

Introduction to Machine Learning

By:

Shyamal Vadera

Agenda

- Introduction
- Basics
- Regression
- Classification
- Clustering
- Decision Trees

In other words...

“Machine learning is a set of techniques, which help in dealing with vast data in the most intelligent fashion (by developing algorithms or set of logical rules) to derive actionable insights (delivering search for users in this case).”

How ML is different from
Simple if-else rules?

Use-Cases

- Spam Email Detection
- Machine Translation (Language Translation)
- Image Search (Similarity)
- Amazon Recommendations : Clustering
- Google News : Classification

Terminology

- Samples
- Features
- Feature vector

Data about Movies

Name	Actors	Director	Budget	Type	Rating
Lagan	Aamir khan	T.R ramana	25 Cr	Drama	7.2
Wanted	Salman khan	Ali Zafar	50 Cr	Action	6.2
3 idiots	Aamir khan	Rajiv irani	35 Cr	Drama	8.4
Queen	Kangana	Abhi chopra	15 Cr	Drama	7.9
Neerja	Sonam kapur	Niraj pandey	20 Cr	Thriller	7.6
Dangal	Aamir khan	Nitesh Tiwari	70 Cr	Sport	8.9
Drishyam	Ajay Devgn	Rahul Prasad	30 Cr	Thriller	8.5

- Features

Name	Actors	Director	Budget	Type	Rating
------	--------	----------	--------	------	--------

- Feature Vector

[178, 5, 2, 25, 1, 7.2]

- Matrix Representation of Data

Let's dig deep into it...

What do you mean by

Apple

Learning



Features:

1. Color: **Radish/Red**
 2. Type : **Fruit**
 3. Shape
- etc...



Features:

1. Sky Blue
 2. **Logo**
 3. Shape
- etc...

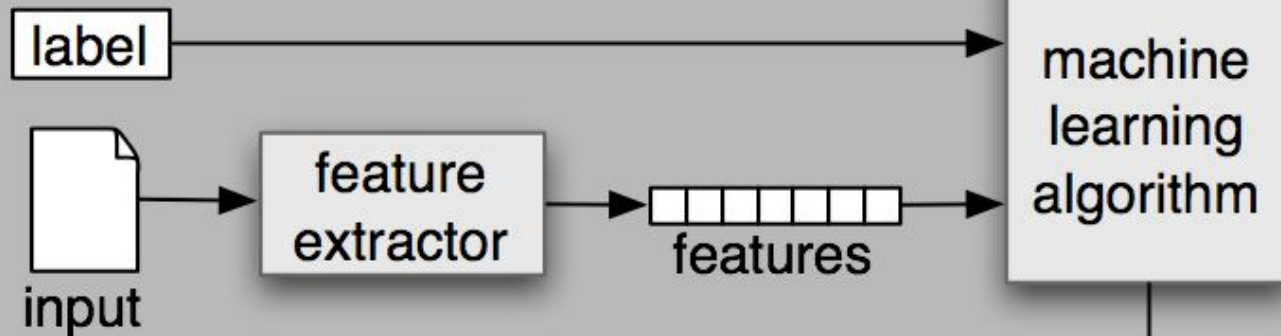


Features:

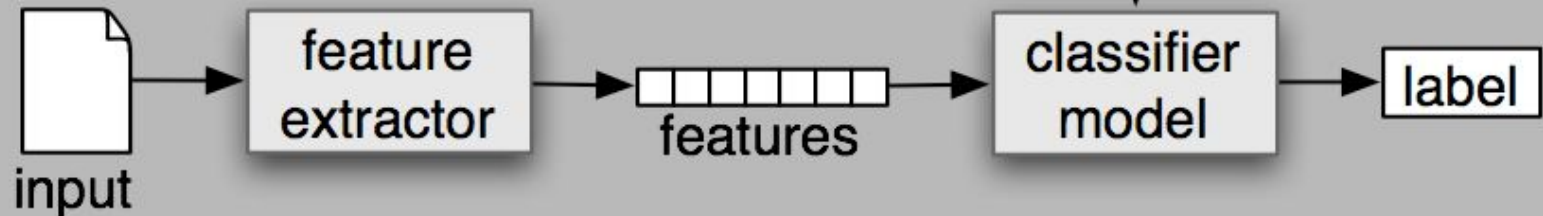
1. **Yellow**
 2. **Fruit**
 3. Shape
- etc...

Workflow

(a) Training



(b) Prediction



Categories

1. Supervised Learning
2. Unsupervised Learning
3. Reinforcement Learning

Supervised Learning

- Predictive Models.
- The correct classes of the training data are known.
- Algorithms: Nearest Neighbor, Naïve Bayes, Decision Trees, Regression.

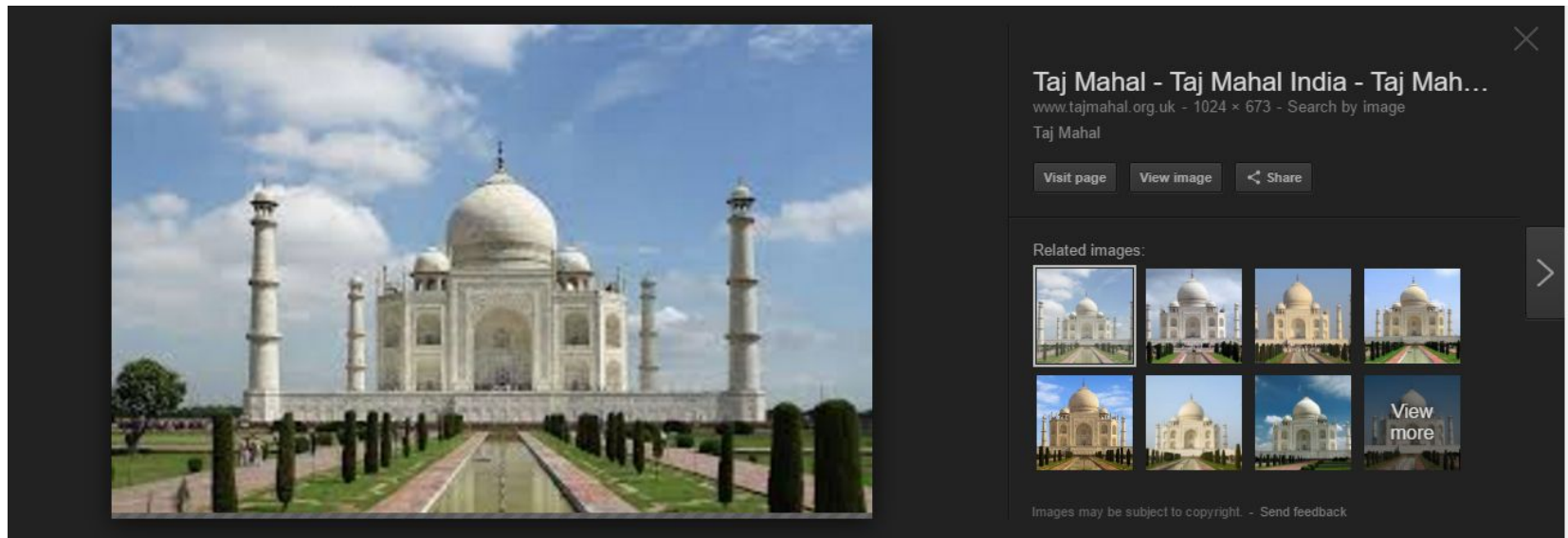
Labels

Name	Actors	Director	Budget	Type	Rating
Lagan	Aamir khan	T.R ramana	25 Cr	Drama	7.2
Wanted	Salman khan	Ali Zafar	50 Cr	Action	6.2
3 idiots	Aamir khan	Rajiv irani	35 Cr	Drama	8.4
Queen	Kangana	Abhi chopra	15 Cr	Drama	7.9
Neerja	Sonam kapur	Niraj pandey	20 Cr	Thriller	7.6
Dangal	Aamir khan	Nitesh Tiwari	70 Cr	Sport	8.9
Drishyam	Ajay Devgn	Rahul Prasad	30 Cr	Thriller	8.5

Unsupervised Learning

- The correct classes of the training data are not known.
- You are trying to find out the structure or relationship between different inputs.
- Algorithms: K-means clustering

Example



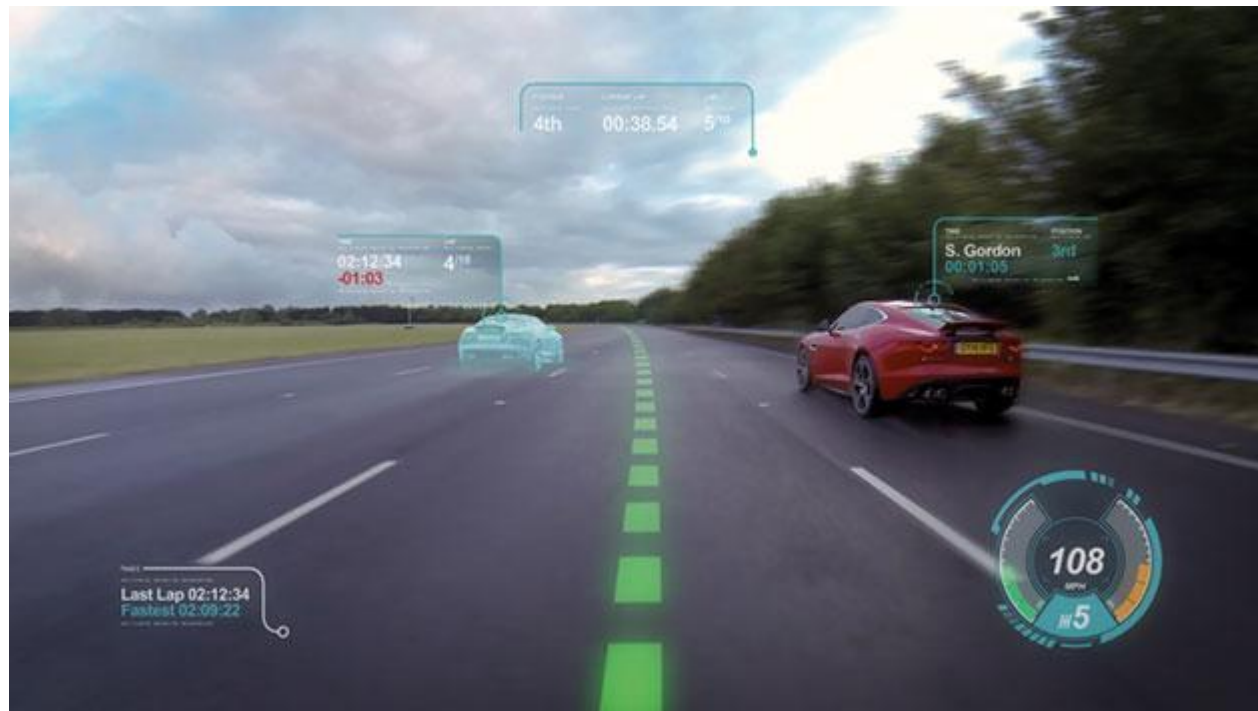
Example

Name	Actors	Director	Budget	Type
Lagan	Aamir khan	T.R ramana	25 Cr	Drama
Wanted	Salman khan	Ali Zafar	50 Cr	Action
3 idiots	Aamir khan	Rajiv irani	35 Cr	Drama
Queen	Kangana	Abhi chopra	15 Cr	Drama
Neerja	Sonam kapur	Niraj pandey	20 Cr	Thriller
Dangal	Aamir khan	Nitesh Tiwari	70 Cr	Sport
Drishyam	Ajay Devgn	Rahul Prasad	30 Cr	Thriller

Reinforcement Learning

- Learning from feedback from the environment.
- Good example of this is Self driving car,
 - It decide continuously from it's environment about which route to take? What speed to drive on?
- Algorithm: Markov Decision Process

Example of RL



Machine Learning Techniques

Techniques

- **Classification**: predict class from observations
- **Regression (prediction)**: predict value from observations
- **Clustering**: group observations into “meaningful” groups

Classification

- Classify a document into a predefined category.
- Documents can be text, images etc.

Spam detection

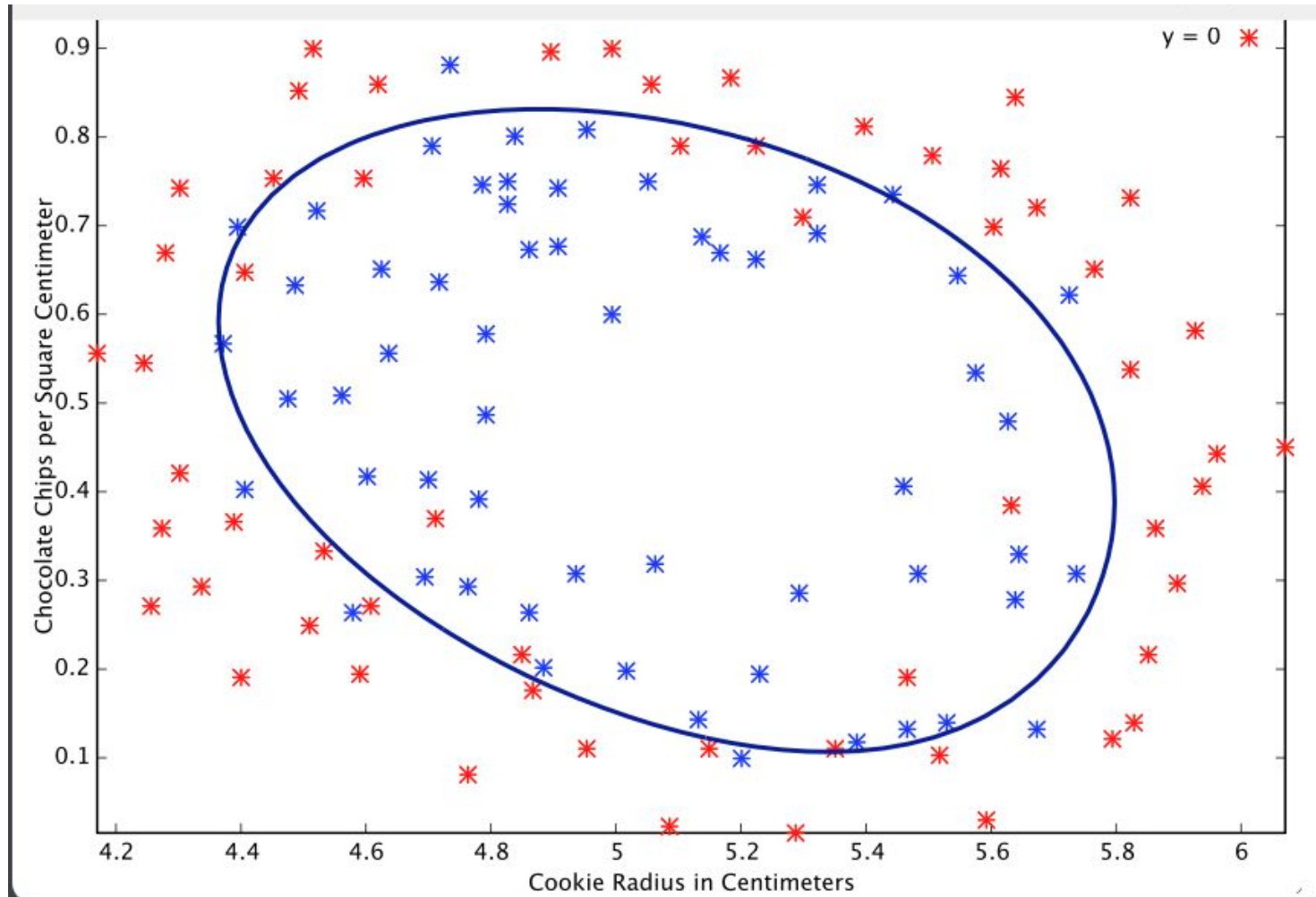
[Delete all spam messages now](#) (messages that have been in Spam more than 30 days will be automatically deleted)

<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	HDFC Bank	LOAN upto Rs 25 lac - Disbursal in 2 days - open in fresh tab -- If you do not want to receive any more newsletters, please click here	9:40 pm
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	iEntry	Welcome iEntry Member - Ultimate Guide To Assessing	9:23 pm
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	New-Zealand-Jobs.067L	Come to New Zealand to find a great job and settle here (Search for all Jobs from diffe... - Search for all Jobs from different kinds of industries Find a Job in Enchanting	8:18 pm
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	CarSizzler	Assured Free Luxurious Ride worth Rs.300 with Uber Cabs - Home Home Buy New Car Buy New Car Sell Car Sell Car Tech Tics Tip & Tale Facebook 41727 others	6:05 pm
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	Supermarket Promotion	Enjoy Rs.1700 voucher valid at any supermarket! - If you are unable to view this mailer Click here HOW TO CONTACT US? BY EMAIL: support@savethedeals.in	4:51 pm
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	Entireweb Newsletter	Hire an SEO the Right Way -- 6 Tips You Must Remember for Life - Unsubscribe me View web version Become a fan on Facebook Follow us on Twitter September 5th, 21	1:24 pm
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	Max Bupa	A policy that understands your family's medical need - open in fresh tab -- If you do not want to receive any more newsletters, please	11:08 am
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	Scoop.it	Your Scoop.it Daily Summary - How to Maximize Your LinkedIn Publishing Exposure SME a... - Scoop.it Facebook Twitter G+ H	9:30 am
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	standard charterer Bank	Instant approval on your Credit Card	7:27 am
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	CAR TRADE	Sell your car at no cost at all - If you are having trouble viewing this email, view web version View this message in your mobile	Sep 4
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	Uday	VPS Web Hosting Services Provider - Dear Sir, I am Uday Sharma, Business development executive. We are providing quality VPS hosting for	Sep 4
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	Mark Regan, SPN	How to Find Your Most Valuable Keywords [Free Guide] - This is a SiteProNews/ExactSeek Webmaster Exclusive Mailing! To drop your subscription, use the link	Sep 4
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	HDFC Bank	LOAN upto Rs 25 lac - Disbursal in 2 days - open in fresh tab You have received this mailer from Shop@Best on behalf of HDFC Bank because you	Sep 4
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	CAR TRADE	Sell your car at no cost at all - open in fresh tab -- If you do not want to receive any further newsletters, please click here	Sep 4
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	ICICI Bank	Home Loan Interest Rate starting from 10.15%*. Get Instant Approval! - open in fresh tab -- If you do not want to receive any more newsletters, please Click Here	Sep 4
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	calculateyourwealth	It's good when your bank helps you manage your wealth and fulfill your ambitions - Calculate Now Dreams you wish to realize in your lifetime require enough wealth. C	Sep 4
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	Angel Broking	Get Low Brokerage - Free Demat & Trading Account - open in fresh tab -- If you do not want to receive any further newsletters, please click here	Sep 4
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	Bankbazaar	7 Minute Instant Online Approval for your PESONAL LOAN - Now get instant online Personal Loan approval in 7 minutes by BankBazaar.com from leading Banks in	Sep 4
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	Jayde	Welcome To The Jayde Newsletter! - Welcome To WebProNews Welcome To The Jayde Newsletter! Before we begin, make sure to add	Sep 4
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	ineedhits noreply	[ineedhits] Your ineedhits Account and Password - ACCOUNT CREATION Account ID : A1588368 Dear Rah, Welcome to ineedhits. Yo	Sep 4
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	Rekha	Mobility Apps for Your Business - While we look at the span of last 20 years, we could broadly look at two distinct eras, - Life in	Sep 4
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	SlideShare Newsletter	Top Tips From the World Champions of PowerPoint - View online version Remember to display images Meet the PowerPoint World Champs Top Tips From the	Sep 4
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	Dilshad Pathan	Feeling Hesitate to Discuss personal Health Queries - My Life Care Follow Us on facebook twitter linkedin Google+ Feeling Hesitate to Discuss personal	Sep 4
<input type="checkbox"/>	<input type="star"/>	<input type="trash"/>	Vaishu	TAKE YOUR PICK. Register in SimplyMarry - TAKE YOUR PICK. Register in SimplyMarry -- Regards Vaishu	Sep 4

Spam

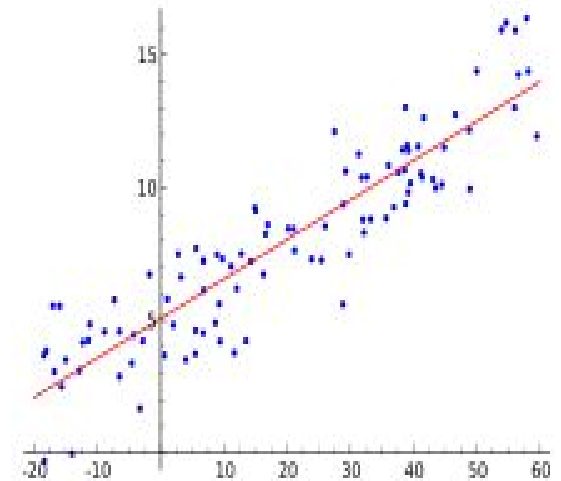
Spam

Classification Example

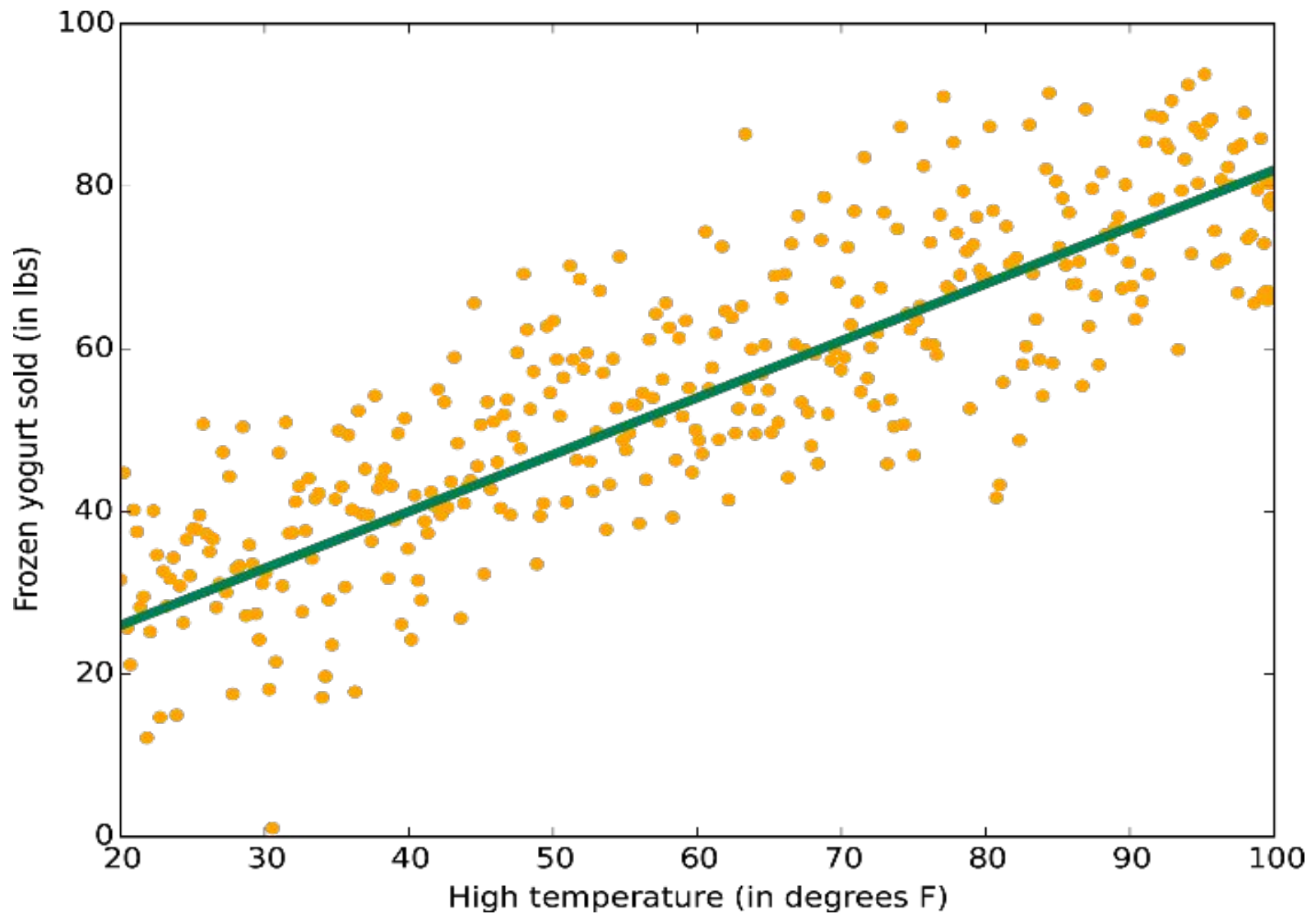


Regression

- **Best-Fit curve**
- **Regression analysis** is a statistical process for estimating the relationships among variables.
- Regression means to **predict** the output value using training data.
- Popular one is Logistic regression (binary regression)



Regression Example



Classification vs Regression

- Classification means to group the output into a class.
- discrete/categorical
- Regression means to predict the output value using training data.
- real num/continuous

Clustering

- **Clustering** is the task of grouping a set of objects in such a way that objects in the same group (called a **cluster**) are more similar to each other
- Objects are not predefined
- For e.g. these keywords
 - “man’s shoe”
 - “women’s shoe”
 - “women’s t-shirt”
 - “man’s t-shirt”
 - can be cluster into 2 categories “shoe” and “t-shirt” or “man” and “women”
- Popular ones are **K-means clustering** and **Hierarchical clustering**

Recommendations

More Items to Consider

You looked at

You might also consider

JavaScript: The Good Parts Paperback by Douglas Crockford
\$29.99 **\$19.79**

JavaScript: The Definitive Guide Paperback by David Flanagan
\$49.99 **\$31.49**

CSS: The Missing Manual Paperback by David McFarland
\$34.99 **\$23.09**

Learning jQuery 1.3 Paperback by John Resig
\$39.99 **\$35.99**

[Find similar items](#)

Related to Items You've Viewed

You looked at

You might also consider

Forms that Work Designing Web Forms for Usability by Steven Krug
\$24.99 **\$19.99**

Don't Make Me Think A Common Sense Approach to Web Usability by Steve Krug
\$19.99 **\$14.99**

Letting Go of the Words Writing Web Content that Sings by Linda Ward Beech
\$24.99 **\$19.99**

Designing Web Interfaces Principles and Best Practices for Creating Seamless, Usable, and Pleasurable Web Experiences by Steve Krug
\$24.99 **\$19.99**

Today's Recommendations For You

Here's a daily sample of items recommended for you. Click here to [see all recommendations](#)

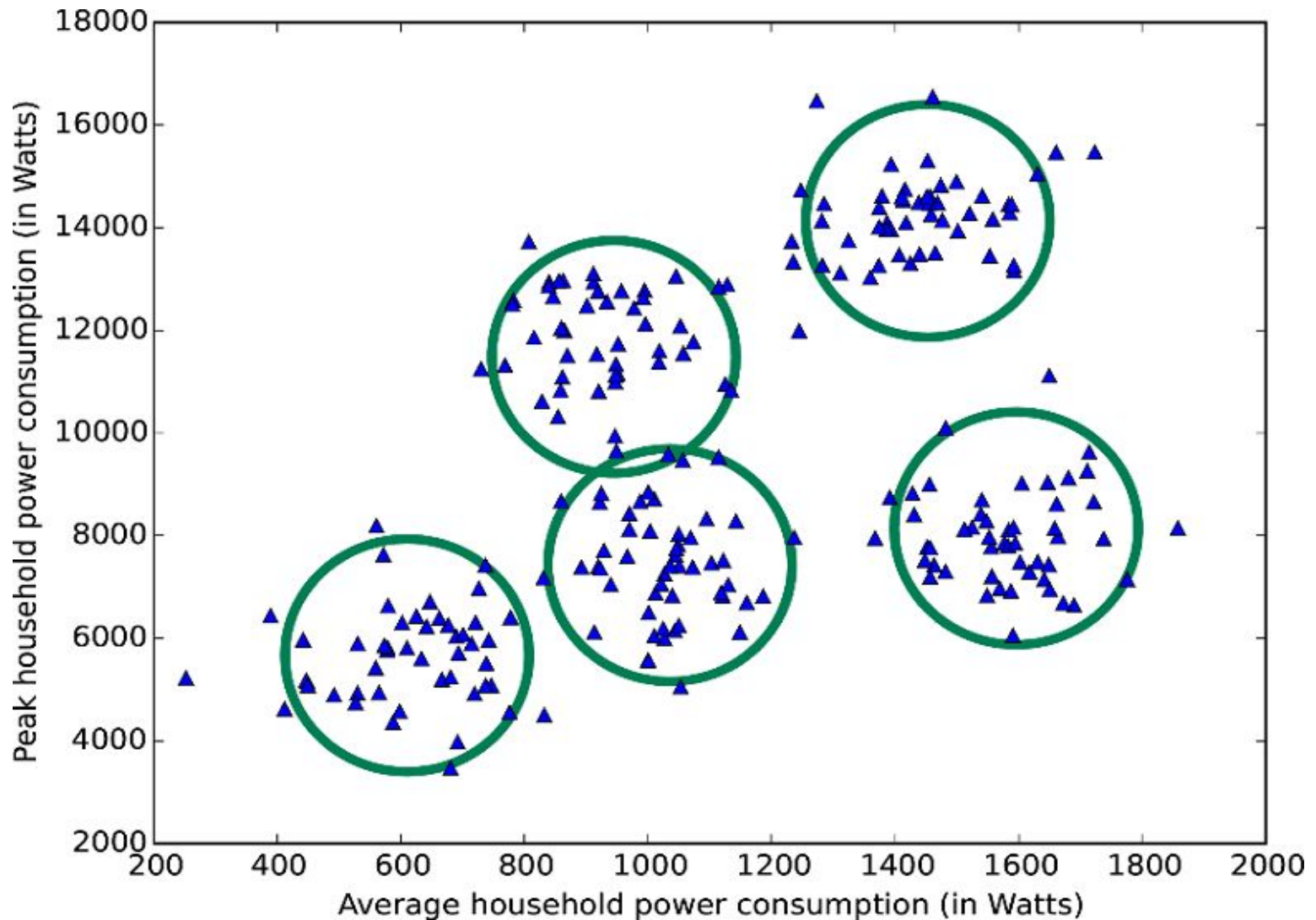
Even Faster Web Sites: Performance Tuning (Paperback) by Steve Souders
★★★★★ (7) \$23.10
[Fix this recommendation](#)

Simply JavaScript (Paperback) by Kevin Yank
★★★★☆ (19) \$26.37
[Fix this recommendation](#)

The Art & Science of JavaScript (Paperback) by John Resig
★★★★★ (3)
[Fix this recommendation](#)

[Any Category](#) Algorithms Boxed Sets Business & Culture Java
Graphic Design Microsoft Networking Networks, Protocols & APIs New SQL

Clustering Example

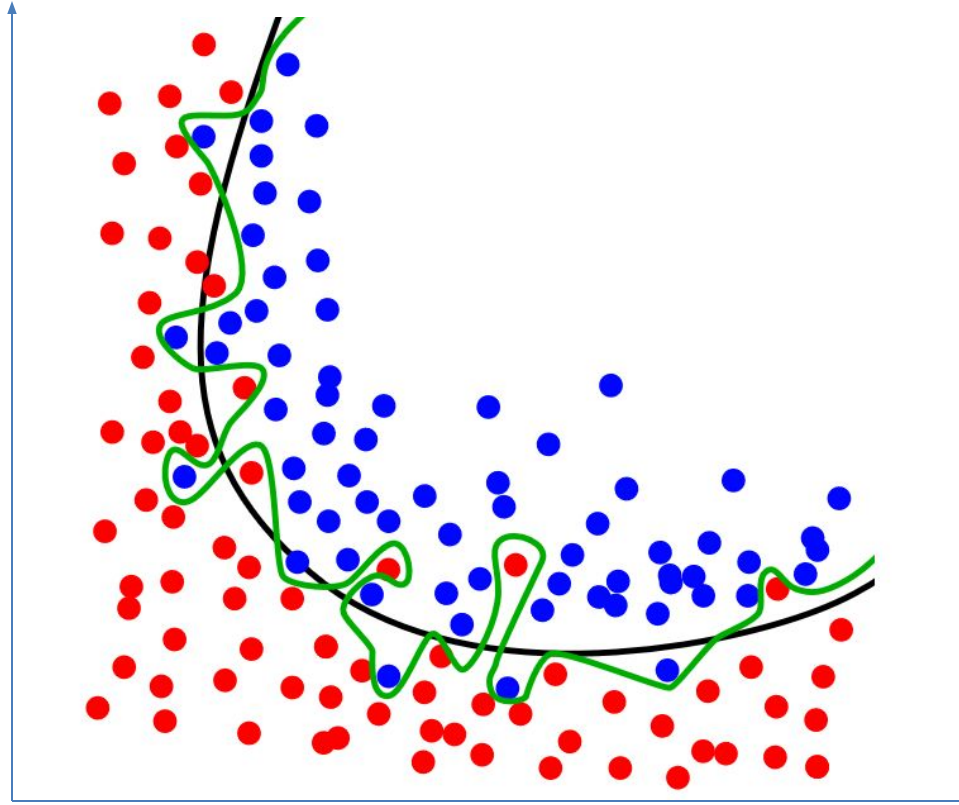


Pre-processing of data

Missing Values

- Reasons for missing values
 - Information is not collected
 - Attributes may not be applicable to all cases
- Handling missing values
 1. Eliminate Data Objects
 2. Estimate Missing Values
 3. Ignore the Missing Value During Analysis
 4. Replace with possible values

Over fitting



Evaluation Of Model

Evaluation of Model

- Partition: Training-and-testing
 - use two independent data sets, e.g., training set (2/3), test set(1/3)
 - used for data set with large number of samples
- Cross-validation
 - divide the data set into k subsamples
 - use k-1 subsamples as training data and one sub-sample as test data --- k-fold cross-validation
 - for data set with moderate size

Metrics for Performance Evaluation

- Confusion Matrix:

ACTUAL CLASS	PREDICTED CLASS	
	Class=Yes	Class=No
	Class=Yes	Class=No
	a (TP)	b (FN)
	c (FP)	d (TN)

$$\text{Accuracy} = \frac{a + d}{a + b + c + d} = \frac{TP + TN}{TP + TN + FP + FN}$$

Cost-Sensitive Measures

$$\text{Precision (p)} = \frac{a}{a + c}$$

$$\text{Recall (r)} = \frac{a}{a + b}$$

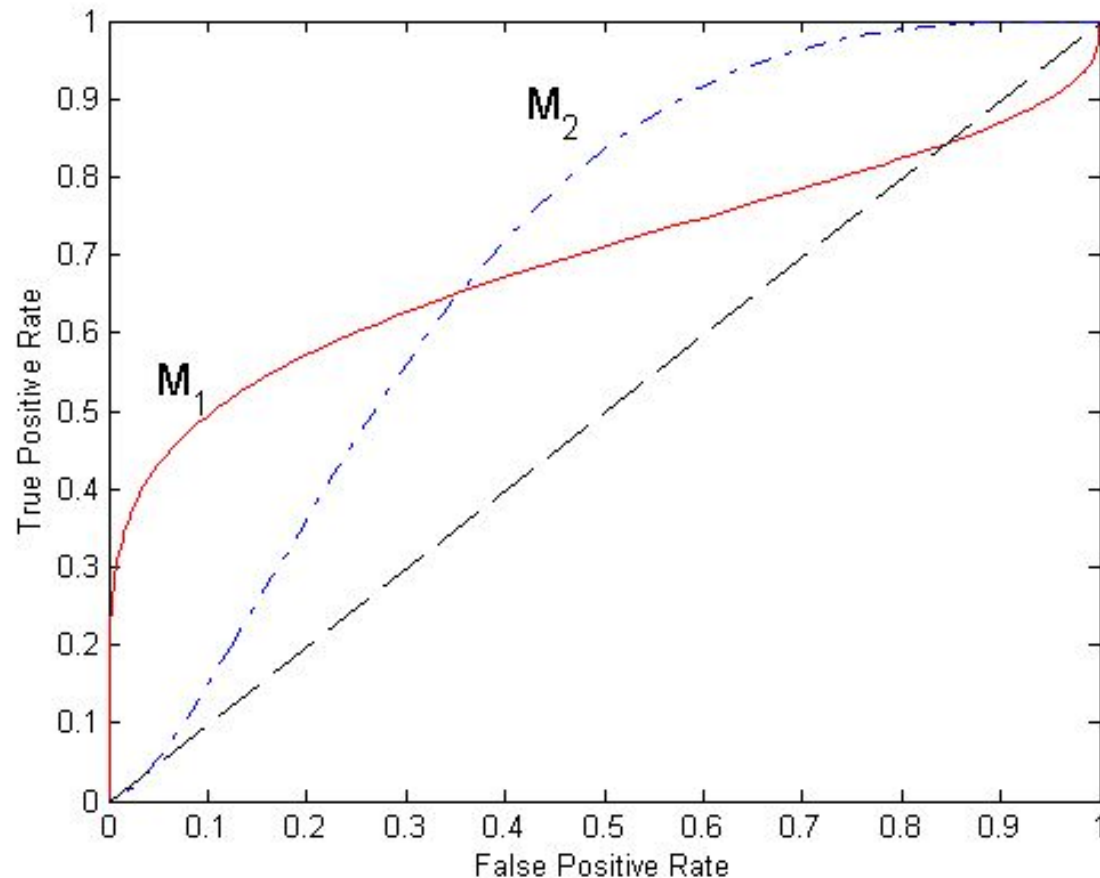
$$\text{F - measure (F)} = \frac{2rp}{r + p} = \frac{2a}{2a + b + c}$$

- Precision is biased towards $C(\text{Yes}|\text{Yes})$ & $C(\text{Yes}|\text{No})$
- Recall is biased towards $C(\text{Yes}|\text{Yes})$ & $C(\text{No}|\text{Yes})$
- F-measure is biased towards all except $C(\text{No}|\text{No})$

$$\text{Weighted Accuracy} = \frac{w_1 a + w_4 d}{w_1 a + w_2 b + w_3 c + w_4 d}$$

ROC Curve

- ROC curve plots TP (on the y-axis) against FP (on the x-axis)
- Area under ROC curve is used to evaluate Any Model
- M_1 is better for small FPR
- M_2 is better for large FPR



Important Libraries

- Numpy
- Pandas : load data
- Sklearn
- Xgboost

Deep Learning Libraries

- Tensorflow
- Theano
- Keras
- Torch
- Nolearn
- Lassagne