

DS501: Machine learning, Part 1

Prof. Randy Paffenroth
rcpaffenroth@wpi.edu

Worcester Polytechnic Institute

2014



WPI

Learning **objectives** for this machine learning class.

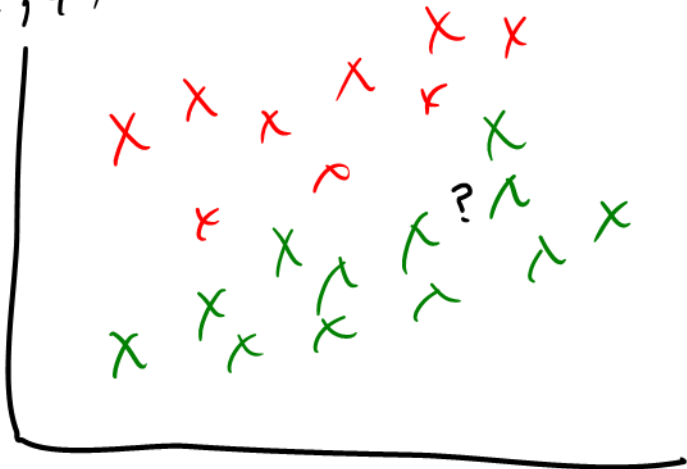
- Overview of machine learning
- Supervised classification.
 - K-nearest neighbors
 - Support vector machines
- Learn some Python packages, including:
 - scikit-learn

Inspired by: http://scikit-learn.org/stable/tutorial/statistical_inference/supervised_learning.html

Basic definitions: Classification vs. Regression

Classification

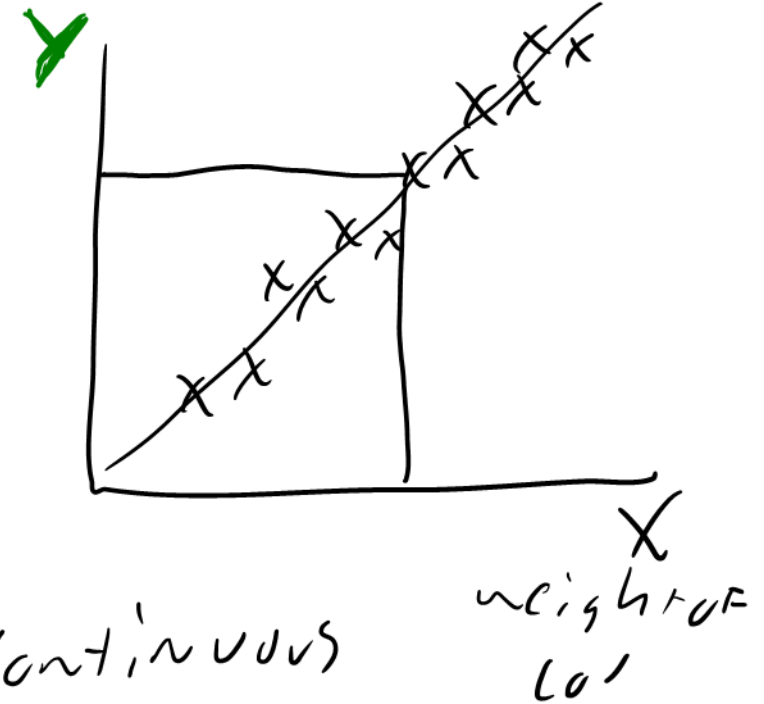
weight



Height

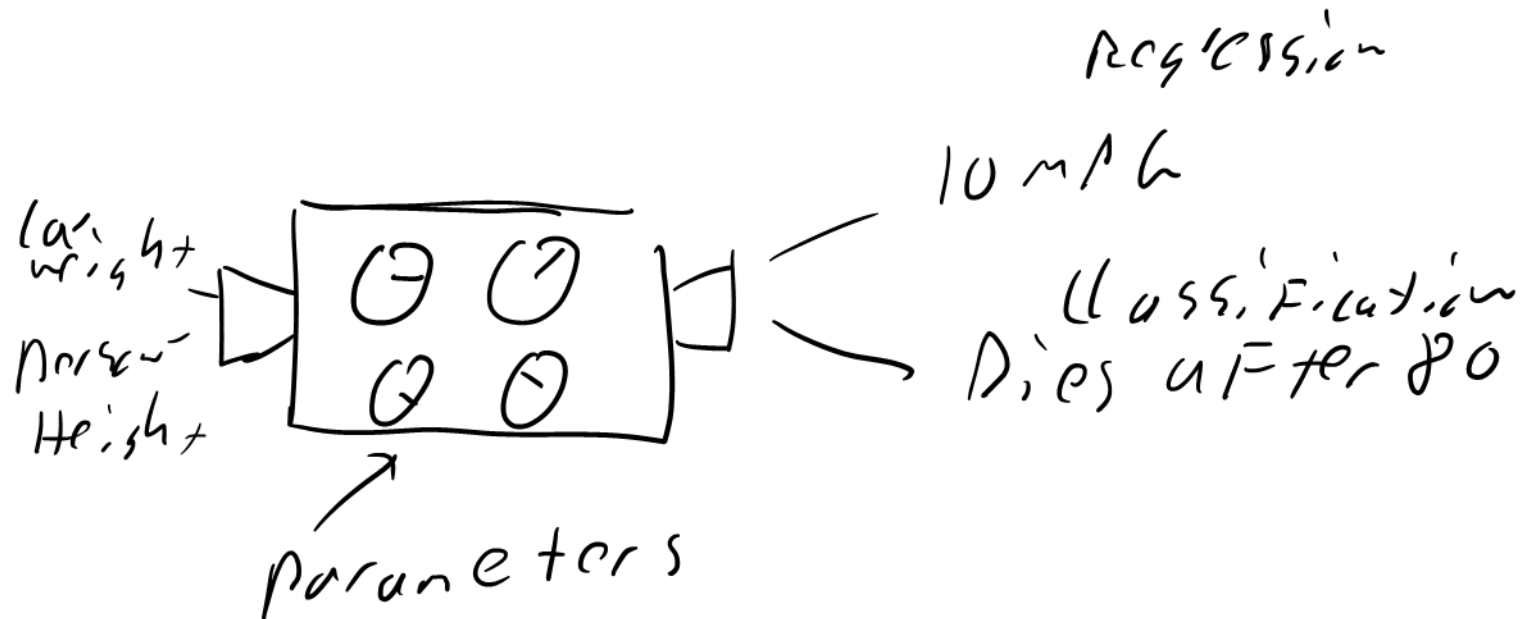
predictions
discrete
attributes
(features
responses)

Regression
continuous



continuous
Response

weight of
load

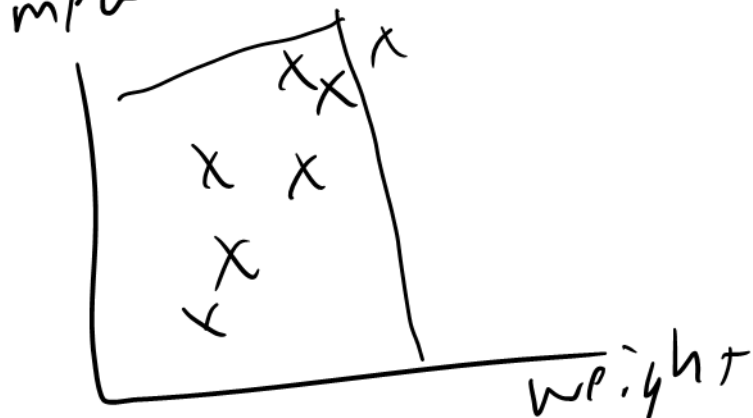


Data is used
to get the
Dials

Basic definitions: Supervised vs. Unsupervised Learning

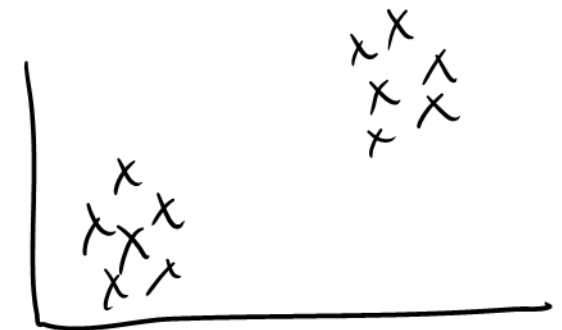
supervised

You have
examples to
learn from



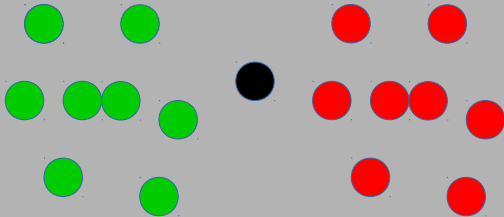
unsupervised

No examples
to learn from!

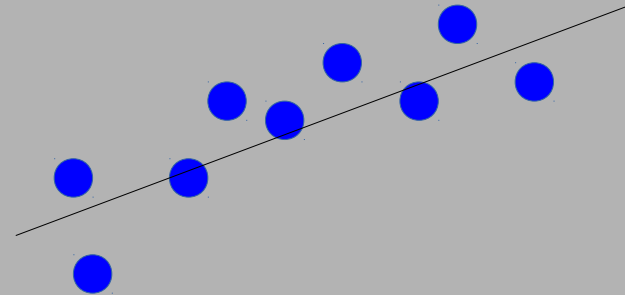


We have four days we will cover four topics.

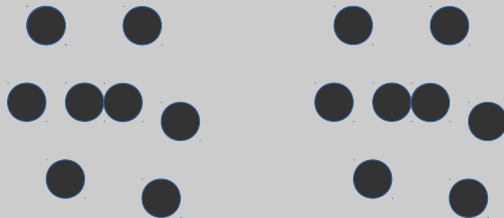
Supervised Classification



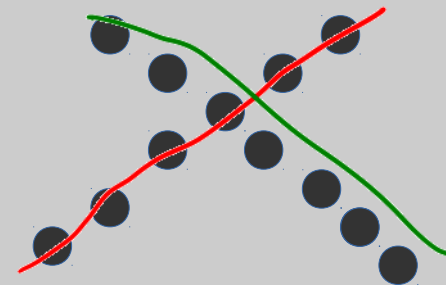
Supervised Regression



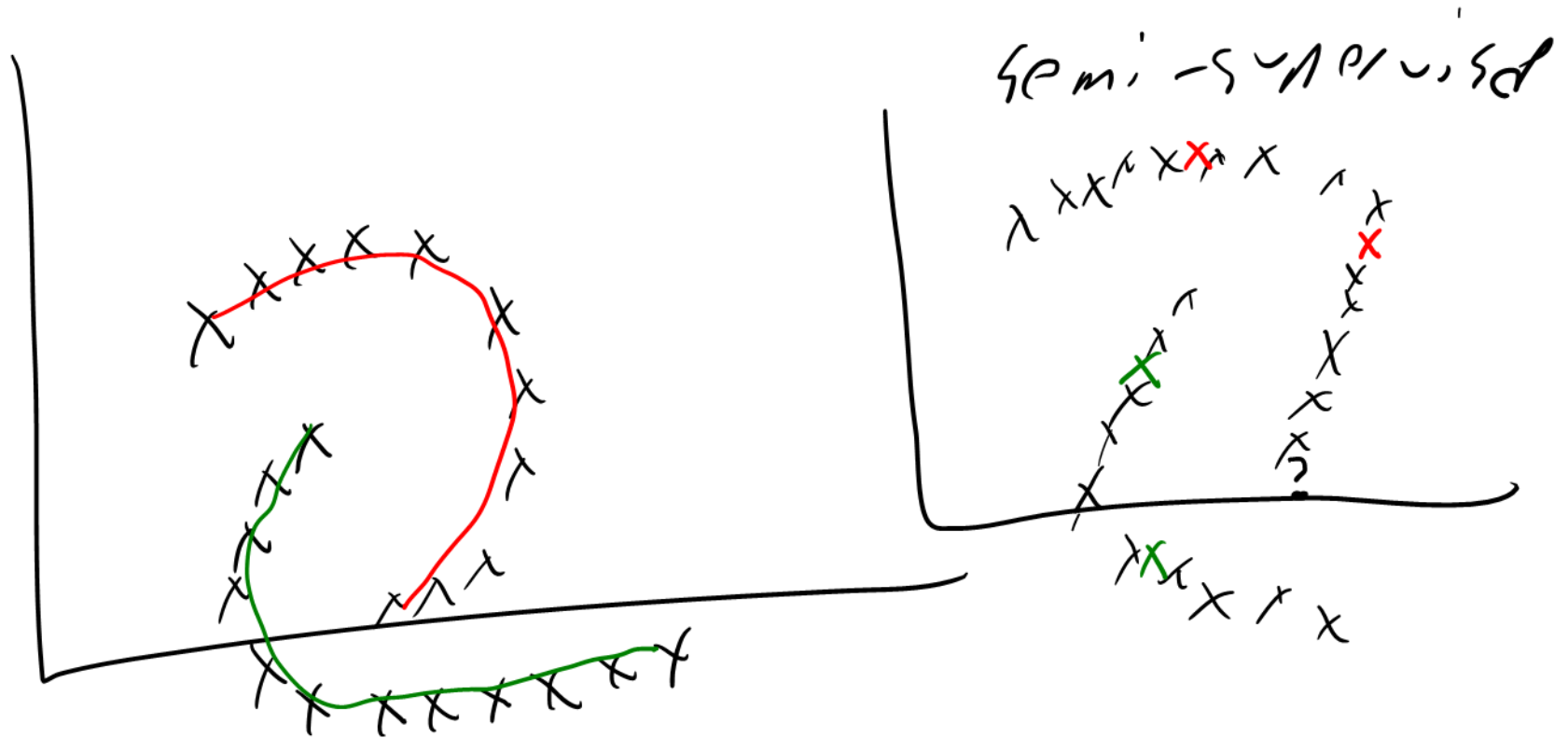
Unsupervised Clustering



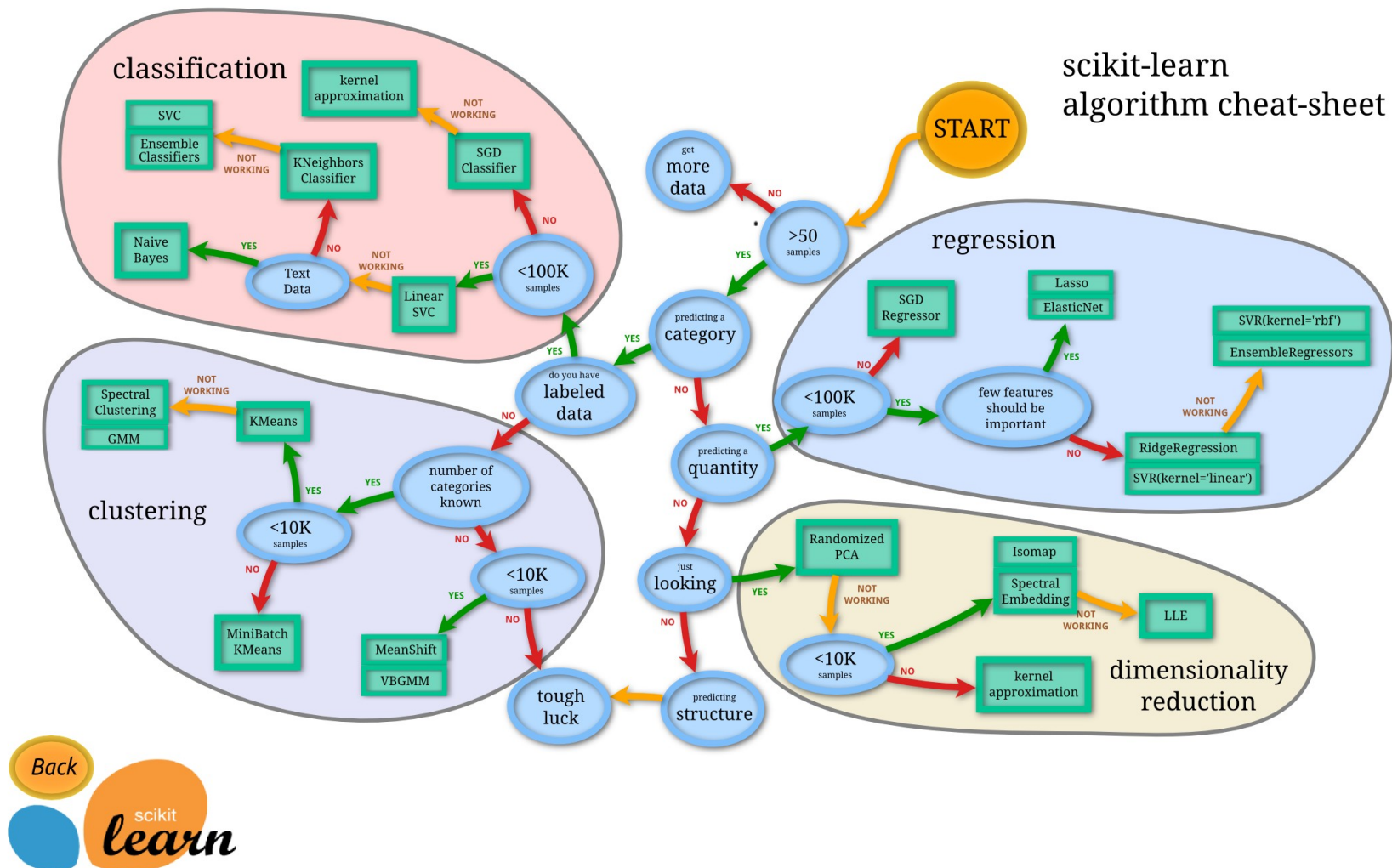
Dimension Reduction



What about Unsupervised Regression?



scikit-learn



Iris data set

Features:

sepal length (cm)
sepal width (cm)
petal length (cm)
petal width (cm)

"Iris virginica" by Frank Mayfield - originally posted to Flickr as Iris virginica shrevei BLUE FLAG. Licensed under Creative Commons Attribution-Share Alike 2.0 via Wikimedia Commons - http://commons.wikimedia.org/wiki/File:Iris_virginica.jpg#mediaviewer/File:Iris_virginica.jpg



Catagories:

setosa
versicolor
virginica

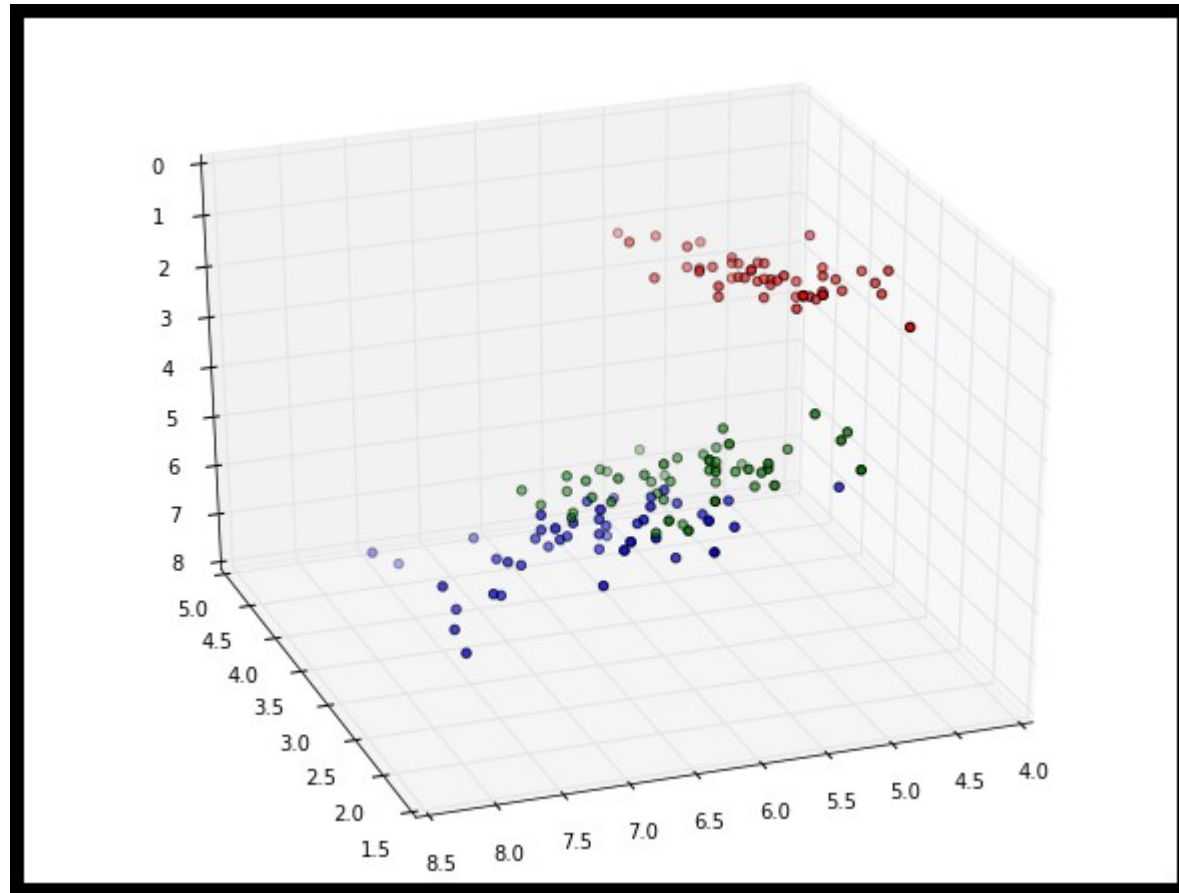


"Kosaciec szczecinkowaty Iris setosa". Licensed under Creative Commons Attribution-Share Alike 3.0 via Wikimedia Commons - http://commons.wikimedia.org/wiki/File:Kosaciec_szczecinkowaty_Iris_setosa.jpg#mediaviewer/File:Kosaciec_szczecinkowaty_Iris_setosa.jpg



"Iris versicolor 3". Licensed under Creative Commons Attribution-Share Alike 3.0 via Wikimedia Commons - http://commons.wikimedia.org/wiki/File:Iris_versicolor_3.jpg#mediaviewer/File:Iris_versicolor_3.jpg

Let's look at it in Python



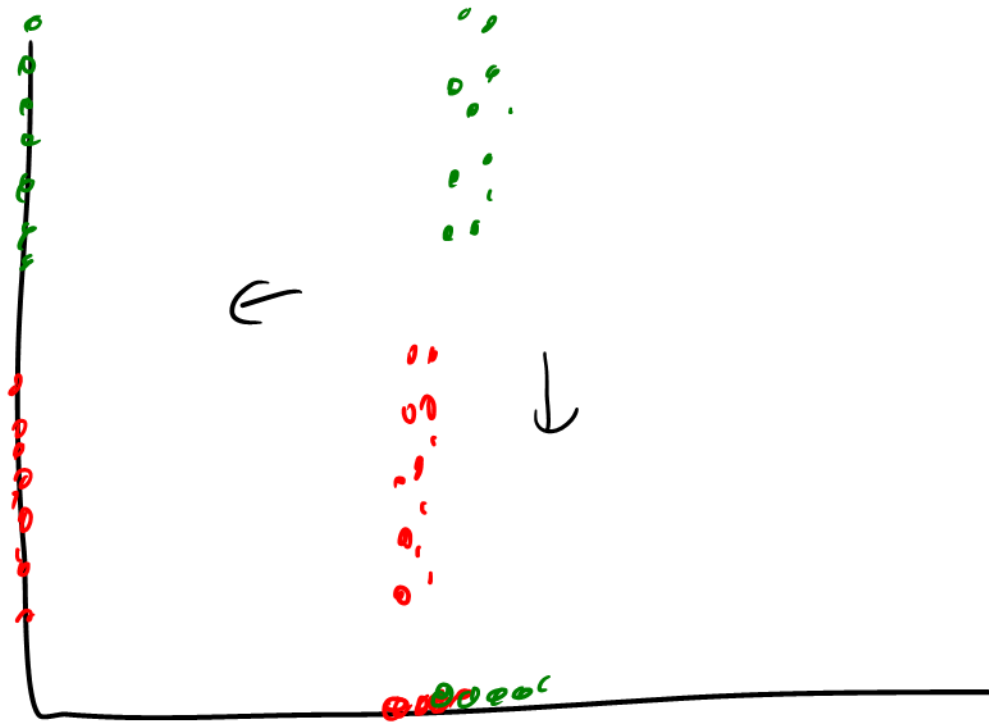
What is PCA?

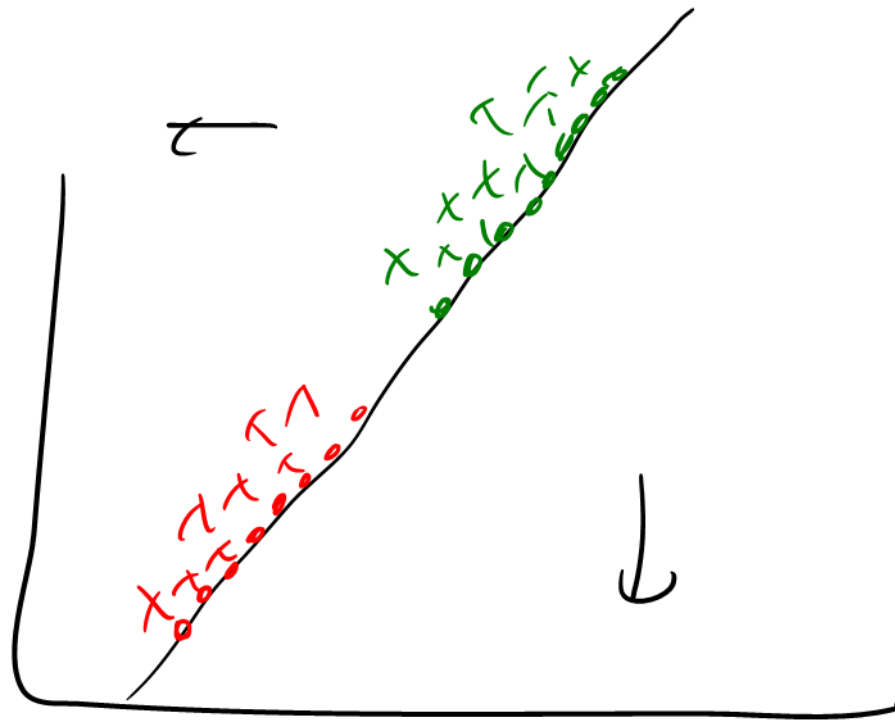
- Principle Component Analysis
 - Commonly used tool for visualization and data pre-processing.

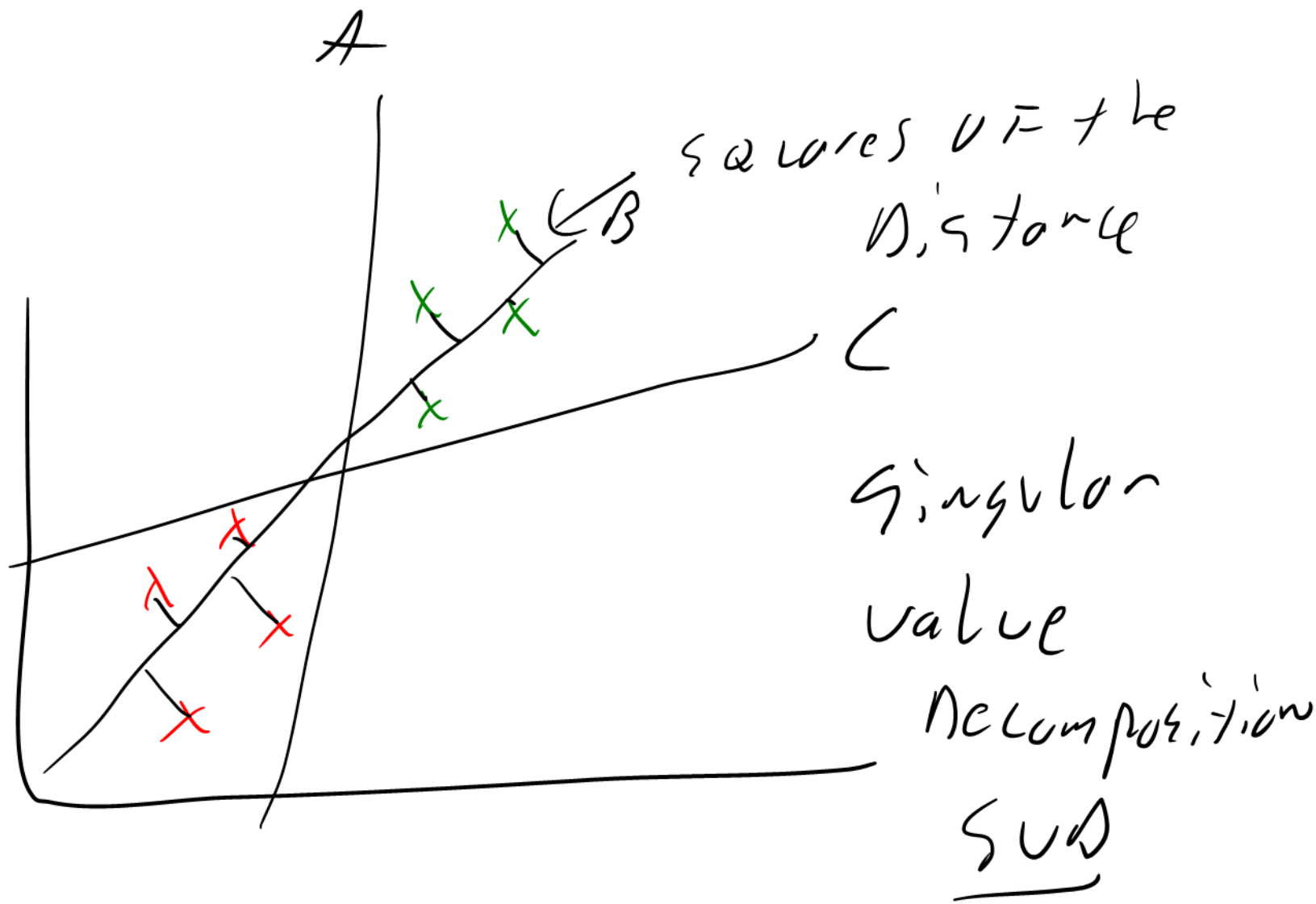


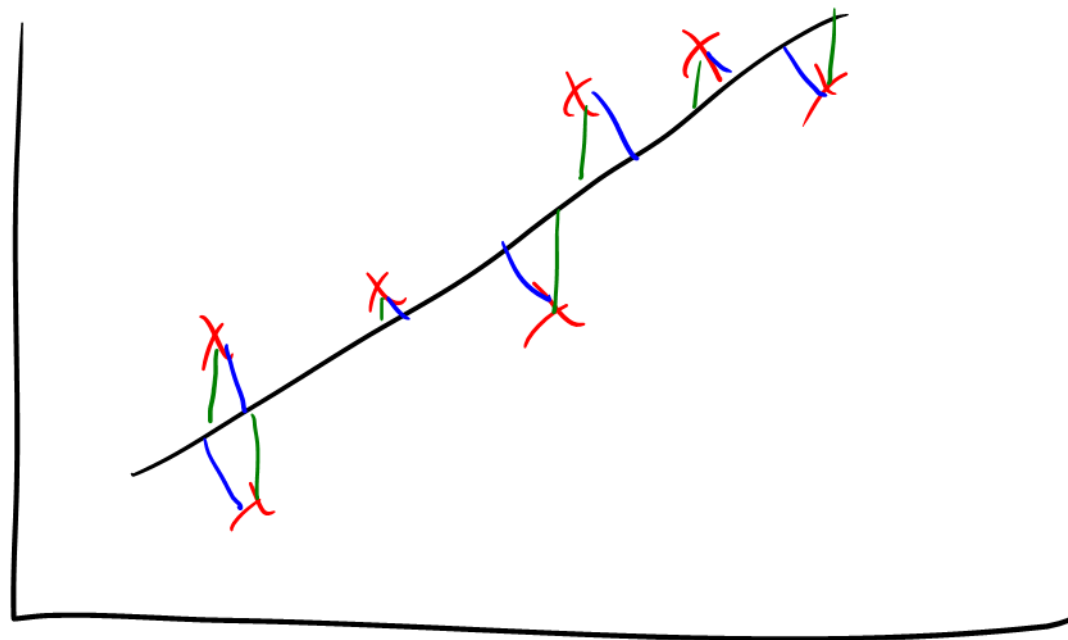
Key idea of PCA

maximize the
variance of
the projection!

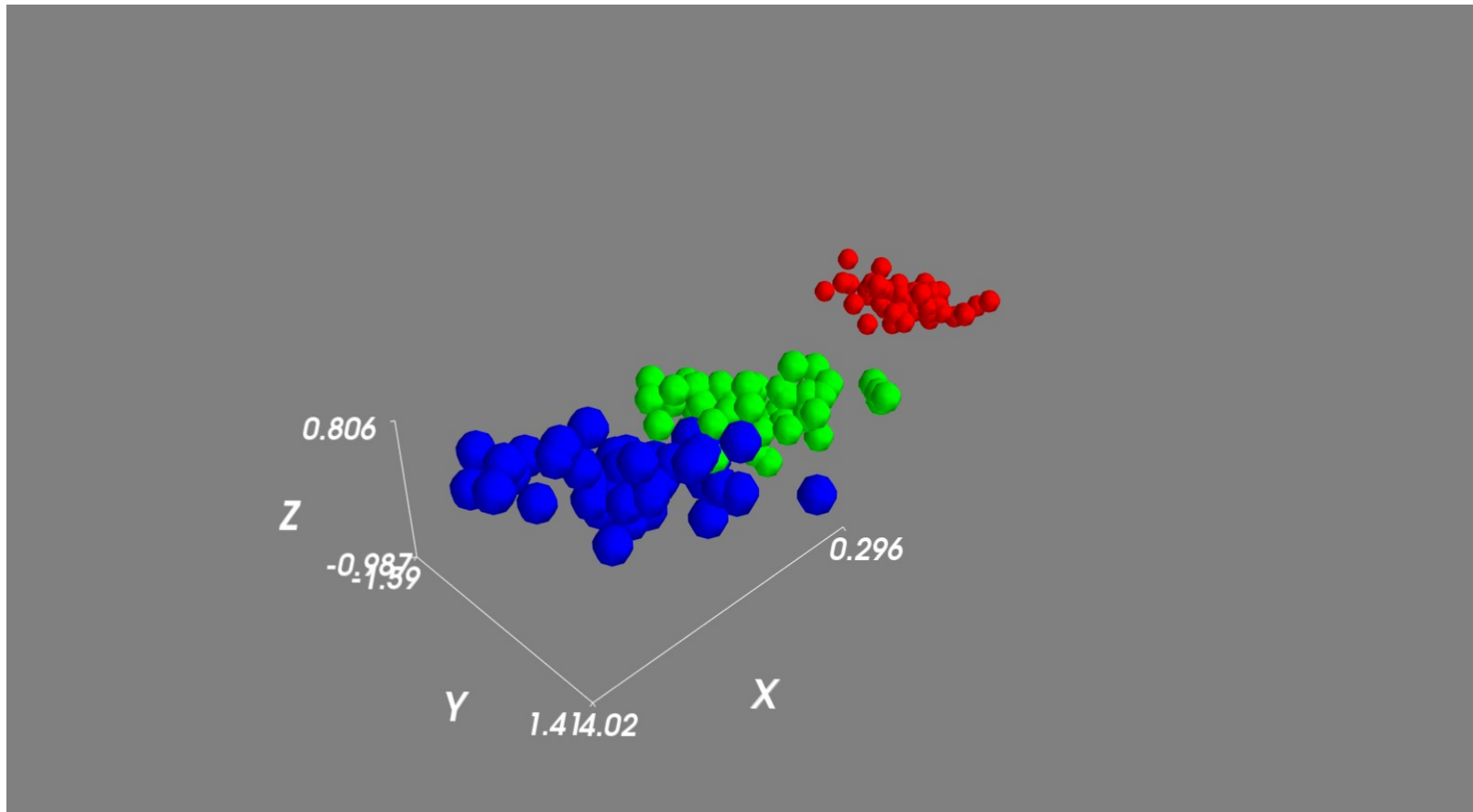


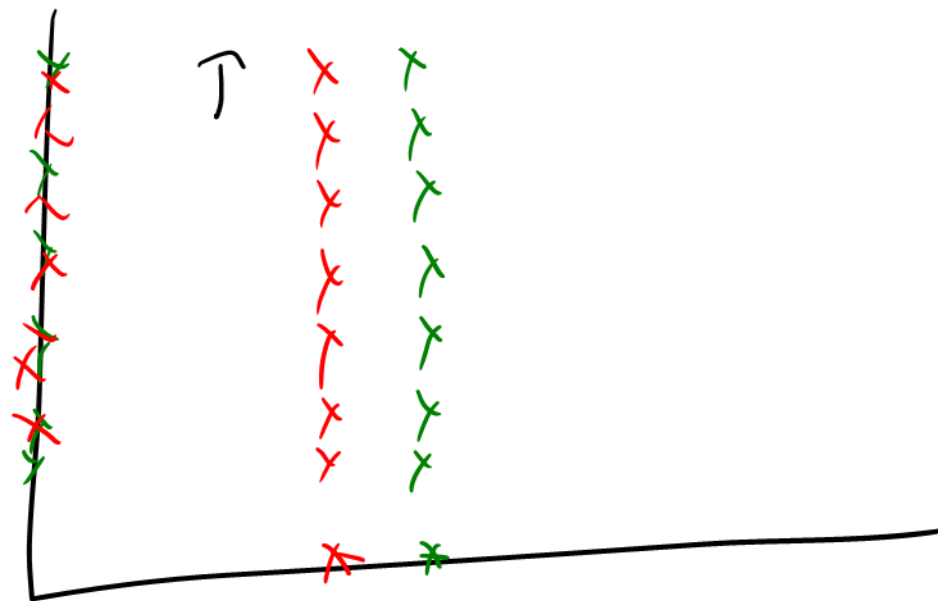






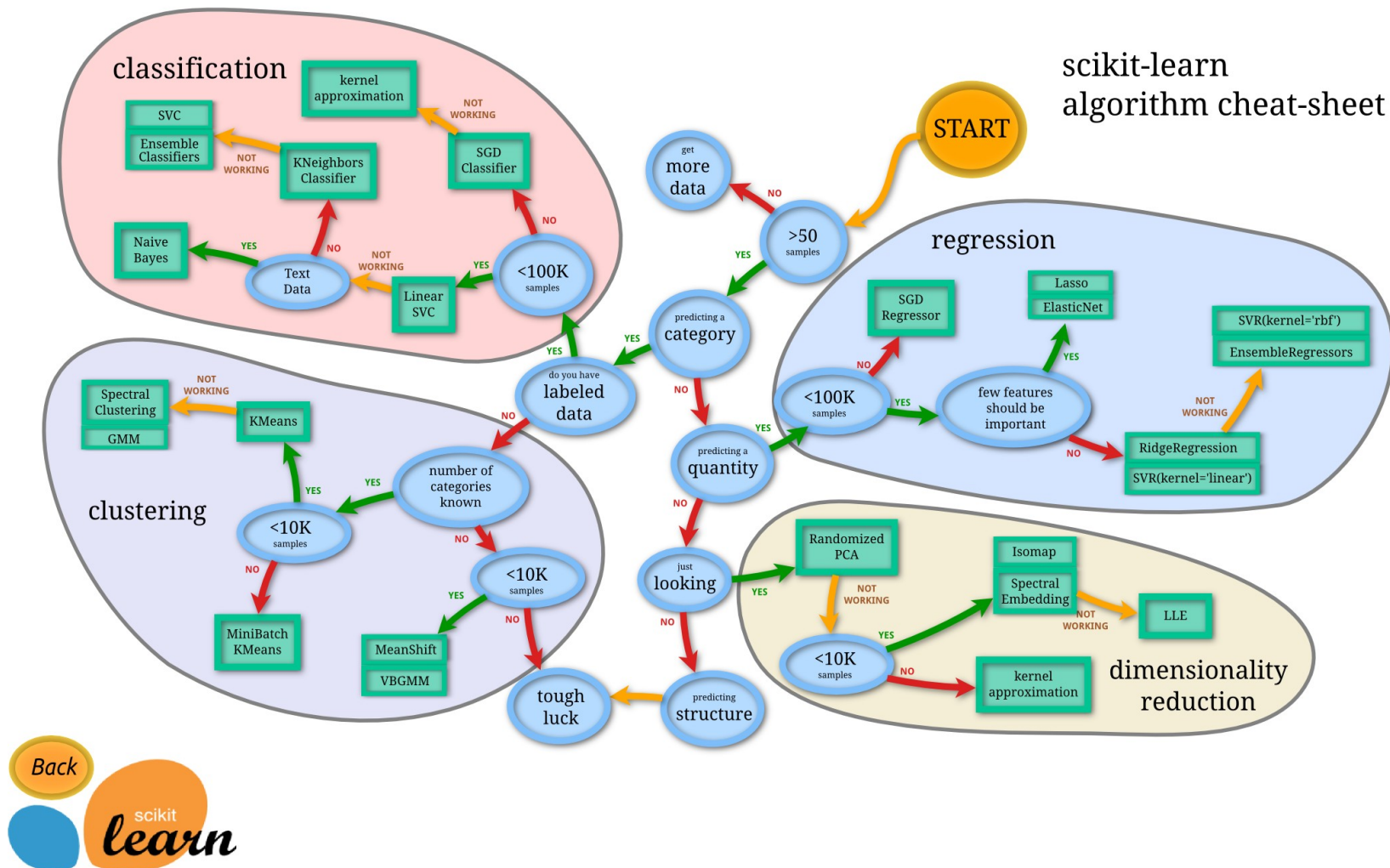
Lets take a look at it in Python





scikit-learn

scikit-learn
algorithm cheat-sheet



http://scikit-learn.org/stable/tutorial/machine_learning_map/index.html

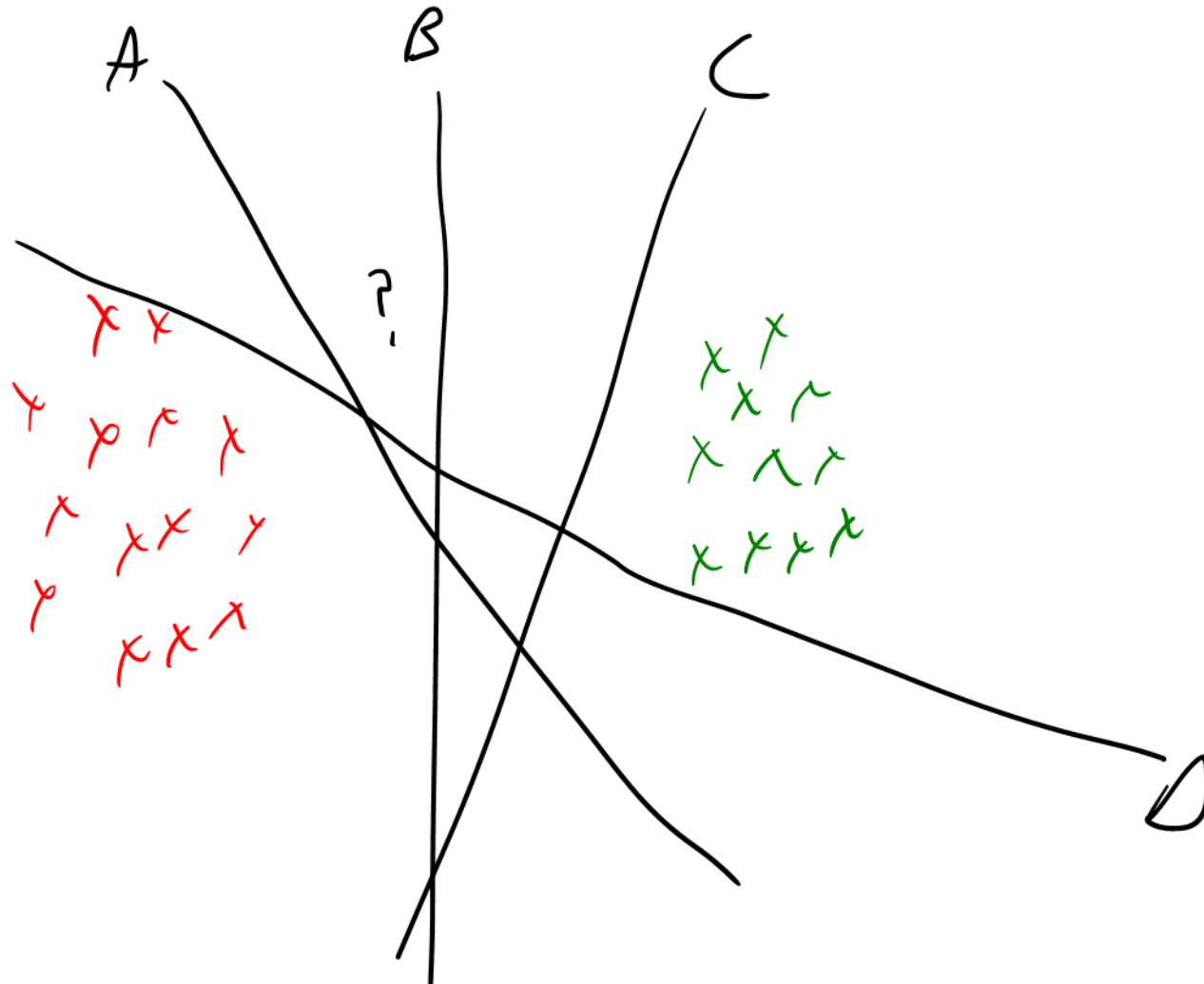


WPI

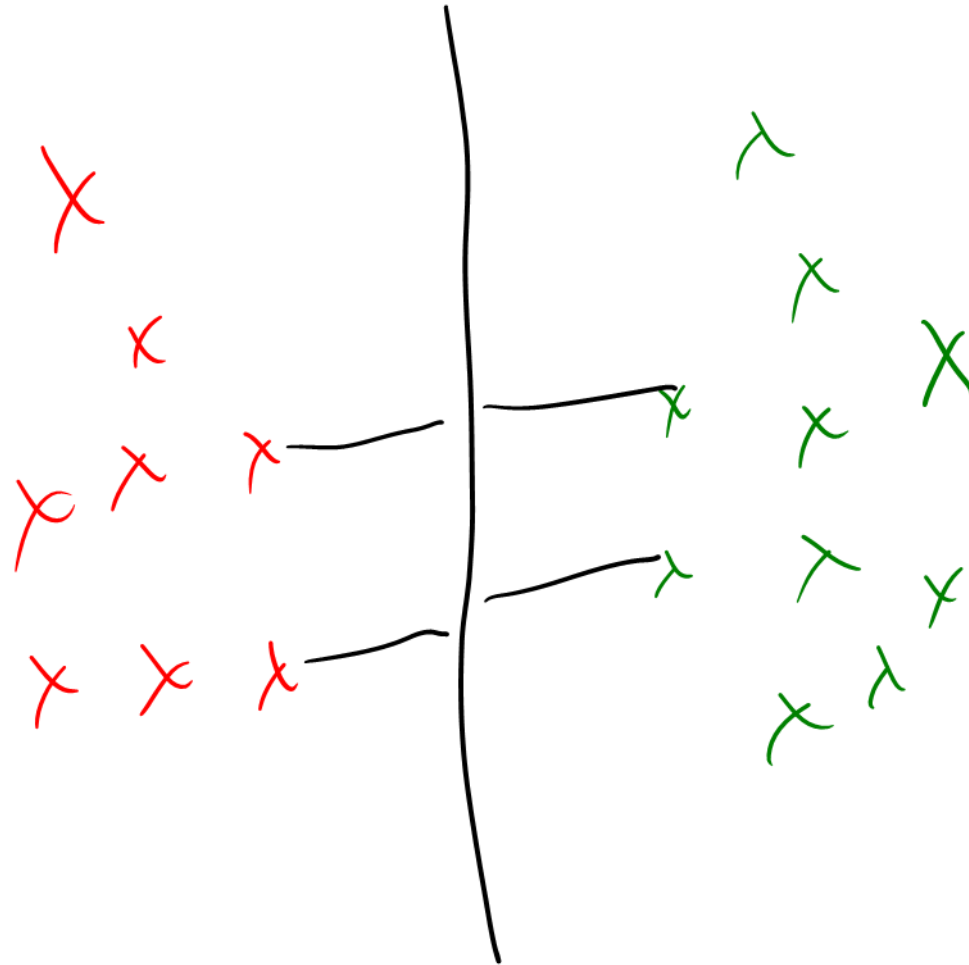
What is Linear Support Vector Machine (SVM)?

- Maximum margin classifier
 - Computes a linear “decision boundary” that splits the data into two regions.
 - Allows one to predict a classification of a point based upon which side of the decision boundary it lay on.

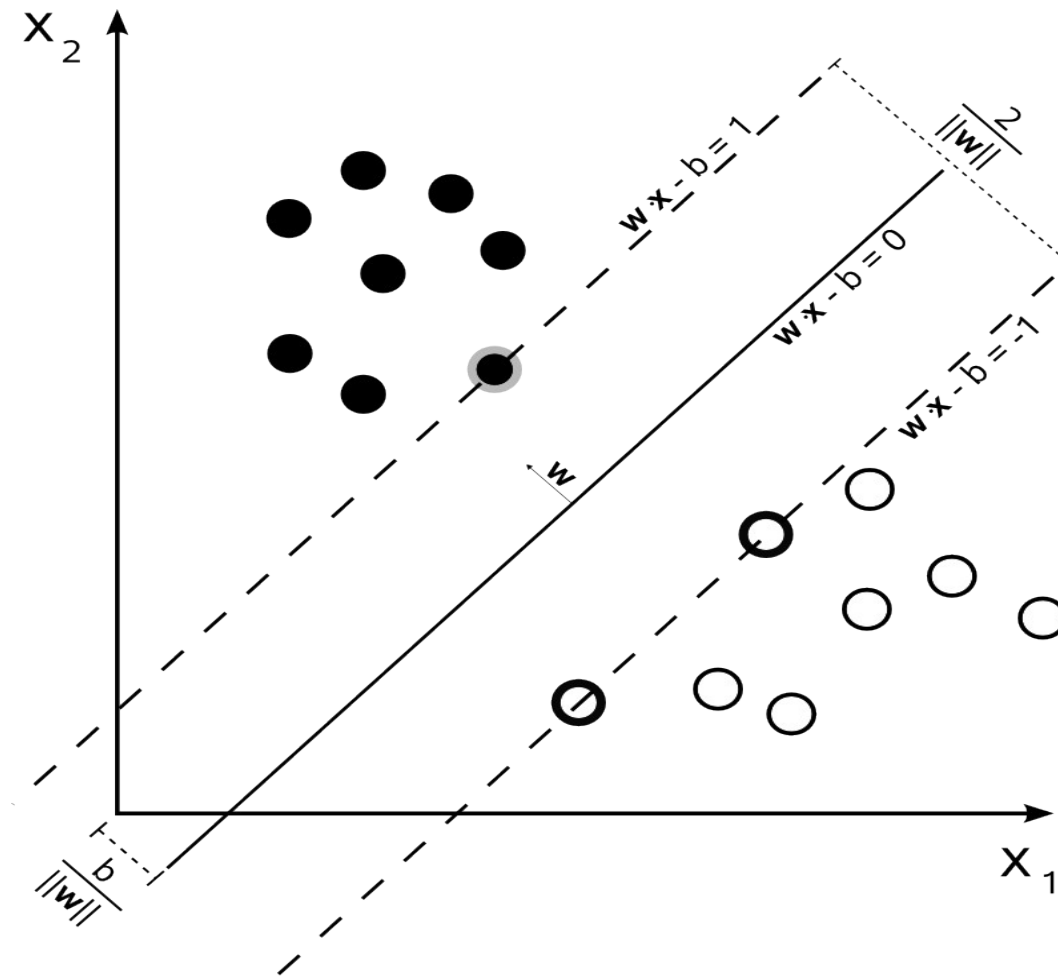
Let's derive SVM



maximize the minimum distance



SVM

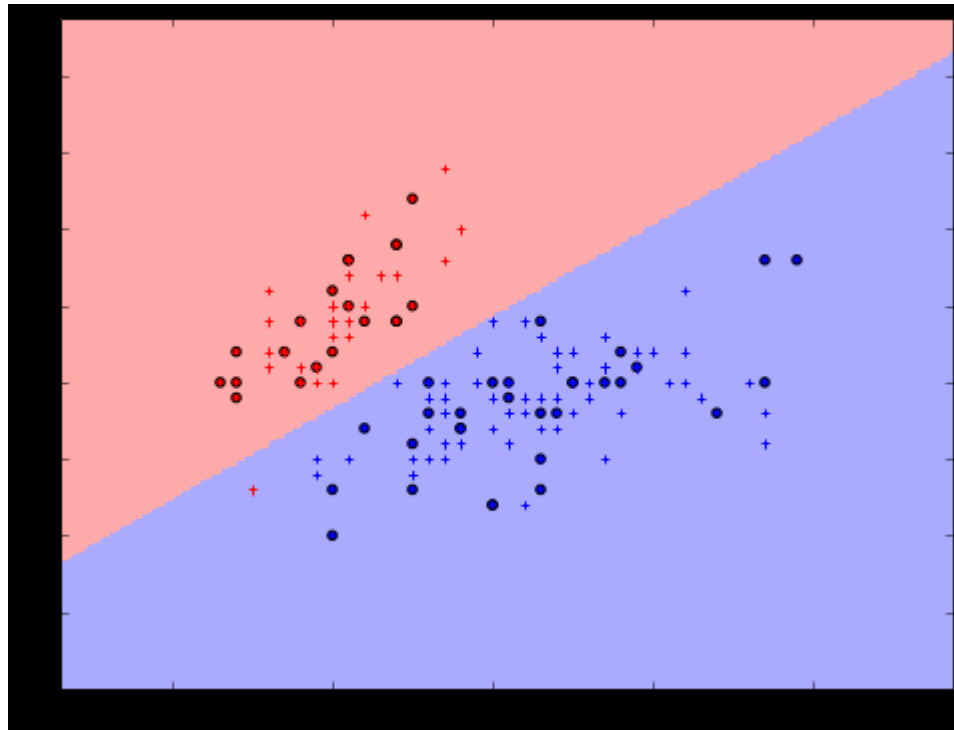


"Svm max sep hyperplane with margin" by Cyc - Own work. Licensed under Public domain via Wikimedia Commons - http://commons.wikimedia.org/wiki/File:Svm_max_sep_hyperplane_with_margin.png#mediaviewer/File:Svm_max_sep_hyperplane_with_margin.png



WPI

Let's look at SVM in Python



But wait! Training vs. testing!

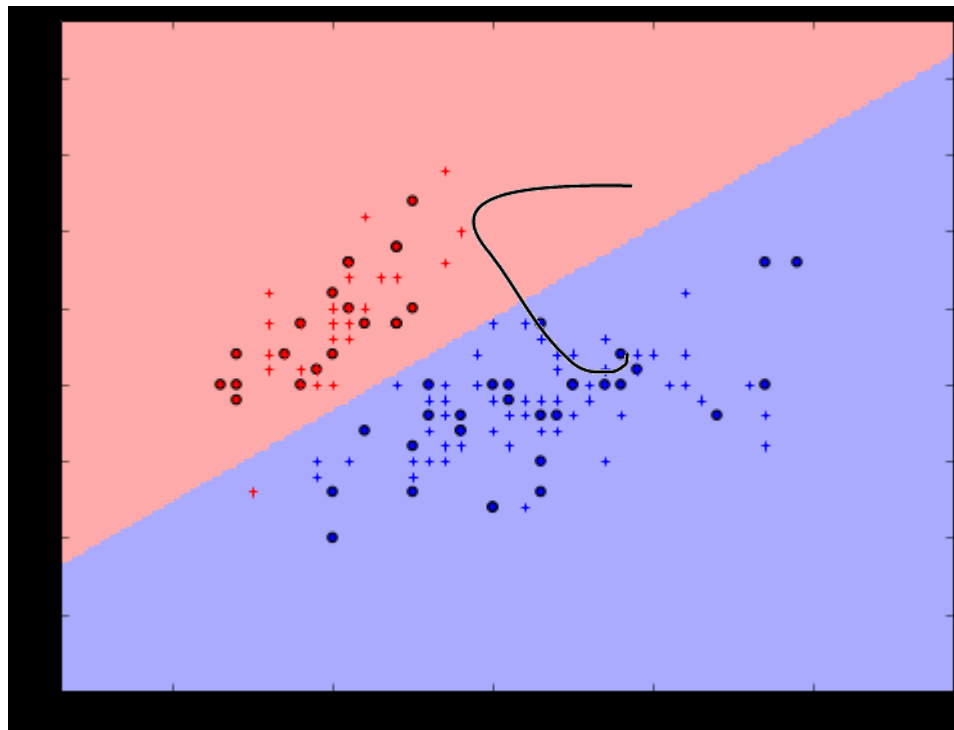
some data For your algorithm

some data For testing your
Algorithm

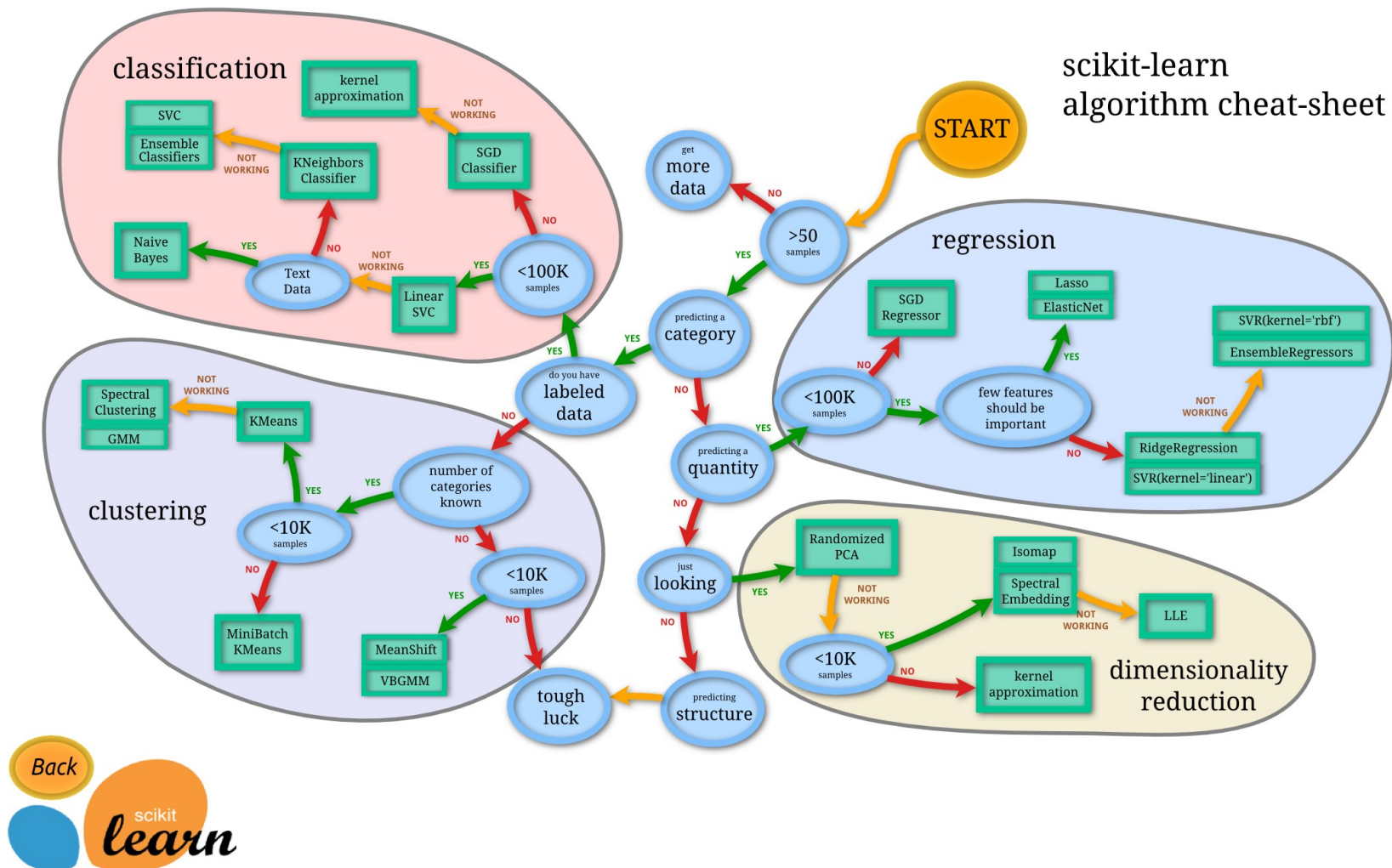
Fundamental

avoid overfitting

Back to Python



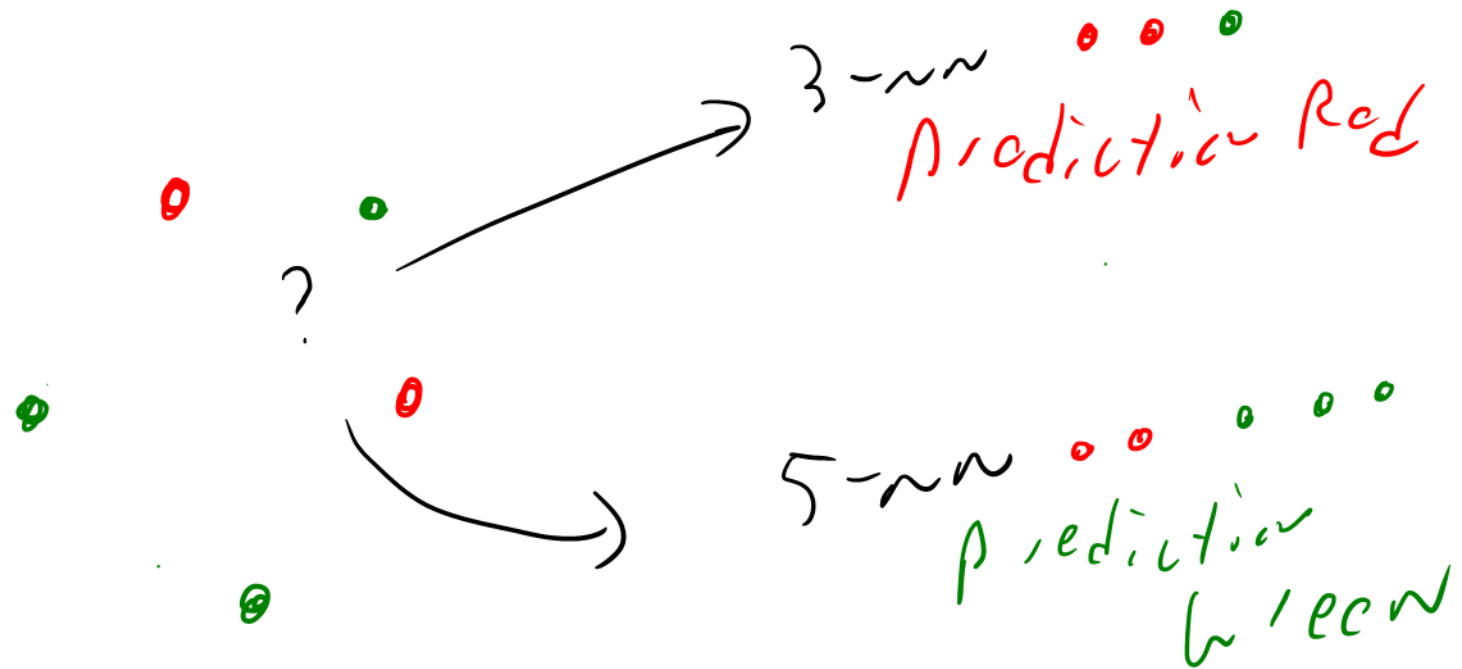
scikit-learn



What is K-NN?

- K-nearest neighbors
- Another common classification algorithm
 - Perhaps the most common

Let's try an example



Let's try K-NN in Python

