

What is generative AI?

→ Generative AI is the class of Artificial Intelligence where these AI do not use to predict According to the patterns from the data that they have learnt what we are actually doing, we were feeding our models with lots of data, and they are using that model for our predictions, but generative AI is not such a case. generative AI is the class of AI where we feed the data but that model actually used to create new data, earlier using machine learning learning, deep learning or NLP we were not be able to create new content or new data or new images. we were only able to do classifications, prediction, sentiment analysis.

Generative AI such a case

class of AI where it is use to create new data, new images, new text, new videos.

How it produces those kind of things?

→ It produces those kind of things, just by learning the patterns of the data.

Generative AI: → It contains two words Generative & AI

Generative → creating/generating new content

AI → using Artificial Intelligence

Generative → create new content (Audio, code, images, text video)

Artificial Intelligence → automatically using a computer program.

→ Generative AI encompasses various approaches, including generative Adversarial networks (GAN), Auto regressive models, variational autoencoders (VAEs).

When we talked about generative AI. that does not mean the LLM one or large language model or chatgpt or the mixture.

generative AI means Any model which can generate new content, new thing or new data. and in this process we have 3-4 kinds of things inside a generative AI.

Auto regressive model:- This is a decoder model which is used by gpt, most of the models which generates next text.

Generative Adversarial Networks:-

Generative means which is a generator which generates new content

Adversarial means there is some discriminator, Discriminator means there will be some neural network which will check generated content of generator. There will be two neural network, one will be for generator and other will be for discriminator.

Generator ek aisa neural network hai, jisko agar kuchh input data diya jaayega aur input data dene ke baad wo kuchh naya cheez generate karega, jab wo naya cheez generate kar lega tab yaha par discriminator aayega jisko (adversarial) bhi bolte hai, wo ye check karega ki jo shi naya data generator ne release kiyा hai wo real data se align kar raha hai ya kitna ka diff. aa raha hai. according to that it will give some scores and in this process it tries to learn. generating \Rightarrow Discreminating

Jaise hi generator (Neural Network) naya data generate karta hai according to input data, aur disreminator ne usko pahchaan liya, jo data generate hua hai, wo real ke jaisa hi hai, tab to generate generator fail ho gaya kyuki disreminator ne Kaise pahchaan liya.

Actually generator ko aisa cheez generate karra hai jo ho to real hi, lekin disreminator ko lage ki real nahi hai, foke hai. egg tarike se generator ek naya aur naye-2 cheez aur real type ka cheez generate kar payega

Variational autoencoders (VAEs) :- ye variational hai aur autoencoders hai.

autoencoders :- Eske paas kuchh encoded input honge jo ki embeddings me honge ya vectors me honge. Jisne kuchh information chhipa hoga aur ye ek neural network use karunga aur kuchh decoder architecture use karega kis, uss embeddings me, uss vectors me jo bhi information chhipa hua hai usko wo decode karke ek data output kar sake.

Autoencoders data ko as it is generate karta hai lekin, variational encoders uske jaisa sara output ko learn kar Jayega $\rightarrow 2, 2, 2, \sqrt{2} \rightarrow$ (jisme bhi uske variations ho sakte hain, wo saare cheez ko learn karayega) -
usecase :- Drug discovery

\Rightarrow Generative AI is not a new concept, we have google translation. It was, day before ULM come into existence. This was also an example of gen AI.
Usse firne jo method use bua usko ULM nahi kahde uska naam (Statistical machine Translation (SMT)) hai;
SMT :- Statistical method use karke the aur uske a/c jo bhi language hai usko translate karke the.

Jo statistical translation model hain esme 3 types ke
model use hote the

- (i) Alignment model
- (ii) Translation model
- (iii) Language model.

i) Alignment model :- Har ek word ko wo corresponding English word se usko relate karta tha. aligned karta tha.

eg:- le → the, chat → cat, noir → black ele

ye alignment huge data par trained tha jisme ek side koi our language the other side koi our language.

ii) Translation model :- esme ek sequence of character lete hain aur wo english me convert kar deتا haj.

eg-

le chat noir → the black cat.

iii) Language model:- ye aise output generate karta hai jo ki natural sounding ki farah ho aur ye sentence our semantic ko capture kar paaye.

Earlier we are also having, Siri, Cortana, Alexa, Bixby, all are the example of generative AI → they were using Hidden markov's models (HMMs) for speech recognition.

Hidden Markov models:- Hidden markov models are used for task like speech recognition where the goal is to understand spoken words.

eske alawa esme supervised learning algorithm hi use ho

Supervised learning :> ESKA matlab ye hai ki jaise hi user ne koi query puchha to wo model jo hai uss tarke se train tha ki wo classify kar pata ki ye query kis category ka hai.

NLP

feature Engg.

Working of GAN's

There is a real face

This is a generator neural network, this neural network will give some input some data inside the generator so generator will create face Θ new data according to the data which has been feed upon once the face has been generated now, it will be given to Discriminator neural network and discriminator neural network will compare Real face and generated face and it will give some score according to the things weather it is a real or fake.

Summary

GAN's me hota ye hai ki ek neural Network hai jisko generator bolte hain. esko humlog kuchh data de le hain uske baad ye naya data generate karta hai ab waha par Descriminator ka kaam hota hai ki uska jo generated data hai wo diya jata hai real data hai wo diya jata hai. ab descimator esme kya karta hai ki esko discriminate karne ka koshish karta hai generated data real hai ki fake hai agar wo generated data ko jaissa usne generate kiya hai agar wo real pancham liya to usko -1 assign karega, agar fake pancham to zero (0) value dega. es duration me training hota hai

Jyada koshish ye kiya jata hai ki generator
esf torah ka data generate karne ki discriminator
USKO pahchaan naa sake tab generator ko award
milega aur discriminator ko penalty.

Auto encoders:-

Autoencoders kya karta hai ki
eske pass 2 neural network hote
the ek hai encoder dusra decoder
encoder me humlog data dete hain data dena ke
baad, ek ye possibility hai ki data higher dimension
ka ho, eske baad data ka dimension reduce karne
ke baad ye ek latent state banata tha

Latent state:- jo bhi data input me numbers me
diya jata jayega, agar wo higher
dimension ka hai, to uske features ko
retain karte hue usko 100 dimension me lekar
aata hai, latent state eski ko bolte hain.

uske baad uss jo reduced dimension hain jiske andar
wo sare features hain, usne koi bhi features loss nahi
kiye hain koi bhi information loose nahi hua hai during
the dimensionality reduction usko usne phir decoder me
diya, phir decoder kya karta hai, jo cheez usme
encoded tha wo usko decode kar ke de deta hai.

eg:-

Input me hamne M.S. Dhoni ka image diya uske baad
hum esko no. me convert karnege uske baad
uska latent state banega, esme uske dimensions ko
reduced kar ke hum usme features ko provide karte
hoye usko latent state me store karne lage. phir usko
ek decoder me denge jo bhi khud bhi ek neural
network hai, wo decoder usko wohi cheez decode kar
ke, jaiga bhi hamara input tha as it is return kar
dega.

Variational AutoEncoders:

eg> jaise hi maine input me m.s.Dhoni ka image pass kiya esne usko no.'s me convert kiya. After conversion Dimension Reduction hua waha par uske features ko retain karte hue latent state banaya fir usse reduced dimension waale features ko decoder me diya: ab ye variational auto encoders kya karega ki ye sirf m.s.Dhoni ke features ko hi learn nahi karega, balki waisa jitter bhi images hain, usi variety ke, usn saare variations ko learn karega aur wo iss tarike ka ability khud me produce kar lega, ki sirf m.s.Dhoni jaisa image nahi, balki usse diff-2 kind of variation create karne ke kabul ho jayega. that is variational Autoencoder.

Architecture:

Encoder Network → Latent state → Decoder Network.

Autoregressive Models

ye two words ka combination hai Auto + Regressive
Regressive means - jo bhi past data hain usi past data ko regress kar dega/learn kar lega, automatically ek new tent generate karega. ye bhi ek statistical model ka ek type hai, eska use time series analysis me bhi ho raha hai, are mostly abhi ke time par sequence generation task me use kar raho hain. eski ke upar based hai large language model.

Large language model: Aisa model jo language ko iss tarike se represent kare/generate kare jo ki coherent bhi ho aur contextually relevant bhi ho usiliye esko language model bola jata hai.

Basically, LLM ek aisa AI hai jo ke Neural Network ya Deep learning technique ko use karta hai toaki wo human type text generate kar sake.

Actually LLM ko iss tarike se trained Kiya gaya hai ki usme usko next word predict karne ka capacity bhi diya gaya hai. Sabhi LLM aise trained hain ki wo next word predict kar sakte hai, looking over the past data. But aaj ke LLM's, QDA, summarizations machine translations bhi karta hai. lekin esko trained Kiya gaya tha sirf next word prediction ke liye. So ye downstreamed task kaise karta hai. for example machine translation, QDA?

Jab hamara model ye samagh leta hai ki how to produce next word, then usko fine-tuned Kiya jata hai; down-streamed task ke liye. fine-tuned ka matlab jaise hum apne LLM ko train Kar diya, uske baad wo next word predict Kar raha ho, QDA bhi kar raha hai summarize bhi kar raha hai iss se sab ke liye jo technique use hua usko kahte hai supervised-fine tuning

Supervised fine tuning → esme large corpus of data liya jata hai jaha par QDA means ek side question & ek side answer aise train Kiya jata hai, jab iss jorah ke data par train Kiya jata hai to eske andar QDA ka ability bhi aa jata hai. Similarly machine translation ke data par bhi train Kiya jata hai.

So downstreamed task ke liye jo bhi aapka base model tha (pretrained model) jo ki next word predict Kar raha tha usko phir aapne fine-tuned Kiya downstreamed task ke liye like QDA, MT, etc.

egs gpt ye sabhi kaam Kar raha hai

Today I tell date transformer is the base model of any LLM.

so understanding the basic architecture is very important to use the LLM in your use cases. Once you don't know the basic architecture of transformers you can use the LLM in your use cases but you will not have a great catch/knowledge may be you can say that may be you can't be able to fine-tune model according to your use cases.

Transformer Architecture me 2 cheeze hain. ek encoder aur ek decoder.

Encoder:- Jo bhi hum textual data input me dete hain usko ye vector me convert karta hai. aur vector me convert karne ke baad jitne bhi uske meanings & information hain usko collect kar ke ek final content vector return kar deta hai.

Sabse pahle hamne input me kuch pass kiya → phir uska embedding banaya, embedding means (No. me convert kiya) → eske baad Positional encoding kiya gaya means kon sa word kis word ke baad aaya hai (ye bahut matter karta hai kyuki humlog sequential data par dealing kar raha hain. agar hum word ka position change kar denge to complete meaning bhi change ho jayega). eski wasah se positional encoding kiya jata hai;

prob! → Ham Kaise pata Karenge ki agar two similar type ke word two diff scenario me use ho raha hai
eg →

The robber went to bank and rob the bank.
The robber went to the bank of river and looking at a tree

yah par hamne positional encoding se meaning define nahi ho payega

soln → Self Attention mechanism / Attention mechanism

Attention mechanism → Har ek word ka relation har ek word se kholta hai; define karta hai aur ek value / output generate karta hai;

eske baad Add & Norm aata hai; Add & Norm ka matlab ye hai ki jo bhi hame output aaya attention se usne kuchh cheez add kar ke normalize karte hain for better output. eske baad Feed Forward Neural Network aata hai (ye ek bahut barba neural network hai jo ki data ke pattern ko samjhta hai). eske baad phir Add & Norm (kuchh cheeze add kiya pichle output par phir usko ~~normalize~~ normalize kar diya). Finally encoder yaha par aapko ek content vector de dega. ab uss content vector ko decode karna hai aur next word predict karna hai decoder ka bhi same architecture hai but decoder encoder se uss tarah se diff. hai ki decoder me naya word predict hota hai. Encoder me naya word predict nahi hoga. encoder me aapne jo text diya hai uska embedding hoga uska pos. encoding hoga uska context samjha jayega uska pattern samjha jayega finally uska vector generate hoga. jo uske meaning ko rakha hoga. isko humlog content vector bolte hain agar usko naya data/text generate karna hai to uss content vector ko hame uss Decoder me dena hoga.

ek Nx layer me multi-head attention Add & Norm — feed forward (Add & Norm). aisa Nx layer open AI Ke GPT me 96 hain. har ek layer me feed forward network hai. aur uss har ek feed forward network me bahut sare hidden layers hain.

Encoder: content vector generate karega jo ki uske meaning ko rakha hoga. aur sabse multi-head attention hota hai multi-head attention me aapne gaya to multi-head attention kaam aata hai. feed forward network kya hai → neural Network huj large no. of hidden layers jo ki data ke pattern ko samjhta hai.

Decoder: New text generate karega.

Neural Network Ka first layer input layer kahata hai; jo last layer hai wo output layer kahata hai; bich bare layer ko hidden layer kahte hai.

e.g., ek house price prediction model banate hain-jisme kuchh features ke saath hamne price predict karne hai.

Ham saare features ko input layer ke perception me bhija ($x_1, x_2, x_3, \dots, x_n$) Sabka relation next hidden layer se hogा aur wo hidden layer ke harenk perception se karega Jis duration me ye sabhi se relate karega uss time me weight initialize kiya jayega. uss time me x_1 feature tha input layer ka uska kitna importance hai jab wo darse (hidden layer) perception me jaa raha hai; ye weight (w_1) batayega agar uska weight hoga jyada hai to matlab wo feature jyada imp. hai for price prediction. jab saare weightage assign ho jaate hain tab hamne ye pata chal jata hai ki kon sa column kitna imp. tha hamne price prediction ke liye. eske baad ek weighted sum next function ho jata hai: $(w_1x_1 + w_2x_2 + \dots)$ uske baad esme ek activation function hota hai; ye activation function em value/output generate karega jo ye predict karega ki house ka iss features aur iss weight ke aise kitna hona chahiye some ek bias bhi add kiya jata hai. bias already saare perception me already included raha hai.

Bias term provide neural networks with the flexibility to model complex relationships in data by shifts & adjustments in the activation functions and intercepts. Jo bhi hamne weight initialize kiya tha ye jannvi nahi hoi ki sahi ho ho sakte hoi ki no of bedrooms ko imp. feature maan ke jyada weightage de diya ho lekin actual me price ~~NO~~ ab bedroom ke badle size of bedroom par depend kar raha ho. esी cheez ko learn karta hoi during back propagation aur phir weight change karta hoi.

Activation function:- Jo bhi hamara data hai uske andar ke non-linearity ko capture karta hai.

Multihead Attention:- multi head me jo multi head hai esme ek head kya kaam karege subject(cat) uska content dekhega aur uske meaning ko dhundhega at a same time dusra head dusre word (sits) verb/action wo mat se kaise relate hai kya content hai eska. iss tarice se saare head apna apna words ke reaction ko identify karta hai to Sab apna ek score generate karta hai. sabhi head ka apna ek score hota hai. finally uska ek avg. calculate kar ke aage forward kija jata hai.

eg:- The cat sits on mat.

Masked multi head attention:-

First previous wale token ko consider karne deta hai jo token aap attention ke liye use kar raho ho uss token se aage wale token ko wo mask kar deta hai;

Output layer me jo probability hai usko bhi select karne ke two methods hain. ek hai greedy decoding and one is beam search

Greedy decoding:- Apko kuchh words milenge next word ke fav par uski probability value milogi aur apko select karna hai; jiski probability value jyada hai aapko usi word ko as a next word predict karna hai; usko greedy decoding bolte hain.

⑪ Beam Search: → esme updated sequence banata hai
phir uske baad ek probability ka multiply
karta hai aur ek combine probability ki value nikal ke
ek final decision leta hai ki kon se sentence to selectra
hai.

Tokenization & Positional Embedding

Tokenization

Tokenization ek process hai jisme aap apne
input sequence ko divide karte ho by the basis of
character, word, sentence. aap apne rules/seed ke
hijaab se karte ho

Training of GPT

① Unsupervised pre training: → Hamne simply data ko liya
aur seed kar diya. ab model
khud se bina human labelled
ke/bina human annotation ke model khud se data ke
pattern ko samjhega, usko structure karega relationship
ko samjhega, without any guidance, without any external
sources aur phir uske pattern ko samajh ke ek apna
go model hai wo ready karega. esi process ko
un-supervised pre training kahte hain.

② Self Supervised learning: esme ham koi sentence input
me diye, us sentence ke kisi bhi
word ko ab truncate kar denge
ab usko predict karne ka try karenge. ESS process
ko bolte hain self supervised learning. Khud ke supervision
me bina kisi external efforts ke. Khud ke supervision me
jab model train hota hai ya model learn karta hai
to esko kahte hain self supervised learning.

Abhi tak ye model bahut achha nahi hua hai. abhi tak hum next word prediction kar paa rakte hain. masking ke bad bhi prediction kar paa rakte hain. esme kuchh error bhi aage hamne error ko minimize kiya ab performance sahi hai, lekin eske aage kuchh downstreamed task the ego Q&A, Summarization, machine Translation iss sab ke liye model ko train karna hoga.

abhi tak haanara model/machine unsupervised se data ke pattern ko samjha hai supervised se data ke next word prediction ko samjha hai masking ko filling karna samjha hai. Q&A, MT - ers sab ke liye ab apne model ko fine tune karna hoga. ab esko fine-tuning supervised type se hoga.

ab hamare paas jo data hoga wo labelled hoga jisme ek side questions doinge aur ek side answers honge.

Same machine translation me bhi hamare pass data hoga jisme ek side koi lang hoga dusre side phir koi ek aur language. hoga

GPT Parameters

- nparams - 125m (meatlab ye hai ki esme kitne weights & bias hain.)
- n layers → Decoder Architecture me kitna layer hain. har Ek layer me → ek feed forward neural network hoga aur ek attention mechanism hoga.)
- dmodel → Embedding ka dimension kitna hain
- n head → ek attention mechanism me kitne head hain (multihead)
- d-head → This refers to the dimensionalities of the query keys & values used in each attention head. Typically d-heads is calculated as d-model divided by n-heads

batch size:- This refers to the no. of input examples (sequences) processed in parallel during each training iteration. A larger batch size can lead to faster training but may require more memory.

GPT Playground

Temperature! \rightarrow Temp. ye batata hai ki hamara jo model hai wo kitna random output generate kar sakte hai (output me randomness hata hai.)

e.g. hamne koi question input kiya aur temp=0 assign kiya hua hai to answer aayega kuchh output me agar again same question aur temp=0 alpit me same answer aayega. lekin agar temp ka value change kar dege to hamesha output me variation aata ratega.

max length! \rightarrow max= hum kitna token generate karna chahte hain.

model! \rightarrow Specifies which version of the gpt model you are using/want to use.

stop sequences! wo kitna token generate karne ke badd rok dega ye parameter hum esme paas karte hain.

inf top p \rightarrow ye ek nucleus sampling hota hai;

hamare paas har ek word ka prob. milta hai for next word. eske baad hum uska cumulative prob. nikalte hain. aur hamne apna stop p ka parameter set kar diya hain apni li cumulative prob. ka value stop p ke parameter ke cross karega usko wo nucleus word select kar leta hai;