

Machine Learning !!

Definition !!

"Study of Computer Algorithms that allows the Computer Program to automatically improve through experience" [Efficient - learning]
[Training]

By Tom Mitchell (founder of ML Department)
School of Computer Science at Carnegie Mellon University.

* "ML is teaching the machine about something"

How !!

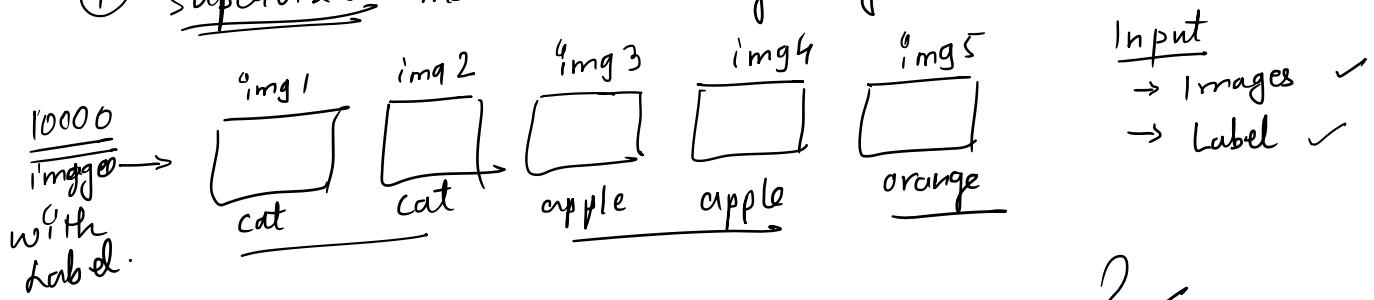
- ① Collect and clean the data
- ② Algorithm (model)
 - ↳ Selected (Readymade)
 - ↳ Built.
- ③ Teach the model essential pattern from data
(Training).
- ④ Export the model to give helpful answer.

Ex To Design a System that determine from MRI Scan, whether Tumor is present or not.

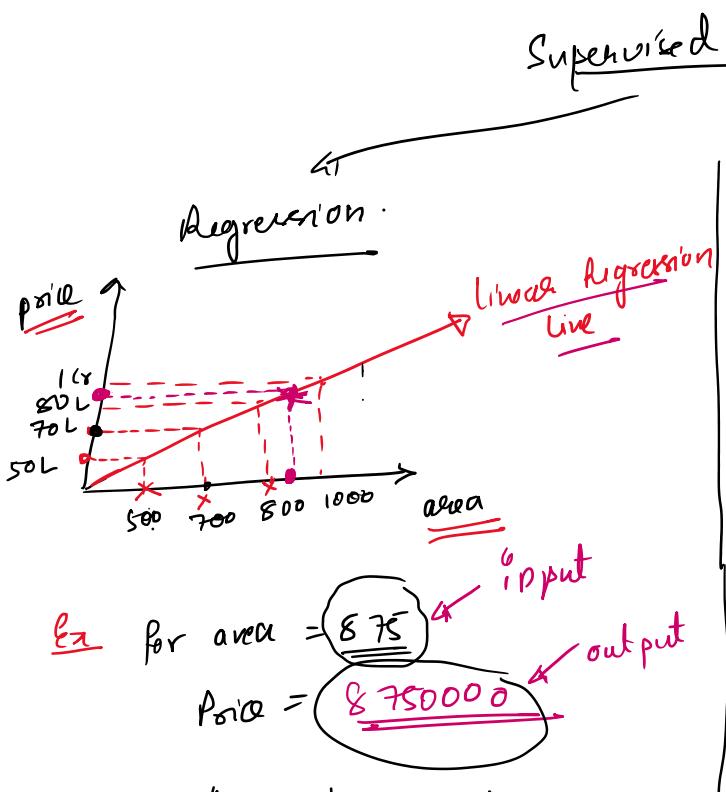
- ① Collect large No of MRI reports
Ex 10000 MRI report
Test 30%. Train 70%.
6000 has Tumor
4000 do not have Tumor.
- ② Build an efficient algorithm that detects presence or absence of Tumor in an MRI Scan.
[Expert Consultation - Radiologist].
- ③ To the Algorithm feed the 1000 MRI scans and allow the model to learn (train).
- ✓ ④ Use around 3000 Images (MRI scans) for testing
- ✓ ⑤ Use this model to determine presence or absence of Tumor from a New Image (MRI Scan)

Types of machine learning Algorithms.

① Supervised Machine learning Algorithms.



- * With lot of images as input, model will be able to identify pattern and will be able to predict.



Here op is not continuous value
 but could be Boolean or
 some class / category as output.

- Classification →
- * Spam Detection
 - Spam Mail
 - Spam Not Mail
 - * MRI Scan
 For Tumor Detection
 - Tumor Mail
 - Tumor Not Mail
 - * - always finds the op belongs to

Here area → Independent variable
 price → Dependent variable

area \Rightarrow Independent variable

price \Rightarrow Dependent variable

Here the dependent and independent Variable can have continuous value.

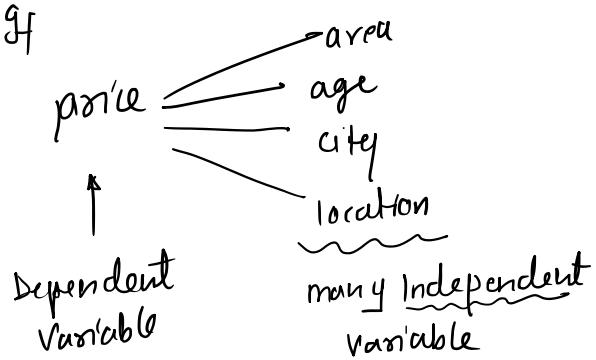
Types \Rightarrow

If price \rightarrow area.

Single Independent variable

\rightarrow Simple Regression

If

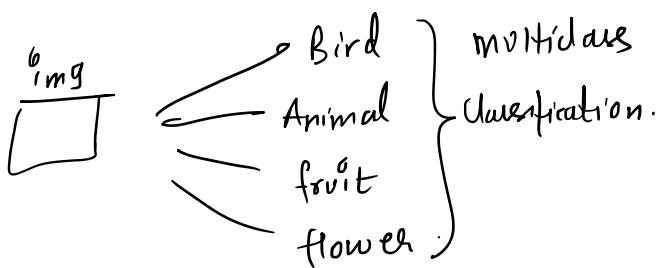


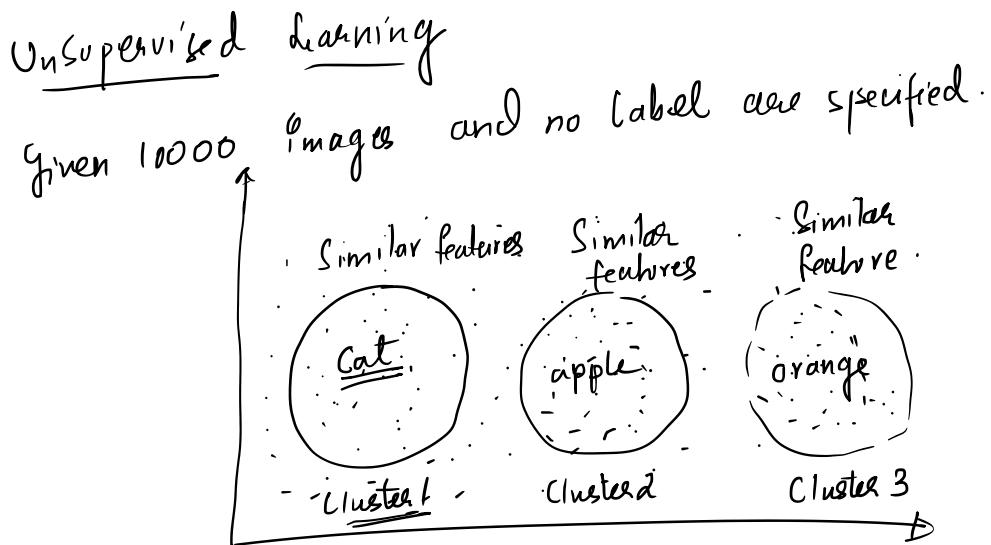
When we have more than one Independent Variable \Rightarrow multiple regression

In above cases the op belongs to one of the two classes.

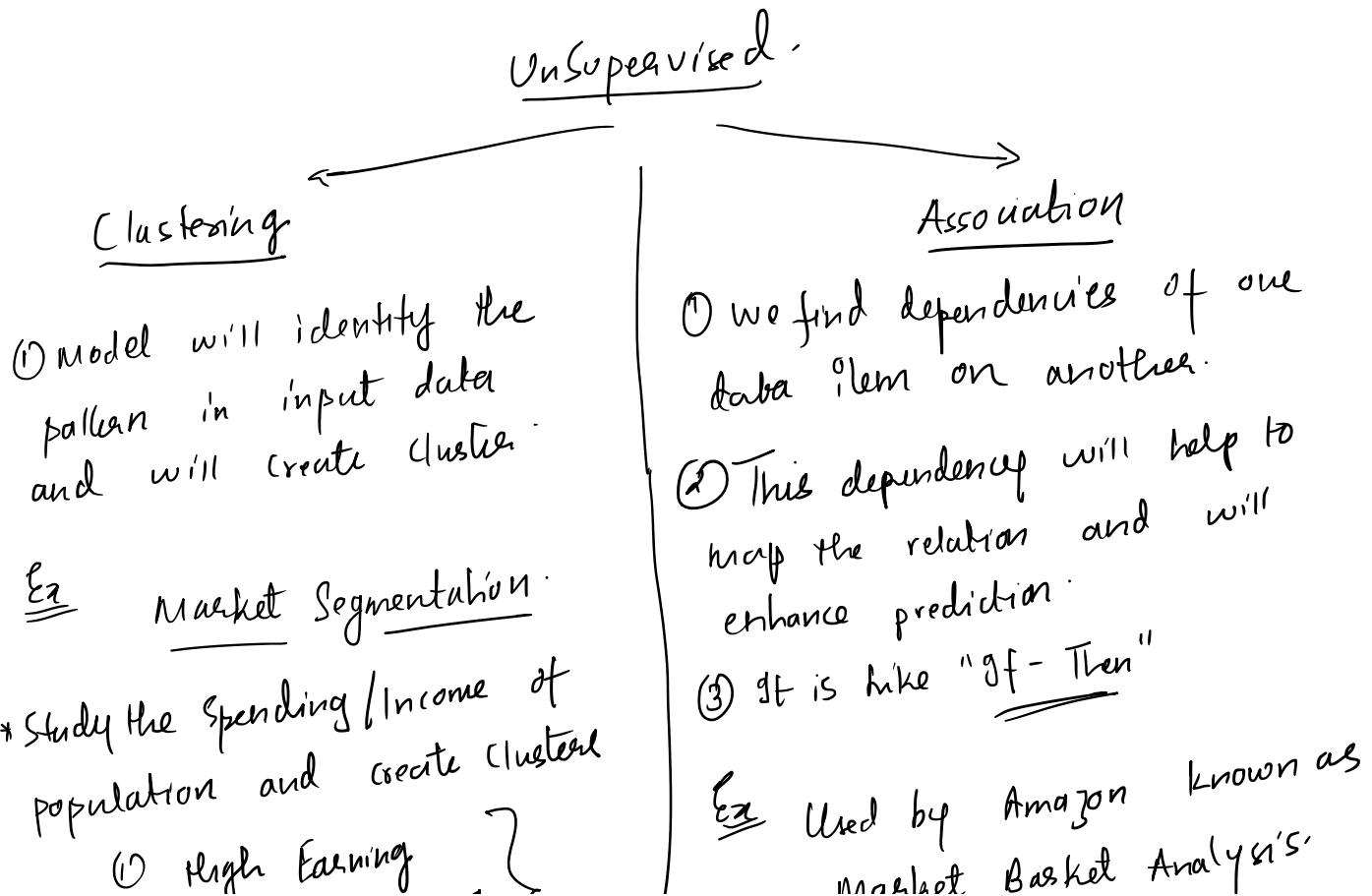
\rightarrow Binary classification.

If op can be one of many classes (more than two) Then it is called as multiclass classification.





- * Here the input is only data and no labels are associated with data.
- * Not sure about type of output
- * UnSupervised Algo will work on 10000 images and will create clusters of images based on the similarities.
- * It is "SELF LEARNING"



populations

- (1) High Earnings
- (2) Medium Earnings
- (3) Low Earnings

e.g. Used by Amazon

Market Basket Analysis

"If a person purchases cellphone than the person has tendency to purchase screen guard & backcover"

If → then

① Apriori Algorithm → is a breadth first

Search based approach which calculates the support for items.

e.g. bread influences the buyer to buy milk and eggs so the mapping helps increases the profit.

② FP Growth Algorithm →

→ Frequency pattern Algo finds the count of the pattern that has been repeated.

Application of UnSupervised learning

- Helps us in understanding pattern which can be used to cluster the data based on various features
- Understanding various defects in dataset which we are not been able to detect initially.
- It helps in mapping the various items based on the dependence of each other.

Real World Appl'n →

- ① Amazon → uses unsupervised learning to learn the customer purchase pattern and recommend the products that are most frequently bought together (Association Rule Mining)

② Credit Card Detection → Can we use unsupervised learning Algo to learn about various patterns of user and their usage of Credit Card.

→ If the card is used in parts that do not match the behaviour, an alarm is generated which could possibly be marked found



Disadvantage of UnSupervised Learning

- ① There is no way of obtaining the method the data is sorted as dataset is unlabelled
- ② They may be less accurate as input data is not known and labelled by humans.

Semi Supervised

Text Document Classification

Eg

1000000 articles and need to classify them
into
News
Literature
Research Paper
Medical Report

Label is not provided.

Manually labelling the 1000000 articles not possible.

→ We will label 10000 articles [Supervised Learning]
→ My model will be trained on these 10000 articles and
will use the pattern identified to classify the remaining
990000 articles [Unsupervised]

- * Uses small amount of labelled data
and large amount of unlabelled data.
- * Benefits of both labelled and unlabelled data.
- * Overcome the challenge of finding large amount of
labelled data.

Ex ① Speech Analysis → we will label audio files
which will need lot of human resource and
then we will use some technique that will help
to improve traditional speech analysis tool.

② Web Content Classification → Organizing the
knowledge available on billions of web pages
will need advance processing, but the
task needs human intervention to classify
the content -

Reinforcement Learning [Experiential Learning]

- Here the agent/model learns how to behave in an Environment by performing action and experiencing the result.

Type

Episodic Learning

- Here we have start and end state, thus an episode is created.
 - Thus the further action is based on feedback of result of earlier action.
- Ex Fear of dog after being bitten is Episodic learning.

Continuous Task Learning

- * There is no terminal state here
- * Agent/model that does Automated Stock Trading goes for Continuous Learning

* Challenges for Machine Learning →

- ① Data Collection → Data plays a key role in any machine learning.
 - * Govt. of work of a data scientist does is collecting relevant data
 - * Repositories like kaggle, ML Repository etc. are good only for initial learning
 - * For real case scenarios data needs to be collected through web-scraping or through clients etc.
 - * Data collection of relevant data for specific problem is challenge.

② Lots Amount of Training Data →

- Two important things that we do while doing a ml project
- ① Select a learning Algo
 - ② train the model with acquired data.

Finding Sufficient amount of relevant data for specific problem is challenge

③ Prov Quality of Data →

We cannot use the collected data directly for training. As the data might not be ready for training & we need to clean the data.

→ Data Preprocessing needs to be done by filtering missing values, extract and rearrange what model needs

④ Irrelevant / Unwanted Features →

In To predict no of hours a person needs to exercise →

[^{collected} ~~need~~ → age, gender, weight, height and location] ↑

→ Remove / filter unwanted features.

⑤ Overshifting the Training Data →

→ Overgeneralizing is something that we do very frequently. This is known as Overfitting

Avoid →

① Gather more training data.

② Select a model with fewer features.

③ Fix data errors ..

⑥ Underfitting the Training Data →

→ It happens when model is too simple to understand the basic structure of data.

→ When we have less information to construct model.

Avoid

① Need to feed better features to learning Algo

② Remove noise from data.

⑦ Offline Learning & Deployment of Model

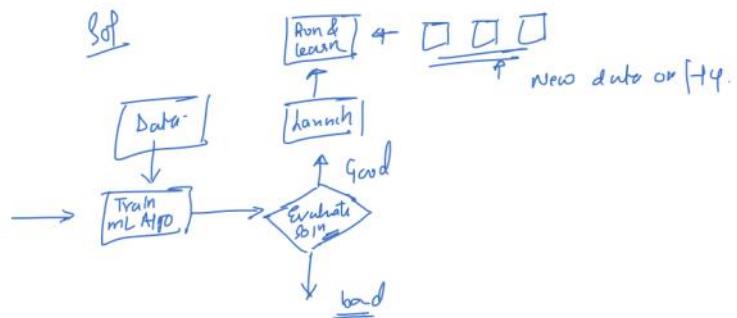
→ Challenges in Deployment

① Lack of practice & dependency issues

② Low understanding of underlying model with business.

③ Understanding of business problems.

④ Data Available online might be irrelevant
So we need to create source.



7 stages of ML are →

① Problem Definition → understand the problem that someone will solve

Question?

- ① What is the business
- ② Why does problem need to be solved.
- ③ What is measurable goal?
- ④ If probabilistic in nature, then does available data allow to model it

② Data Collection → Identify the source

Questions →

- ① What kind of data do I need
- ② Where is the data available
- ③ How can I obtain it
- ④ How to store and access it efficiently.

③ Data Preparation → can take from 70 - 90% of the work time.

steps involved like →

- Data Filling
- Data Validation & cleansing
- Data Formatting
- Data Aggregation & Reconciliation.

④ Data Visualization → perform Exploratory Data Analysis (EDA)

→ helps in identifying pattern.

ways of Visualization →

- Area chart
- Bar chart
- Bubble cloud
- Heat Map, Histogram
- Network Diagram, Wordcloud etc

w/w 1. Diagram, Wordcloud etc

⑤ ML Model → More magic Happens

→ finding of pattern in data using supervised or unsupervised.

→ Clustering, Regression, Classification etc.

→ Apply mathematical model to train ML algorithms that will help to make predictions.

⑥ Feature Engineering → it is collection of methods

for identifying an optimal set of inputs to ML Algorithm.

Characteristics of Good features

- ① Data represented must be unambiguous.
- ② Ability to capture linear & non linear relationship among data points
- ③ Capture context detail.

⑦ Model Deployment → to put the model in production environment to make data driven decisions in automated way.

ML models can be deployed as

① SaaS (Software as Service) ⇒

- cloud Application servⁱc
 - Use Internet to deliver applⁿ which are managed by third party vendor to its user
 - Executed on browser and don't requires download or installation on client side
- Ex → Google Apps, Dropbox, ParkMyCloud
Salesforce

② PaaS (Platform as a Service)

- provides cloud Components to certain c/w used for applⁿ.
- It delivers a framework for developers that can build upon and use to create Customized Applications

Ex → AWS Elastic, RedHat OpenShift, Windows Azure.
Vmware, IBM Bluemix

③ IaaS (Infrastructure as a Service)

→ Cloud Infrastructure services known as IaaS are made of highly scalable and Automated Computing Resources.

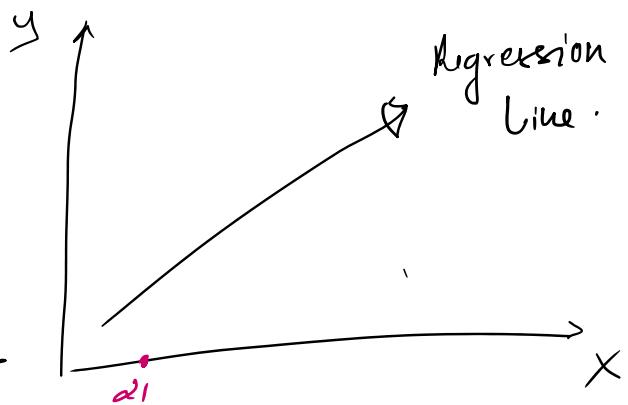
- It is fully self service for accessing & Monitoring things like Compute, Storage, N/w and other infrastructure related service on purchase on demand basis.

Ex → AWS, Azure, GCP(Google Cloud Platform),
IBM Cloud.

SR.NO	INDEPENDENT X	DEPENDENT Y	Actual value - $n=6$	
			$\sum x^2$	$\sum xy$
1	43	99	1849	4257
2	21	65	441	1365
3	25	79	625	1975
4	42	75	1764	3150
5	57	87	3249	4959
6	59	81	3481	4379

$$\sum x = 247 \quad \sum y = 466 \quad \sum xy = 11409 \quad \sum x^2 = 20485$$

NOW predict is $X = 55$, what is $\underline{\underline{y}}$?



$$y = a + bx + e$$

↑ ↑
 Dependent Intercept
 Variable } decides impact of x on y .

a = Intercept

b = Slope (coefficient of Independent Variable)

e = Error

[if $x \uparrow y \uparrow$ +ve Impact]

[if $x \uparrow y \downarrow$ -ve Impact].

Let us assume $e = 0$

$$y = a + bx \quad \text{--- (1)}$$

To find a and b .

Take \sum on both side of 1

$$\sum y = \sum a + \sum bx$$

$$\sum y = a \sum 1 + b \sum x$$

$$\sum y = an + b \sum x \quad \text{--- (2)}$$

Multiply Eq 2 with Independent Variable X .

$$\frac{1}{n} \sum \dots \sum x^2 \quad \text{--- (3)}$$

Multiply Eq 2 with Σx^2

$$\underline{\underline{\Sigma xy}} = a \underline{\Sigma x} + b \underline{\underline{\Sigma x^2}} \quad - (3)$$

Here $\frac{\text{Eq 2}}{\text{Eq 3}}$ $486 = a6 + b247 \quad \checkmark$
 $20485 = a247 + b11409 \quad \checkmark$

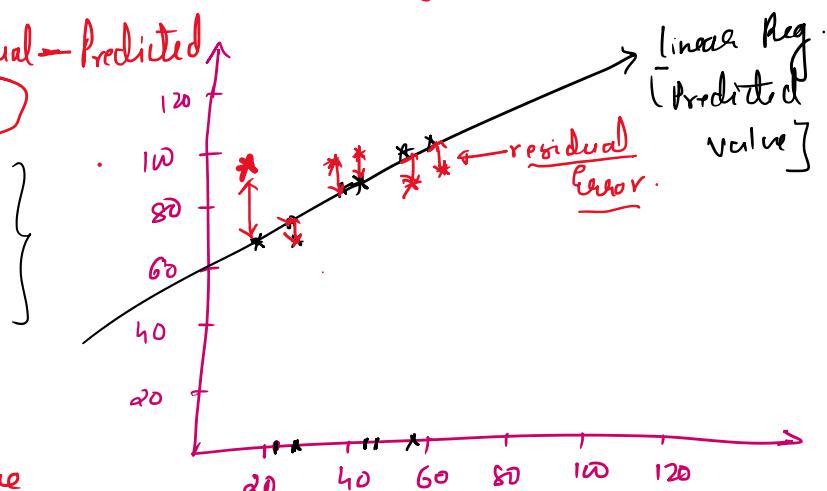
$$a = 65.14$$

$$b = 0.385$$

$$\boxed{y = 65.14 + 0.385x}$$

Using this let us calculate Predicted Y for every X

SR.NO	INDEPENDENT	DEPENDENT(ACTUAL) Y	PREDICTED Y	Absolute Difference
	X			
1	43	99	81.695	-17.305
2	21	65	73.225	8.225
3	25	79	74.765	-4.235
4	42	75	81.31	6.31
5	57	87	87.085	0.085
6	59	81	87.855	6.855



Residual Error = Error between the Actual Value and Predicted Value.

Jaad Rakho

$$\underline{\underline{y = a + bx}}$$

$$a = \frac{(\Sigma y)(\Sigma x^2) - (\Sigma x)(\Sigma xy)}{n(\Sigma x^2) - (\Sigma x)^2}$$

$$b = \frac{n(\Sigma xy) - (\Sigma x)(\Sigma y)}{n(\Sigma x^2) - (\Sigma x)^2}$$

Quantitative Analysis

$$b = \frac{n(\bar{xy}) - (\bar{x})(\bar{y})}{n(\bar{x}^2) - (\bar{x})^2}$$

Given $X \& Y$ Re^n $y = a + b x$
 Find a, b.

If more than one Independent Variable \Rightarrow

$$y = a + b_1 x_1 + b_2 x_2 + b_3 x_3 \dots$$

\uparrow impact
 of x_1 on y \uparrow impact
 of x_2 on y \uparrow impact
 of x_3 on y .