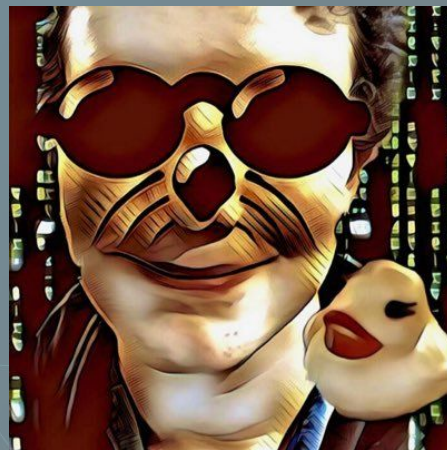# How Much Food Coloring Can Your Robot Handle?

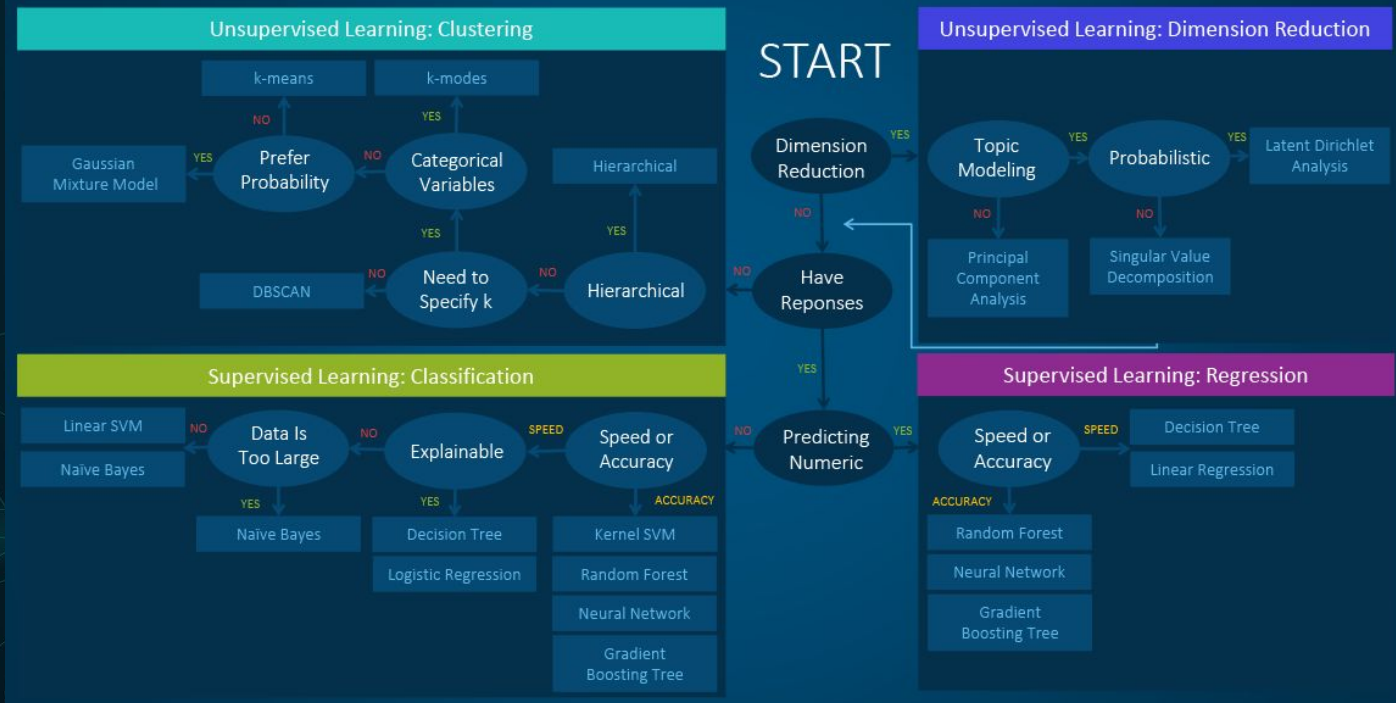An Intro to Poisoning Machine Learning Systems

# Who Am I?

- Corbin Frisvold (@QuesoSec)
- GitHub: @Kousei03
- Maker.godshell.com
- 17 years old
- Whitewater kayaker by day
- Hacker, maker, mathematician, and scientist by night
- Currently in a one year math degree program
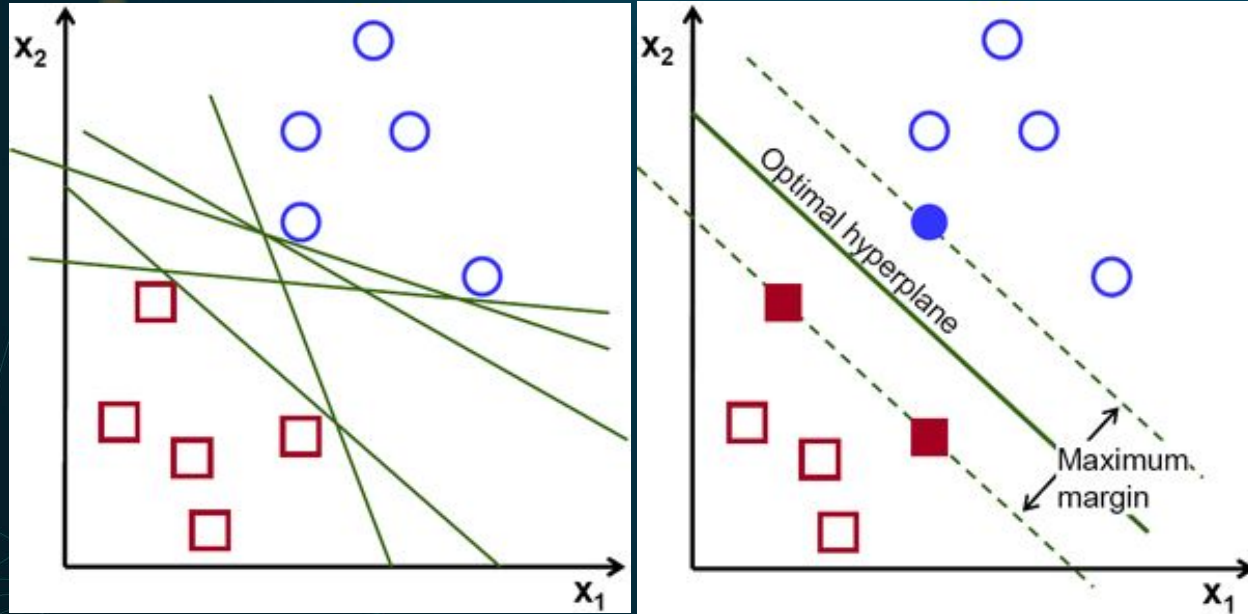- Researcher at Lafayette, UVM, and Harvard
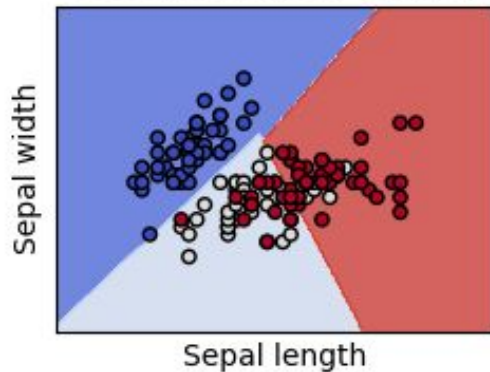
# Brief Intro

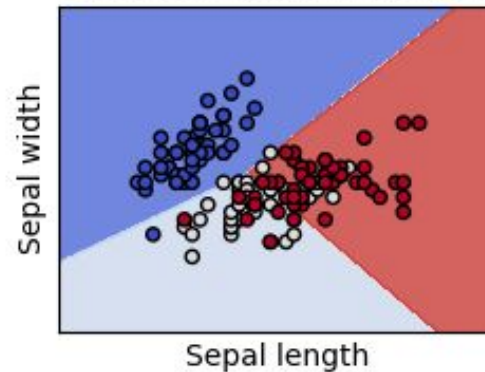# Average Workflow



**Machine Learning Algorithms Cheat Sheet**

**Unsupervised Learning: Clustering**

- k-means
- k-modes

Gaussian Mixture Model — YES — Prefer Probability — NO — Categorical Variables — Hierarchical

NO (k-means), YES (k-modes)

DBSCAN — NO — Need to Specify k — NO — Hierarchical

YES, YES

**START**

Dimension Reduction — YES — **Unsupervised Learning: Dimension Reduction**

Topic Modeling — YES — Probabilistic — YES — Latent Dirichlet Analysis

NO — Principal Component Analysis

NO — Singular Value Decomposition

NO

Have Responses — NO

YES

Predicting Numeric

**Supervised Learning: Classification**

- Linear SVM
- Naïve Bayes

Data Is Too Large — NO — Explainable — SPEED — Speed or Accuracy

YES — Naïve Bayes

YES — Decision Tree, Logistic Regression

ACCURACY — Kernel SVM, Random Forest, Neural Network, Gradient Boosting Tree

**Supervised Learning: Regression**

Speed or Accuracy — SPEED — Decision Tree, Linear Regression

ACCURACY — Random Forest, Neural Network, Gradient Boosting Tree
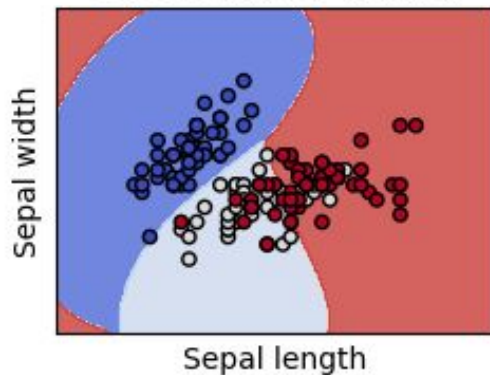
# Support Vector Machines (SVMs)
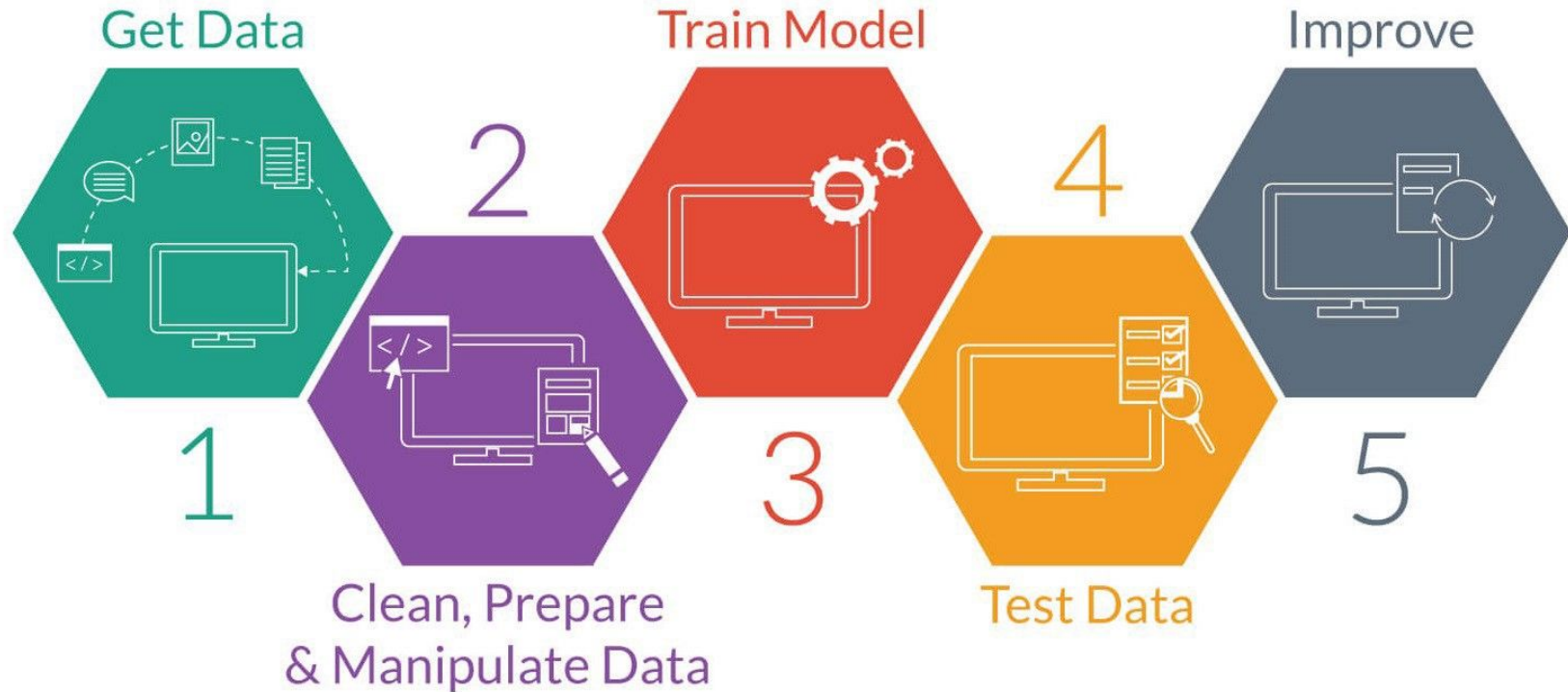
SVC with linear kernel

LinearSVC (linear kernel)

SVC with RBF kernel

SVC with polynomial (degree 3) kernel

# Typical Model Flow



**Get Data** — 1

**Clean, Prepare & Manipulate Data** — 2

**Train Model** — 3

**Test Data** — 4

**Improve** — 5

# Overview

## Attacks

- Poisoning
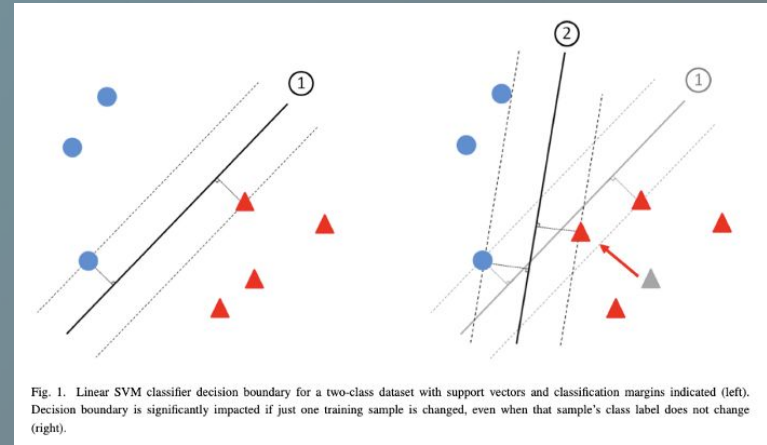- Evasion
- Trojan Attacks

## Defenses

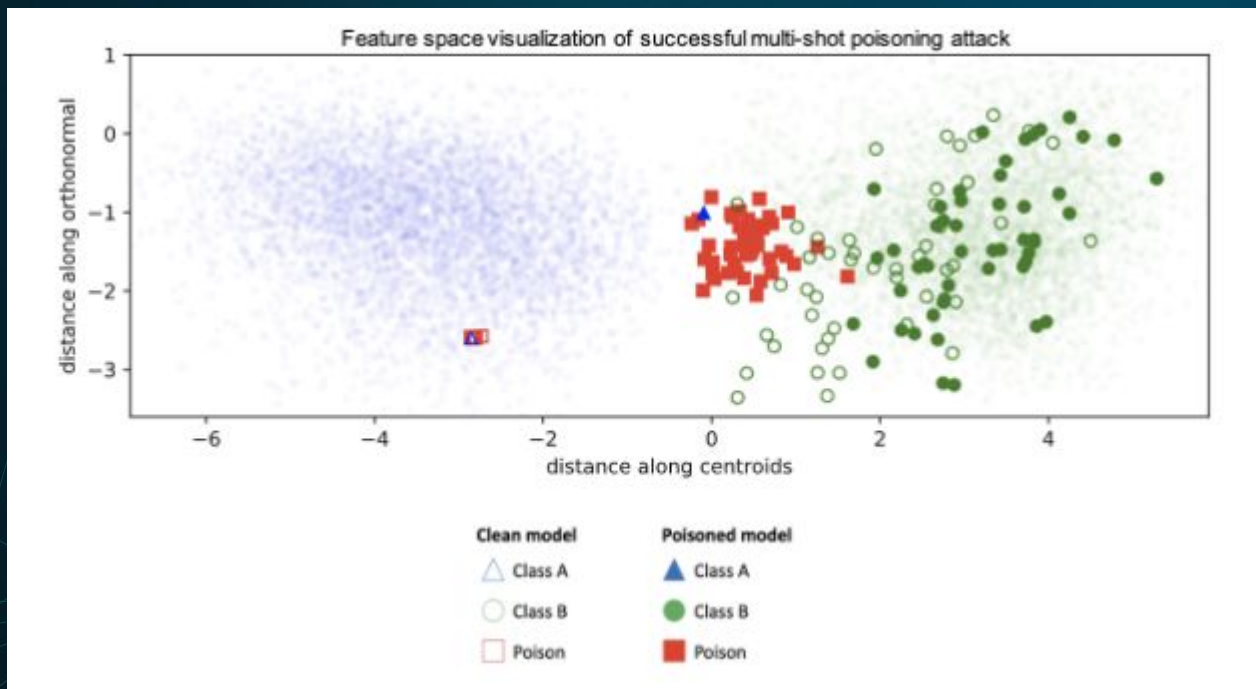- Improved Data Sanitization
- Adversarial Training
- Noise Detection

# Attacks

# Poisoning Attacks

- Performed during training (mostly)
- Goal is to ruin training data



Fig. 1. Linear SVM classifier decision boundary for a two-class dataset with support vectors and classification margins indicated (left). Decision boundary is significantly impacted if just one training sample is changed, even when that sample's class label does not change (right).

# Successful Attack



Feature space visualization of successful multi-shot poisoning attack

# Backdoor Attacks

- Input some form of unknown data to the model
- Malware detection algorithms can be a good example here

# Evasion Attacks

- Also referred to as adversarial attacks
- Performed after training, when model is in production



$x$

"panda"
57.7% confidence

$+ .007 \times$

$\text{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$

"nematode"
8.2% confidence

$=$

$\boldsymbol{x} + \epsilon\text{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$

"gibbon"
99.3 % confidence

Figure 2: A dodging attack by perturbing an entire face. Left: an original image of actress Eva Longoria (by Richard Sandoval / CC BY-SA / cropped from https://goo.gl/7QUvRq). Middle: A perturbed image for dodging. Right: The applied perturbation, after multiplying the absolute value of pixels' channels ×20.



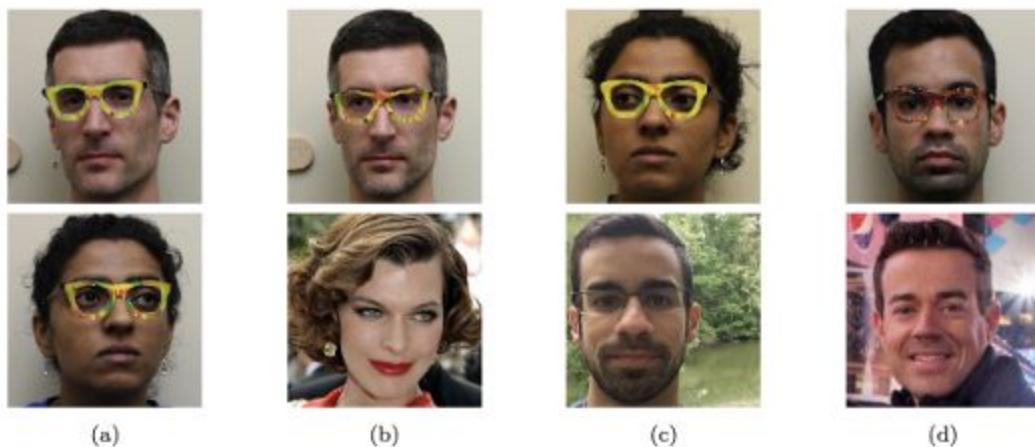Figure 5: The eyeglass frames used by $S_C$ for dodging recognition against $DNN_B$.

Figure 4: Examples of successful impersonation and dodging attacks. Fig. (a) shows $S_A$ (top) and $S_B$ (bottom) dodging against $DNN_B$. Fig. (b)–(d) show impersonations. Impersonators carrying out the attack are shown in the top row and corresponding impersonation targets in the bottom row. Fig. (b) shows $S_A$ impersonating Milla Jovovich (by Georges Biard / CC BY-SA / cropped from https://goo.gl/GlsWlC); (c) $S_B$ impersonating $S_C$; and (d) $S_C$ impersonating Carson Daly (by Anthony Quintano / CC BY / cropped from https://goo.gl/VfnDct).

# Defenses

# Noise Detection and Data Sanitization

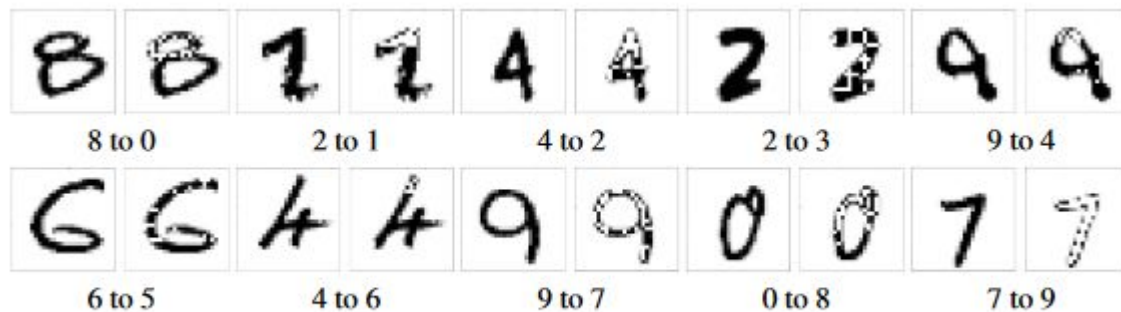- Safety Verification of Deep Neural Networks - Xiaowei Huang, et al.



Fig. 8. Adversarial examples for a neural network trained on MNIST

8 to 0  2 to 1  4 to 2  2 to 3  9 to 4

6 to 5  4 to 6  9 to 7  0 to 8  7 to 9

# Adversarial Training

- Generate adversarial data examples and retrain network to increase robustness

# 1. Train a model



Training data

Empty space between data points

Section of the data manifold you're trying to fit (top down view)

# Real World Examples

# An Example!

hopefully..

# The Future

- MIT approach
- Standard security approach

## Adversarial Examples Are Not Bugs, They Are Features

Andrew Ilyas*
MIT
ailyas@mit.edu

Shibani Santurkar*
MIT
shibani@mit.edu

Dimitris Tsipras*
MIT
tsipras@mit.edu

Logan Engstrom*
MIT
engstrom@mit.edu

Brandon Tran
MIT
btran115@mit.edu

Aleksander Mądry
MIT
madry@mit.edu

# Fin.

Questions?