learning terms versus inferential statistics terms



What are all these archaic, outmoded and confusing terms?

What are all these new fangled and confusing terms?

- Feature
- Input
- Target
- Object

- Variable
- Independent Variable
- Dependent Variable
- Observation

§sas

# Terminology

## What are Machine Learning terminology?

- In statistics we predict a Y or a dependent variable.

- In data mining, Y is called a target.

- In machine learning, a target is called a label.

- In statistics and data mining our inputs are called X's.

- In machine learning our inputs are called features.

- In statistics and data mining we transform our X's.

- In machine learning we do feature creation.

§sas

# AI or ML?
## What's the Difference

**AI systems perform tasks that typically require human-level intelligence**

- Understanding Language
- Recognizing images and patterns
- Making Decisions
- Learning from the past

**Machine Learning uses data & algorithms to learn and make decisions**

- ML may be part of the brains in an AI system, or it may be used in a stand-alone usage
- Generally, we think of predictive modeling

§sas

# How do Models Learn?
## Distinguish apple from orange

# How do Models Learn?
## Distinguish Granny Smith apple from Fuji apple

# How do Models Learn?
Finding the rotten apple

# How do Models Learn – New Data
## Predictions

# How do Models Learn – New Data
## Predictions

# How do Models Learn – New Data
## Predictions

# How Does Machine Learning Work?

## Unsupervised Learning

Trained on unlabeled examples

# Machine Learning
## Unsupervised Learning



Outlier

## Clustering

Groups of similar data
(e.g. related products on
supermarket basket)

## Anomaly Detection

Identifying outliers
(e.g. Abnormal credit card
transactions)

§sas

# Machine Learning
## Unsupervised Learning Use Cases

Segmentation

Recommendation Engines

Anomaly Detection

Predictive Maintenance

Dimensionality Reduction

§sas

# How Does Machine Learning Work?

## Supervised Learning

Trained on labeled examples

§sas

# Machine Learning
## Supervised Learning



## Regression

Predict a numerical value
(e.g. price of a house,
demand for milk)

## Classification

Predict a label or future
event
(e.g. Cat or Dog, Probability
of loan default)

# Machine Learning
## Supervised Learning Use Cases

Risk Modeling

Fraud Detection

Forecasting

Customer Retention

Document Classification

§sas

# Machine Learning Algorithms

Available in SAS Viya

# Regression

## What Is It?

– Used to identify the relationship between a dependent variable and one or more independent variables

– Many types – linear, logistic, quantile, polynomial, stepwise, ridge, lasso, ElasticNet, etc…

– Oldie but goodie



**Fit Plot for Weight**

| | |
|---|---|
| Observations | 19 |
| Parameters | 2 |
| Error DF | 17 |
| MSE | 126.03 |
| R-Square | 0.7705 |
| Adj R-Square | 0.757 |

Fit — 95% Confidence Limits ----- 95% Prediction Limits

§sas

# Decision Trees

## What Is It?

- Linear separation of data using "if then else" logic

- Separation is performed via an exhaustive search of splitting points for each variable.

- Many different architectural variations based on the above architecture

- Users might refer to them as

  - CHAID Trees

  - CART Trees

  - C4.5 Trees

  - C5.0 Trees.

  - Each of the above is simply a variation on the tree a

# Decision Tree

- Easy to Visualize

# Decision Trees
## Multivariate Step Function

# Random Forest

## What Is It?

- A combination of several "decision trees."

- A random forest consists of a forest of fully trained decision trees.

- The random forest averages the output of all the decision trees in the "forest."

ALL DATA

Random Subset

Tree

Random Subset

Tree

Random Subset

Tree

Random Subset

Tree

Combine:
Average/Vote

§sas

# Random Forest

## Algorithm

- Select a number of trees in the random forest.
- For each tree in the forest, use the following split algorithm:
  - Select a random sample of data.
  - Select a random subset of variables.
  - Determine the best split from the sample of data and the sample of variables.
  - Keep selecting random data and random subsets of variables until the maximum number of trees is trained.
- When all the trees are built, the prediction is the average of all trees.

§sas

# Gradient Boosting
## What Is It?

- A combination of several "decision trees."

- Gradient boosting consists of a **forest** of **small** decision trees ("**shrubs**", "stumps").

- Each **shrub** is poor at predicting target, but each subsequent shrub tries to fit the remaining error.

- Eventually converges to good solution.

# Gradient Boosting

Example: Iterations=0

# Gradient Boosting

Example: Iterations=1

# Gradient Boosting

Example: Iterations=10

# Gradient Boosting

Example: Iterations=25

# Gradient Boosting

Example: Iterations=50

# Gradient Boosting

Example: Iterations=75

# Gradient Boosting

Example: Iterations=100

# Gradient Boosting

Example: Iterations=200

# Gradient Boosting

Example: Iterations=300

# Gradient Boosting

# Neural Network

## What Is It?

– Non-linear relationship between inputs and output

– Prediction more important than ease of explaining model

– Requires a lot of training data

– Users can specify the number of hidden layers, the number of hidden neurons, and associated activation functions for each layer

– Users can configure Input and Target Standardizations, Target Error, and Activation Functions

Many types...
- Feedforward Neural Network
- Radial Basis Function Neural Network
- Multilayer Perceptron
- Convolutional Neural Network (CNN)
- Recurrent Neural Network (RNN)
- Modular Neural Network.
- Sequence-To-Sequence Models

§sas

# Support Vector Machines

## What Is It?

- Enables the creation of linear and nonlinear support vector machine models

- Constructs separating hyperplanes that maximize the margin between two classes

- The vectors (cases) that define the hyperplane are the support vectors

- Enables use of a variety of kernels: linear, polynomial, radia basis function, and sigmoid function. The node also provides interior point and active set optimization methods.

# Clustering
## What Is It?



- Goal: The goal of clustering is to partition data into groups so that the observations within a group are as similar as possible to each other, and as dissimilar as possible to the observations in other groups.

- Many types - Hierarchical, k-means, SOM, etc..

# Ensemble Modeling

## What Is It?

- **Two or more** predictive models **combined** to create a potentially more accurate model

- Works better when model predictions are uncorrelated

- Creates new models by combining the posterior probabilities (for class targets) or the predicted values (for interval targets) from multiple predecessor models.

- 3 Methods
  - Average
  - Maximum
  - Voting

# Machine Learning in SAS Viya

# Interfaces
## Building a Model from Scratch in the Visual Reporting Interface

# Interfaces

## Build Models Using Pipelines in Model Studio

- Drag-and-drop pipelines including preprocessing and machine learning techniques

- Customizable and portable nodes and SAS best practice pipelines (Toolbox)

- Support for SAS coding (macro, data step, procs, batch Enterprise Miner) within pipelines

- Collaboration using the "Toolbox" – a collection of SAS Best Practice Pipelines, in addition to user-generated templates



Example Code for Pipeline

§sas

# SAS® Visual Data Mining and Machine Learning
## Pipelines

˅ Data Mining Preprocessing

- Anomaly Detection
- Clustering
- Feature Extraction
- Feature Machine
- Filtering
- Imputation
- Interactive Grouping
- Manage Variables
- Reject Inference
- Replacement
- Text Mining
- Transformations
- Variable Clustering
- Variable Selection

˅ Supervised Learning

- Batch Code
- Bayesian Network
- Decision Tree
- Forest
- GLM
- Gradient Boosting
- Linear Regression
- Logistic Regression
- Model Composer
- Neural Network
- Quantile Regression
- Score Code Import
- SVM

˅ Postprocessing

- Ensemble

˅ Miscellaneous

- Data Exploration
- Open Source Code
- SAS Code
- Save Data
- Score Data
- Scorecard
- Segment Profile

§sas

# Building Pipelines
## Use prebuilt templates or automatically generate the pipeline

# Automated Pipelines



- ✓ Repository of best practice pipelines
- ✓ Models by SAS or by end-user
- ✓ Dynamically reads thru data
- ✓ Fixes data quality issues w/ ML
- ✓ Performs Data transformations
- ✓ Recommends & builds models
- ✓ Optimizes across models
- ✓ Fully editable, no black-box

§sas

# Automatically Create Features
## Feature Machine Node

# Automatically Tune Hyperparameters for Multiple Model Types

## Model Composer Node

# Interpretable Machine Learning

## Why is it important?



Flu

age
weight
sneeze
headache
no fatigue

sneeze
headache
no fatigue

Model                    Data & Prediction              Explanation              Human Decision

§sas

# Interpretable Machine Learning
## Popular Approaches

1. Variable Importance

2. Partial Dependency Plots

3. Individual Conditional Expectation (ICE)

4. Local Interpretable Model-Agnostic Explanation (LIME)

5. SHapley Additive exPlanations (SHAP)

§sas

# Automated Insights & Interpretability
## Description in simple language

# Automated Insights & Interpretability
## Model Interpretability Charts

- Variable Importance Plots and Rankings
- Partial Dependence (PD) Plots
- LIME (Local Interpretable Model-agnostic Explanations)
- ICE (Individual Conditional Expectation) Plots
- Kernel SHAP Method (Shapley Values)

# Automated Insights & Interpretability
## Model Interpretability Charts

Each interpretability chart has insights included

# Interfaces
## Building a Model Using SAS Studio Tasks

# SAS Visual Data Mining and Machine Learning
## Programming Tasks in SAS Studio

# Interfaces
## Building a Model Using SAS Studio Snippets

# Interfaces
## Building a Model Using Open Source

# Review

## What we covered today

💡 What is Machine Learning?

💡 Machine Learning Terminology

💡 Intro to ML Modeling Algorithms

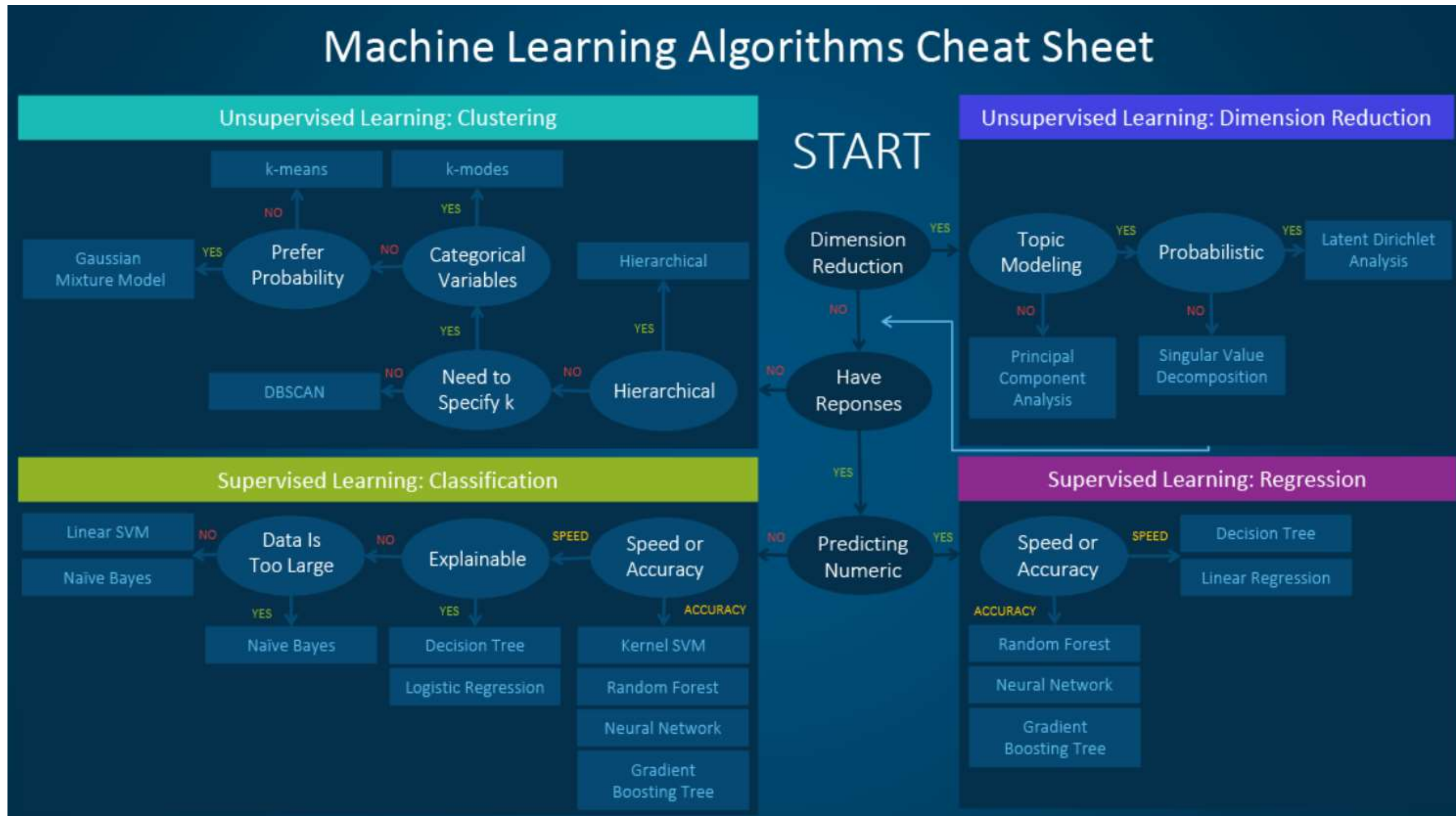💡 Machine Learning in SAS Viya

§sas

# Resources

Where to learn more

§sas

# Machine Learning Algorithms Cheat Sheet



Access Here

# Recommended Resources

An Overview of SAS® Visual Data Mining and Machine Learning on SAS® Viya
https://support.sas.com/resources/papers/proceedings17/SAS1492-2017.pdf

Video  - Automated Machine Learning at Scale
http://www.sas.com/en_us/webinars/automated-machine-learning-scale.html

Machine learning - what it is and why it matters (reading)
http://www.sas.com/en_us/insights/analytics/machine-learning.html

Live web and classroom training - Big Data, Data Mining, and Machine Learning
Big Data course

§sas

# SAS Tutorial

Videos

How to Choose a Machine Learning Algorithm
https://youtu.be/-oZcf0QEzYM

Transforming variables in SAS
https://communities.sas.com/t5/SAS-Data-Mining-and-Machine/New-video-Transforming-Variables-in-SAS/m-p/710687#M8553

§sas

# Thank you for your time and attention!

Introduction to Machine Learning

Connect with me:
LinkedIn: https://www.linkedin.com/in/melodierush
Twitter: @Melodie_Rush

§sas