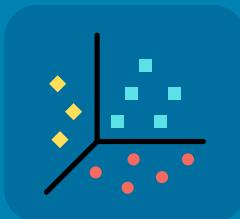
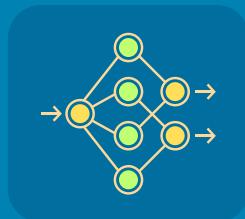


2024 EDITION

FREE

DATA SCIENCE

FULL ARCHIVE



Daily Dose of
Data Science

Avi Chawla
DailyDoseofDS.com

How to make the most out of this book and your time?

Reading time of this book is about 9-10 hours. But not all chapters will be of relevance to you. This 2-minute assessment will test your current expertise and recommend chapters that will be most useful to you.

Which Chapters from the Full Archive Are of Most Relevance to You?

Answer 21 yes/no questions and we'll email you the list of chapters that must read to improve your skillset.

[Take the Assessment Now!](#)



Start The Assessment

Name *

Email *

[Start the Assessment](#)

Scan the QR code below or open this link to start the assessment. It will only take 2 minutes to complete.



<https://bit.ly/DS-assessment>

Table of contents

Section #1) Deep learning.....	7
1.1) Learning Paradigms.....	8
Transfer Learning, Fine-tuning, Multitask Learning and Federated Learning.....	8
Introduction to Federated Learning.....	12
Building Multi-task Learning (MTL) Models.....	15
Active Learning.....	19
1.2) Run-time and Memory Optimization.....	27
Momentum.....	27
Mixed Precision Training.....	33
Gradient Checkpointing.....	41
Gradient Accumulation.....	46
4 Strategies for Multi-GPU Training.....	52
1.3) Miscellaneous.....	56
Label Smoothing.....	56
Focal loss.....	59
How Dropout Actually Works?.....	64
Issues with Dropout in CNNs.....	70
What Hidden Layers and Activation Functions Actually Do.....	75
Shuffle Data Before Training.....	79
1.4) Model compression.....	82
Knowledge Distillation for Model Compression.....	82
Activation Pruning.....	89
1.5) Deployment.....	91
Deploy ML Models from Jupyter Notebook.....	91
4 Ways to Test ML Models in Production.....	96
Version Controlling and Model Registry.....	99
1.6) LLMs.....	102
Where Did the GPU Memory Go?.....	102
Full-model Fine-tuning vs. LoRA vs. RAG.....	108
5 LLM Fine-tuning Techniques.....	113

Section #2) Classical ML..... 116

2.1) ML Fundamentals.....	117
Training and Inference Time Complexity of 10 ML Algorithms.....	117
25 Most Important Mathematical Definitions in Data Science.....	119
How to Reliably Improve Probabilistic Multiclass-classification Models.....	123
Your Entire Model Improvement Efforts Might Be Going in Vain.....	128
Loss Function of 16 ML Algos.....	131
10 Most Common Loss Function.....	132
How to Actually Use Train, Validation and Test Set.....	133
5 Cross Validation Techniques.....	137
What To Do After Cross Validation?.....	140
Double Descent vs. Bias-Variance Trade-off.....	145
2.2) Statistical Foundations.....	148
MLE vs. EM — What's the Difference?.....	148
Confidence Interval and Prediction Interval.....	153
Why is OLS Called an Unbiased Estimator?.....	159
Bhattacharyya Distance.....	165
Why Prefer Mahalanobis Distance Over Euclidean distance?.....	169
11 Ways to Determine Data Normality.....	173
Probability vs. Likelihood.....	178
11 Key Probability Distributions in Data Science.....	183
A Common Misinterpretation of Continuous Probability Distributions.....	189
2.3) Feature Definition, Engineering and Selection.....	195
11 Types of Variables in a Dataset.....	195
Cyclical feature encoding.....	204
Feature Discretization.....	208
7 Categorical Data Encoding Techniques.....	212
Shuffle Feature Importance.....	215
The Probe Method for Feature Selection.....	218
2.4) Regression.....	221
Why Mean Squared Error (MSE)?.....	221
Sklearn Linear Regression Has No Hyperparameters.....	225
Poisson Regression vs. Linear Regression.....	229
How to Build Linear Models?.....	232

Dummy Variable Trap.....	235
Visually Assess Linear Regression Performance.....	236
Statsmodel Regression Summary.....	238
Generalized Linear Models (GLMs).....	245
Zero-inflated Regression.....	249
Huber Regression.....	252
2.5) Decision Trees and Ensemble Methods.....	256
Condense Random Forest into a Decision Tree.....	256
Transform Decision Tree into Matrix Operations.....	262
Interactively Prune a Decision Tree.....	272
Why Decision Trees Must Be Thoroughly Inspected After Training.....	275
Decision Trees ALWAYS Overfit!.....	279
OOB Validation in Random Forest.....	282
Train Random Forest on Large Datasets.....	286
A Visual Guide to AdaBoost.....	289
2.6) Dimensionality Reduction.....	295
The Utility of 'Variance' in PCA.....	295
KernelPCA vs. PCA.....	298
PCA is not a Visualization Technique.....	302
t-SNE vs. SNE — What's the difference?.....	306
How To Avoid Getting Misled by t-SNE Projections?.....	312
Accelerate tSNE with GPU.....	315
Scale tSNE to Millions of Data Points With openTSNE.....	318
PCA vs. t-SNE.....	320
2.7) Clustering.....	323
Beyond KMeans: 6 Must-Know Types of Clustering Algorithms.....	323
Intrinsic Measures for Clustering Evaluation.....	324
Breathing KMeans vs KMeans.....	330
How Does MiniBatchKMeans Works?.....	335
ANN-driven KMeans with Faiss.....	340
KMeans vs. Gaussian Mixture Models.....	342
DBSCAN++: A Faster and Scalable Alternative to DBSCAN.....	345
HDBSCAN vs. DBSCAN.....	348

2.8) Correlation analysis.....	350
Correlation != Predictiveness.....	350
Beware of Summary Statistics.....	353
Pearson correlation can only measure linear association.....	355
Correlation with Ordinal Categorical Data.....	358
2.9) Drift.....	361
Multivariate Covariate Shift.....	361
Using Proxy-Labelling to Identify Drift.....	371
2.10) kNN.....	375
Using kNNs on Imbalanced Datasets.....	375
Approximate Nearest Neighbor Search Using Inverted File Index.....	381
2.11) Kernels.....	385
Why is Kernel Trick Called a "Trick"?.....	385
The Mathematics Behind RBF Kernel.....	388
2.12) Missing Data.....	391
3 Types of Missing Values.....	391
MissForest and kNN Imputation.....	396
2.13) Pitfalls and Misconceptions.....	401
When is Random Splitting Fatal for ML Models?.....	401
Feature Scaling is NOT Always Necessary.....	406
A Misconception About Log Transform.....	408
The True Purpose of Feature Scaling and Standardization.....	411
L2 Regularization is Not Just Used for Regularization.....	414
2.14) Miscellaneous.....	420
Is Your Model Data Deficient?.....	420
Bayesian Optimization.....	424
Train and Test-time Data Augmentation.....	426
2.15) Data Analysis.....	430
15 Pandas \leftrightarrow Polars \leftrightarrow SQL \leftrightarrow PySpark Translations.....	430
2 Alternatives to Pandas' Describe.....	433
Accelerate Pandas with GPU Using RAPIDS cuDF.....	436
Missing Data Analysis with Heatmaps.....	438
DataFrame Styling.....	442

8 Automated EDA Tools.....	443
2.16) Data Visualisation.....	444
Most Important Plots in Data Science.....	444
How are QQ Plots Created?.....	451
8 Elegant Alternatives to Traditional Plots.....	456
Interactive Controls.....	464
Mosaic Plots.....	465
Enrich Matplotlib Plots with Inset Axis and Annotations.....	466
Professionalize Matplotlib Plots.....	469
Sankey Diagrams.....	471
Ridgeline Plots.....	474
Sparkline Plots.....	477
2.17) SQL.....	479
Grouping Sets, Rollup and Cube in SQL.....	479
Semi, Anti, and Natural Joins.....	485
Use SQL "NOT IN" With Caution.....	491
2.18) Python OOP.....	496
Getters and Setters.....	496
Descriptors in Python.....	499
20 Most Common Magic Methods.....	509
Memory Efficient Class Objects using Slots.....	510
Why Don't We Invoke model.forward() in PyTorch?.....	515
True OOP Encapsulation is Missing From Python.....	520
A Common Misconception About __init__()......	524
Function Overloading in Python.....	527

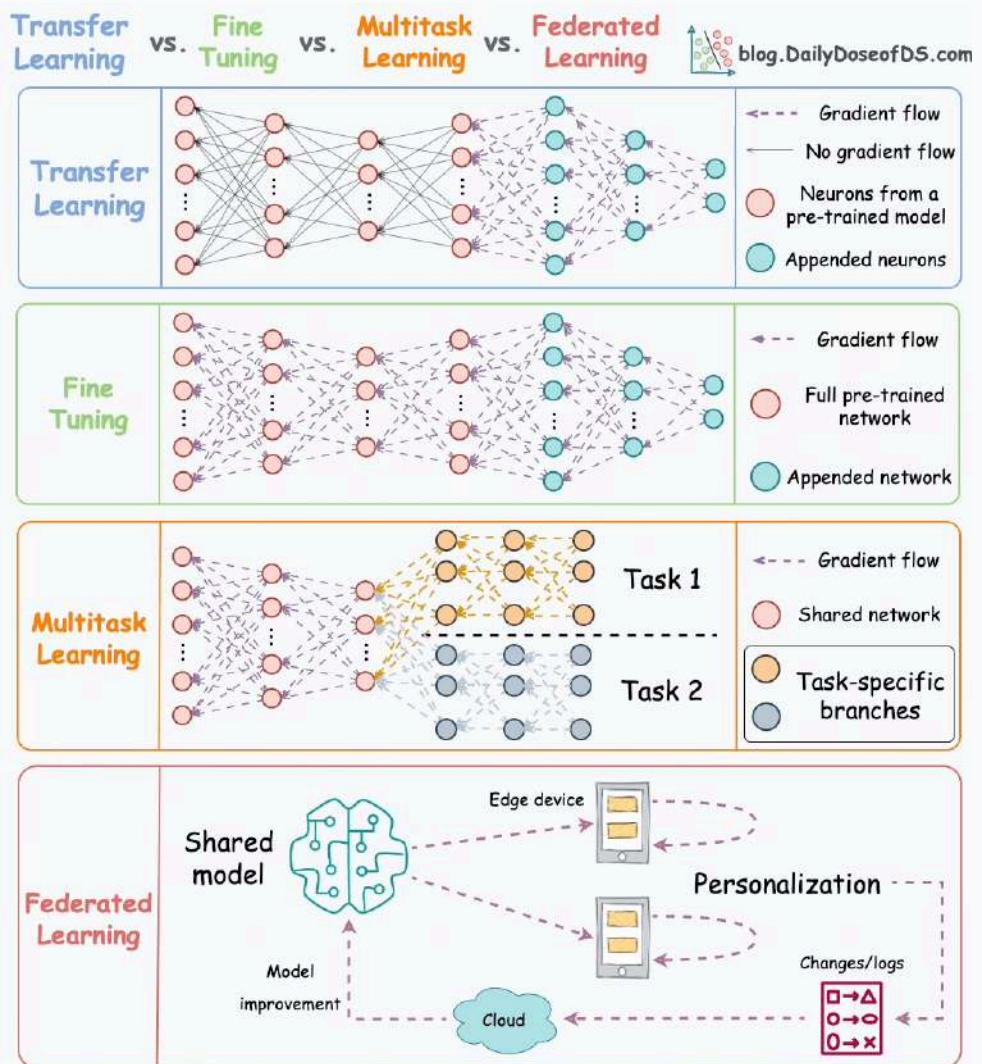
Deep learning

Learning Paradigms

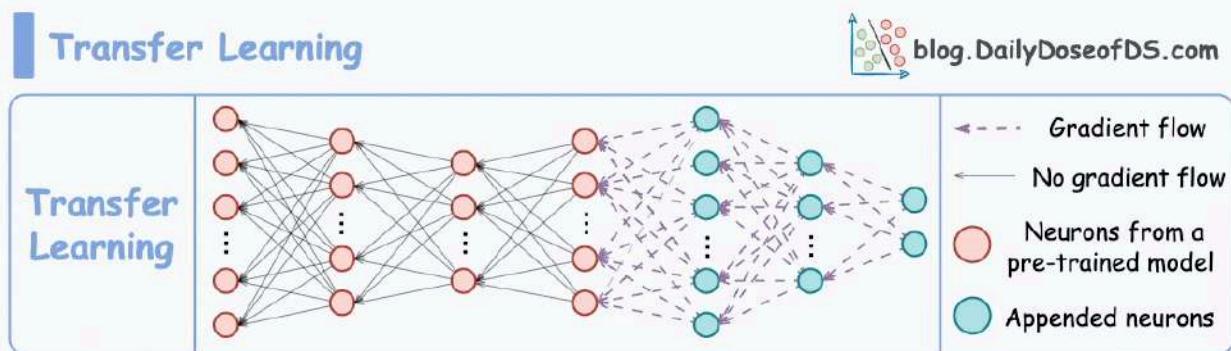
Transfer Learning, Fine-tuning, Multitask Learning and Federated Learning

Most ML models are trained independently without any interaction with other models. However, in the realm of real-world ML, there are many powerful learning techniques that rely on model interactions to improve performance.

The following image summarizes four such well-adopted and must-know training methodologies:



#1) Transfer learning



This is extremely useful when:

- The task of interest has less data.
- But a related task has abundant data.

This is how it works:

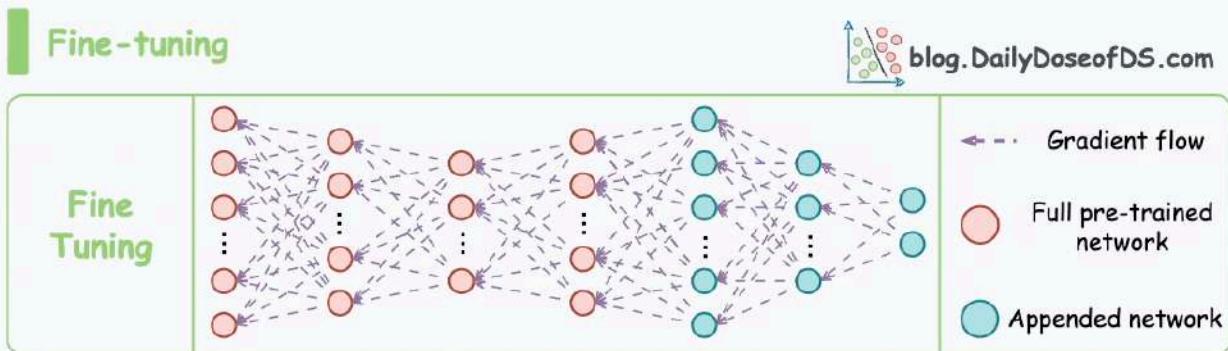
- Train a neural network model (base model) on the related task.
- Replace the last few layers on the base model with new layers.
- Train the network on the task of interest, but don't update the weights of the unreplaced layers during backpropagation.

By training a model on the related task first, we can capture the core patterns of the task of interest. Later, we can adjust the last few layers to capture task-specific behavior.

Another idea which is somewhat along these lines is knowledge distillation, which involves the “transfer” of knowledge. We will discuss it in the upcoming chapters.

Transfer learning is commonly used in many computer vision tasks.

#2) Fine-tuning

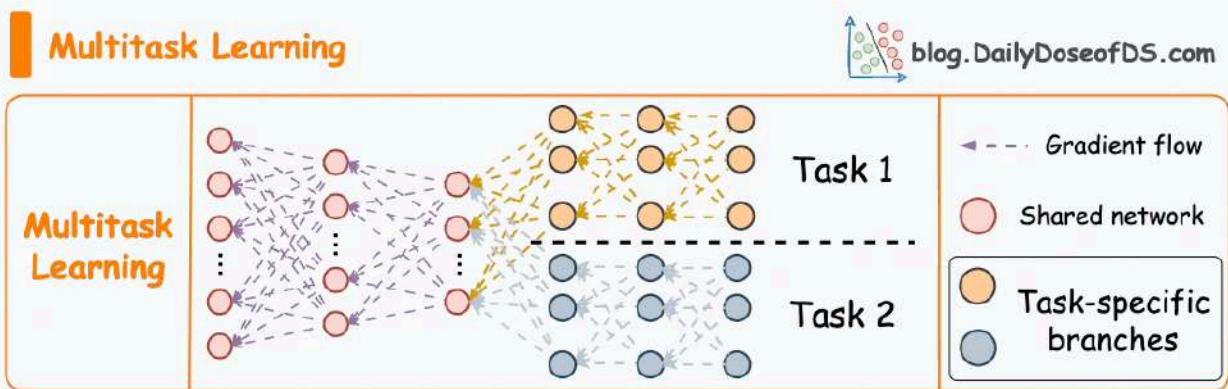


Fine-tuning involves updating the weights of some or all layers of the pre-trained model to adapt it to the new task.

The idea may appear similar to transfer learning, but in fine-tuning, we typically do not replace the last few layers of the pre-trained network.

Instead, the pretrained model itself is adjusted to the new data.

#3) Multi-task learning



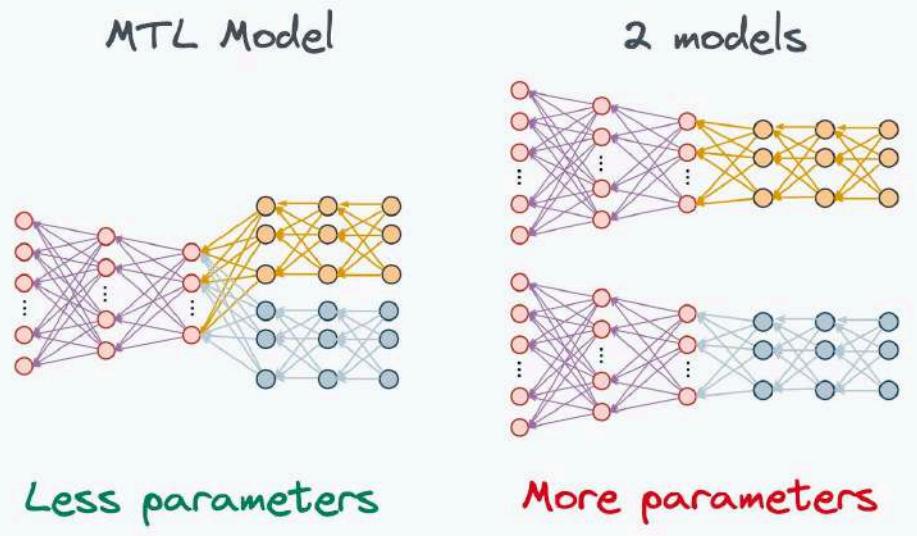
As the name suggests, a model is trained to perform multiple tasks simultaneously.

The model shares knowledge across tasks, aiming to improve generalization and performance on each task.

It can help in scenarios where tasks are related, or they can benefit from shared representations.

In fact, the motive for multi-task learning is not just to improve generalization.

We can also save compute power during training by having a shared layer and task-specific segments.



- Imagine training two models independently on related tasks.
- Now compare it to having a network with shared layers and then task-specific branches.

Option 2 will typically result in:

- Better generalization across all tasks.
- Less memory utilization to store model weights.
- Less resource utilization during training.

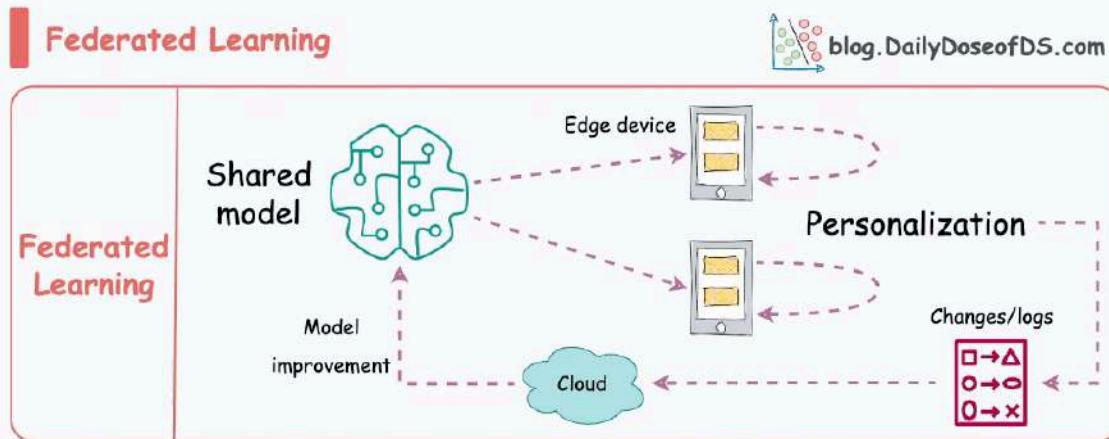
#4) Federated learning

Let's discuss it in the next chapter.

Introduction to Federated Learning

In my opinion, federated learning is among those very powerful ML techniques that is not given the true attention it deserves.

Here's a visual that depicts how it works:



Let's understand this topic in this chapter and why I consider this to be an immensely valuable skill to have.

The problem

Modern devices (like smartphones) have access to a wealth of data that can be suitable for ML models.

To get some perspective, consider the number of images you have on your phone right now, the number of keystrokes you press daily, etc.

That's plenty of data, isn't it?

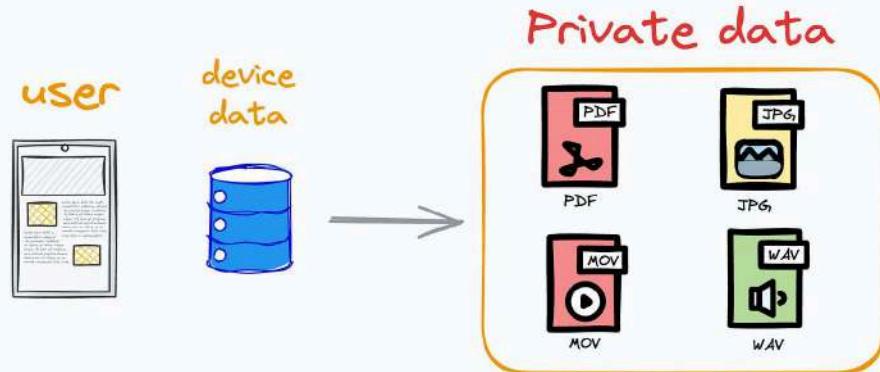
And this is just about one user — you.

But applications can have millions of users. The amount of data we can train ML models on is unfathomable.

So what is the problem here?

The problem is that almost all data available on modern devices is private.

- Images are private.
- Messages you send are private.
- Voice notes are private.



Being private, it is likely that it cannot be aggregated in a central location, as traditionally, ML models are always trained on centrally located datasets.

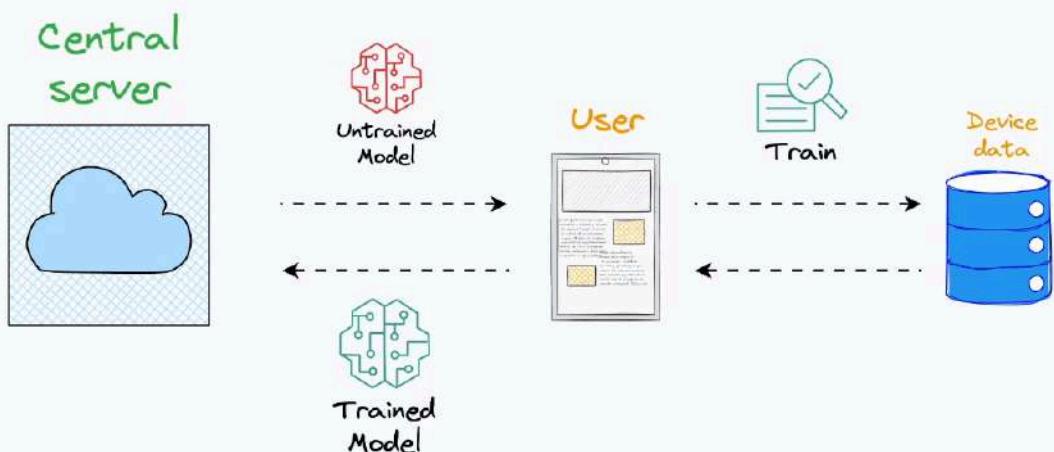
But this data is still valuable to us, isn't it?

We want to utilize it in some way.

The solution

Federated learning smartly addresses this challenge of training ML models on private data.

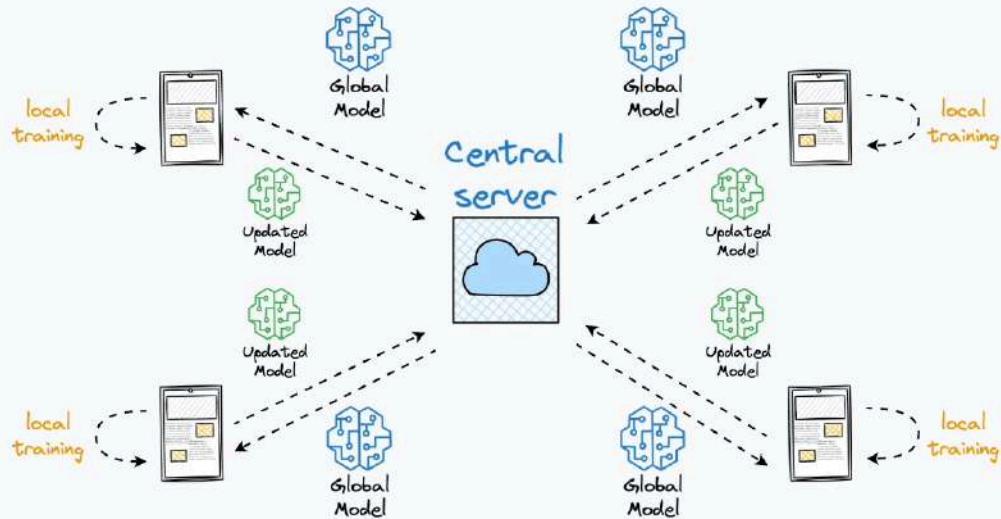
Here's the core idea:



- Instead of aggregating data on a central server, dispatch a model to an end device.
- Train the model on the user's private data on their device.
- Fetch the trained model back to the central server.
- Aggregate all models obtained from all end devices to form a complete model.

That's an innovative solution because each client possesses a local training dataset that remains exclusively on their device and is never uploaded to the server.

Yet, we still get to train a model on this private data.



Send a global model to the user's device, train a model on private data, and retrieve it back.

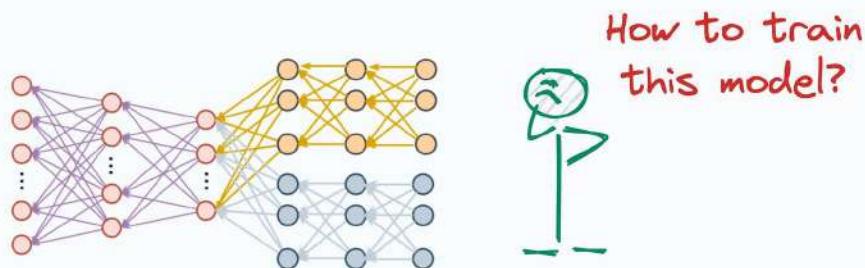
Furthermore, federated learning distributes most computation to a user's device.

As a result, the central server does not need the enormous computing that it would have demanded otherwise.

This is the core idea behind federated learning.

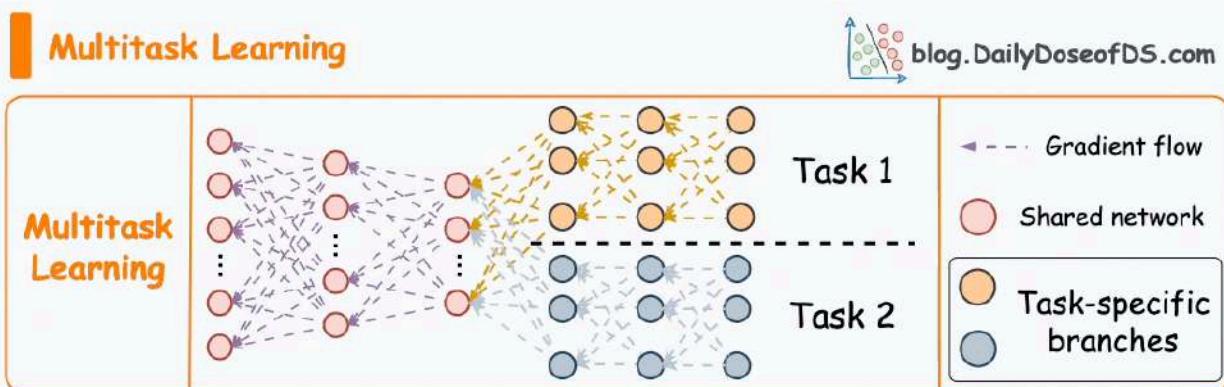
Building Multi-task Learning (MTL) Models

Most ML models are trained on one task. As a result, many struggle to intuitively understand how a model can be trained on multiple tasks simultaneously.



So let's discuss it in this chapter.

To reiterate, in MTL, the network has a few shared layers and task-specific segments. During backpropagation, gradients are accumulated from all branches, as depicted in the animation below:



Let's take a simple example to understand its implementation.

Consider we want our model to take a real value (x) as input and generate two outputs:

- $\sin(x)$
- $\cos(x)$

This can be formulated as an MTL problem.

First, we define our model class using PyTorch.

```

class Net(nn.Module):
    def __init__(self, h):
        super(Net, self).__init__()

        self.model = torch.nn.Sequential(
            torch.nn.Linear(1, h),
            torch.nn.ReLU(),
            torch.nn.Linear(h, h),
            torch.nn.ReLU()
        )

        self.model_sin = torch.nn.Sequential(
            torch.nn.Linear(h, h),
            torch.nn.ReLU(),
            torch.nn.Linear(h, 1)
        )

        self.model_cos = torch.nn.Sequential(
            torch.nn.Linear(h, h),
            torch.nn.ReLU(),
            torch.nn.Linear(h, 1))
    
```

- We have some fully connected layers in `self.model` → These are the shared layers.
- Furthermore, we have the output-specific layers to predict $\sin(x)$ and $\cos(x)$.

Next, let's define the forward pass in the class above:

- First, we pass the input through the shared layers (`self.model`).
- The output of the shared layers is passed through the sin and cos branches.
- We return the output from both branches.

```

class Net(nn.Module):
    def __init__(self, h):
        super(Net, self).__init__()

        self.model = ...
        self.model_sin = ...
        self.model_cos = ...

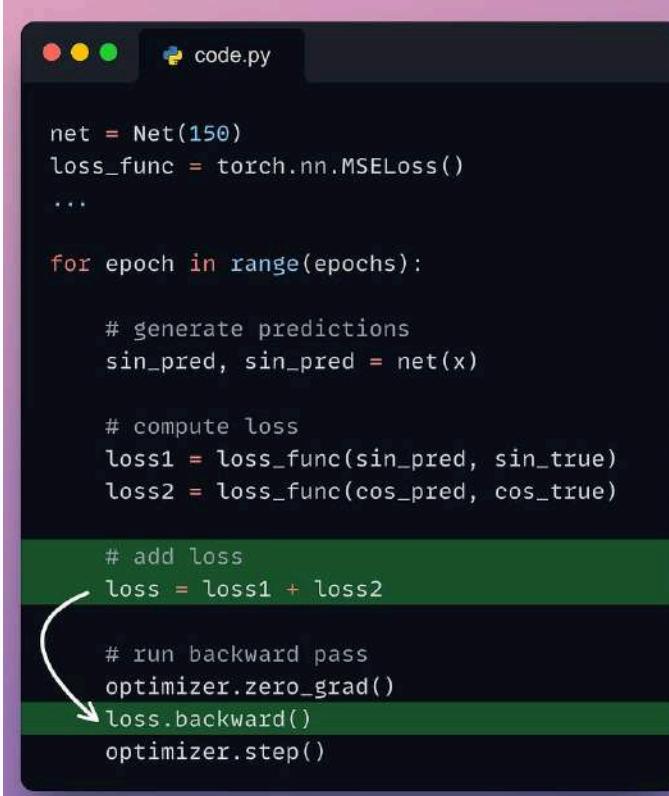
    def forward(self, inputs):
        # pass through shared layers
        x1 = self.model(inputs)

        # generate sin(x) prediction
        output_sin = self.model_sin(x1)

        # generate cos(x) prediction
        output_cos = self.model_cos(x1)

        # return both predictions
        return output_sin, output_cos
    
```

We are almost done. The final part of this implementation is to train the model. Let's use mean squared error as the loss function. The training loop is implemented below:



Network training loop

```

net = Net(150)
loss_func = torch.nn.MSELoss()
...

for epoch in range(epochs):

    # generate predictions
    sin_pred, cos_pred = net(x)

    # compute loss
    loss1 = loss_func(sin_pred, sin_true)
    loss2 = loss_func(cos_pred, cos_true)

    # add loss
    loss = loss1 + loss2

    # run backward pass
    optimizer.zero_grad()
    loss.backward()
    optimizer.step()

```

Decreasing loss

Epoch	Loss
0	1.014
5	0.947
10	0.803
15	0.543
20	0.438
25	0.181
30	0.132
35	0.077
40	0.045
45	0.04
50	0.026

- We pass the input data through the model.
- It returns two outputs, one from each segment of the network.
- We compute the branch-specific loss values (loss1 and loss2) using true predictions.
- We add the two loss values to get the total loss for the network.
- Finally, we run the backward pass.

With this, we have trained our MTL model. Also, we get a decreasing loss, which depicts that the model is being trained.

And that's how we train an MTL model. You can extend the same idea to build any MTL model of your choice.

Do remember that building an MTL model on unrelated tasks will not produce good results.

Thus, “task-relatedness” is a critical component of all MTL models because of the shared layers. Also, it is NOT necessary that every task must equally contribute to the entire network’s loss.

We may assign weights to each task as well, as depicted below:

$$\text{loss} = c1 * \text{loss_task1} + c2 * \text{loss_task2}$$

weighted loss

The weights could be based on task importance.

Or...

At times, I also use dynamic task weights, which could be inversely proportional to the validation accuracy achieved on that task.

$$c1 \propto \frac{1}{\text{accuracy_task1}} \quad c2 \propto \frac{1}{\text{accuracy_task2}}$$

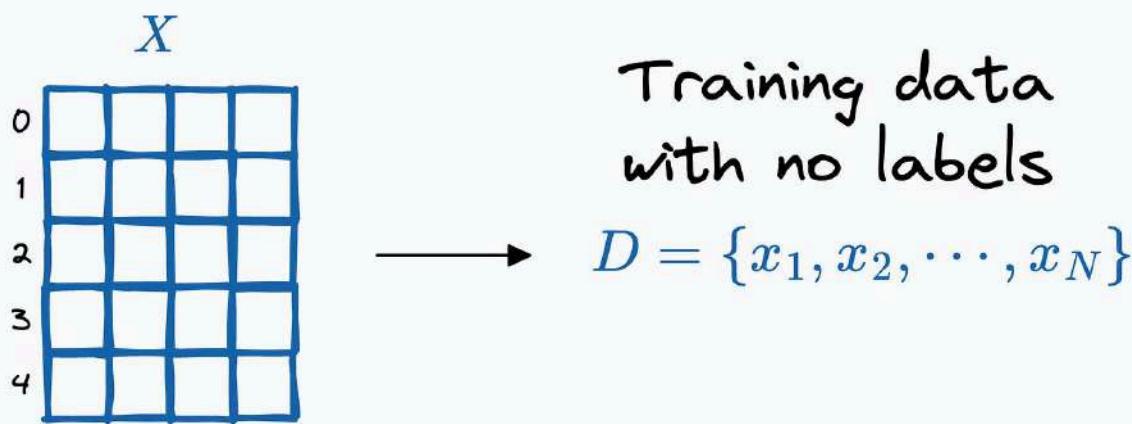
My rationale behind this technique is that in an MTL setting, some tasks can be easy while others can be difficult.

If the model achieves high accuracy on one task during training, we can safely reduce its loss contribution so that the model focuses more on the second task.

You can download the notebook for this chapter here: <https://bit.ly/3ztY5hy>.

Active Learning

There's not much we can do to build a supervised system when the data we begin with is unlabeled.



Using unsupervised techniques (if they fit the task) can be a solution, but supervised systems are typically on par with unsupervised ones.

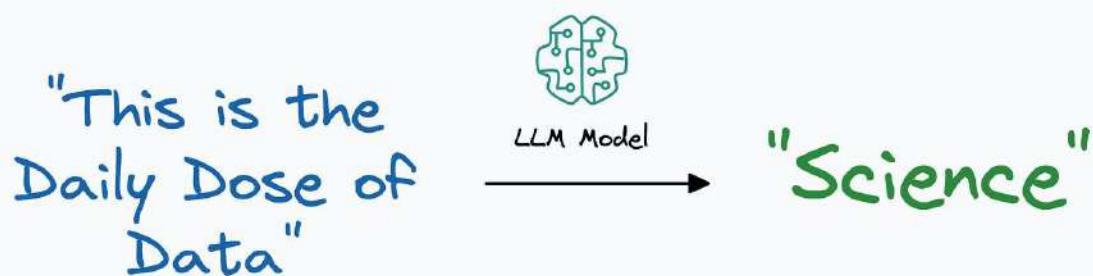
Another way, if feasible, is to rely on self-supervised learning.

Self-supervised learning is when we have an unlabeled dataset (say text data), but we somehow figure out a way to build a supervised learning model out of it.

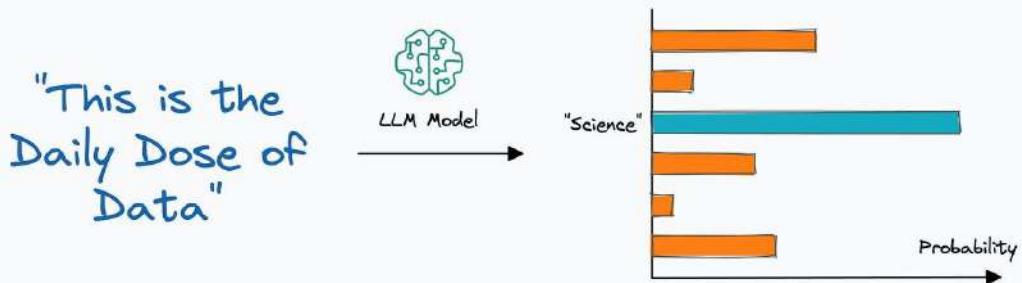
This becomes possible due to the inherent nature of the task.

Consider an LLM, for instance.

In a nutshell, its core objective is to predict the next token based on previously predicted tokens (or the given context).

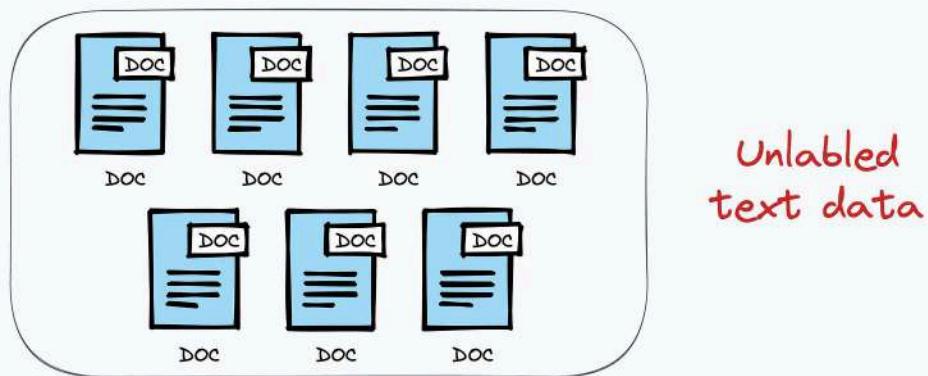


This is a classification task, and the labels are tokens.



But text data is raw. It has no labels.

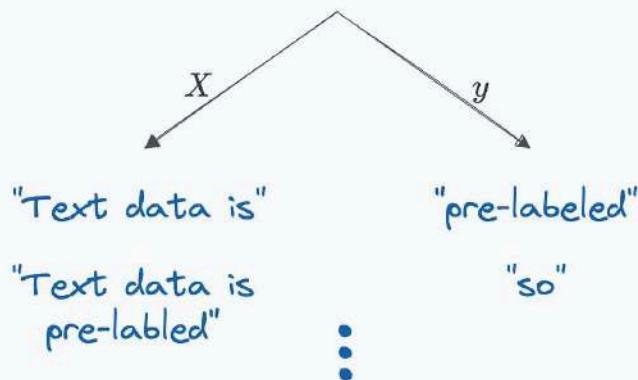
Then how did we train this classification task?



Self-supervised techniques solve this problem.

Due to the inherent nature of the task (next-token prediction, to be specific), every piece of raw text data is already self-labeled.

"Text data is pre-labeled so it does not need any labeling"



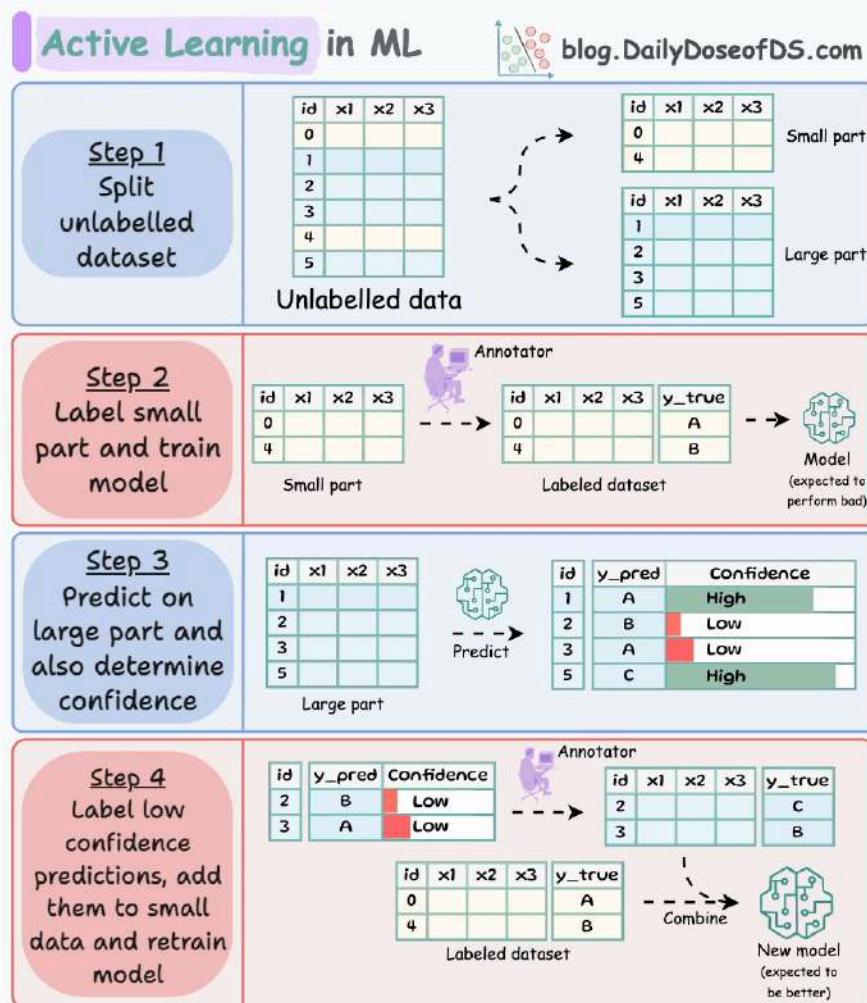
The model is only supposed to learn the mapping from previous tokens to the next token.

This is called self-supervised learning, which is quite promising, but it has limited applicability, largely depending on the task.

At this stage, the only possibility one notices is annotating the dataset. However, data annotation is difficult, expensive, time-consuming, and tedious.

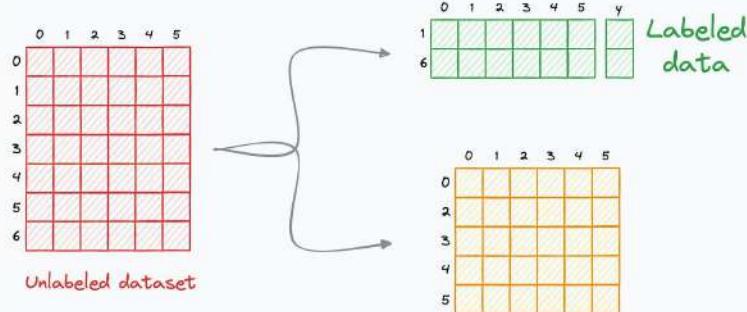
Active learning is a relatively easy, inexpensive, quick, and interesting way to address this.

As the name suggests, the idea is to build the model with active human feedback on examples it is struggling with. The visual below summarizes this:



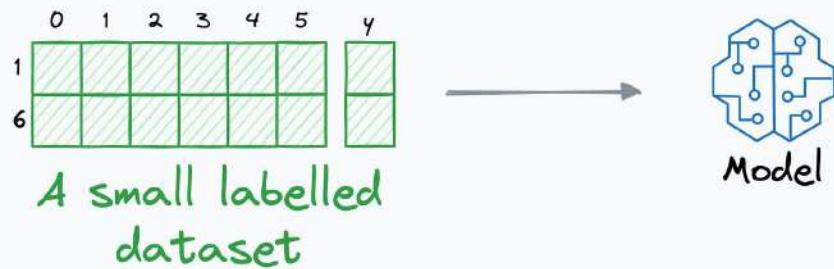
Let's get into the details.

We begin by manually labeling a tiny percentage of the dataset.

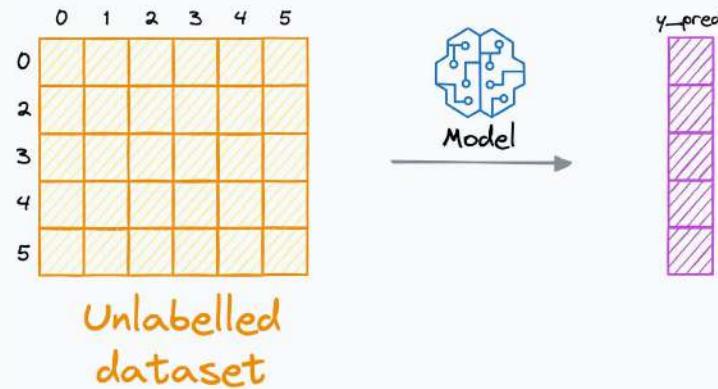


While there's no rule on how much data should be labeled, I have used active learning (successfully) while labeling as low as ~1% of the dataset, so try something in that range.

Next, build a model on this small labeled dataset.



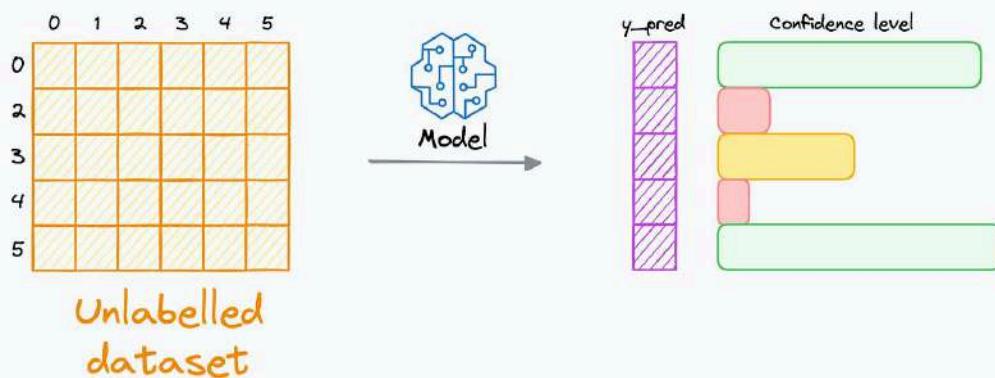
Of course, this won't be a perfect model, but that's okay. Next, generate predictions on the dataset we did not label:



It's obvious that we cannot determine if these predictions are correct as we do not have any labels.

That's why we need to be a bit selective with the type of model we choose.

More specifically, we need a model that, either implicitly or explicitly, can also provide a confidence level with its predictions.

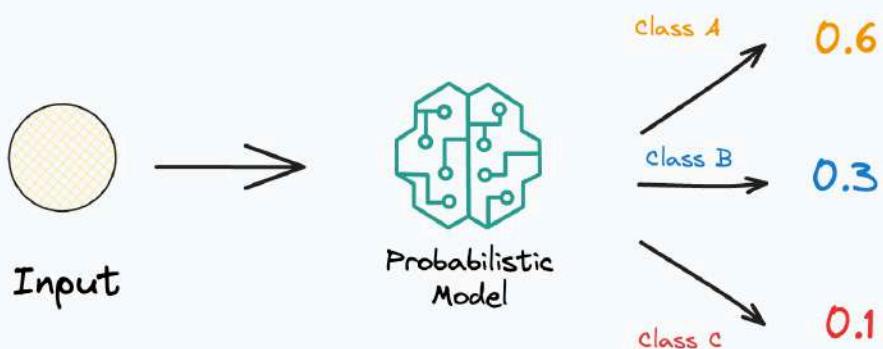


As the name suggests, a confidence level reflects the model's confidence in generating a prediction.

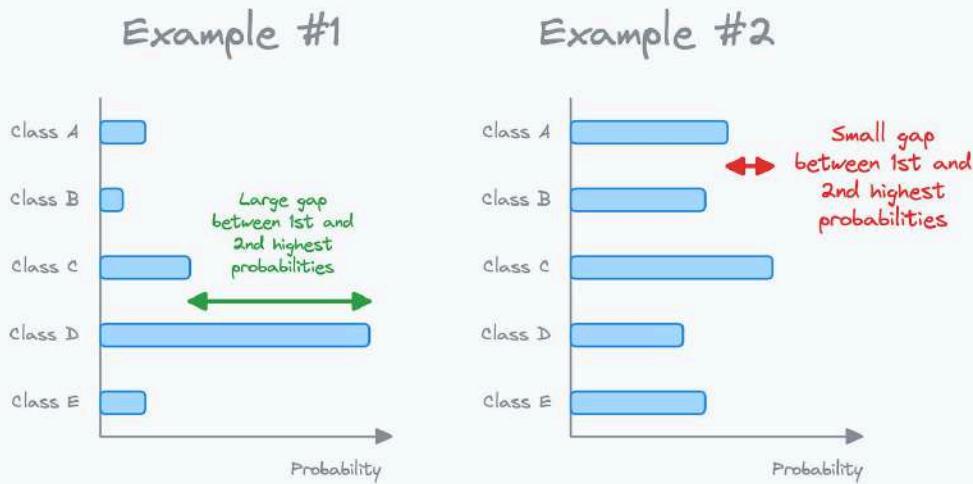
If a model could speak, it would be like:

- I am predicting a “cat” and am 95% confident about my prediction.
- I am predicting a “cat” and am 5% confident about my prediction.
- And so on...

Probabilistic models (ones that provide a probabilistic estimate for each class) are typically a good fit here.



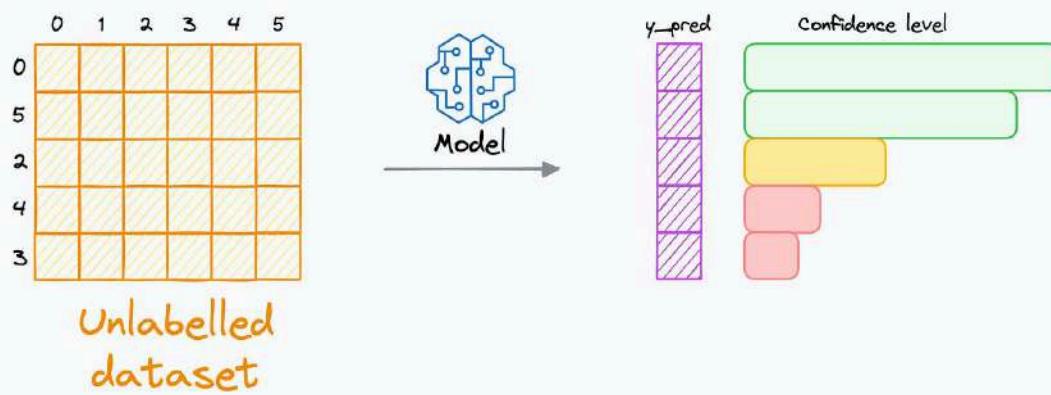
This is because one can determine a proxy for confidence level from probabilistic outputs.



In the above two examples, consider the gap between 1st and 2nd highest probabilities:

- In example #1, the gap is large. This can indicate that the model is quite confident in its prediction.
- In example #2, the gap is small. This can indicate that the model is NOT quite confident in its prediction.

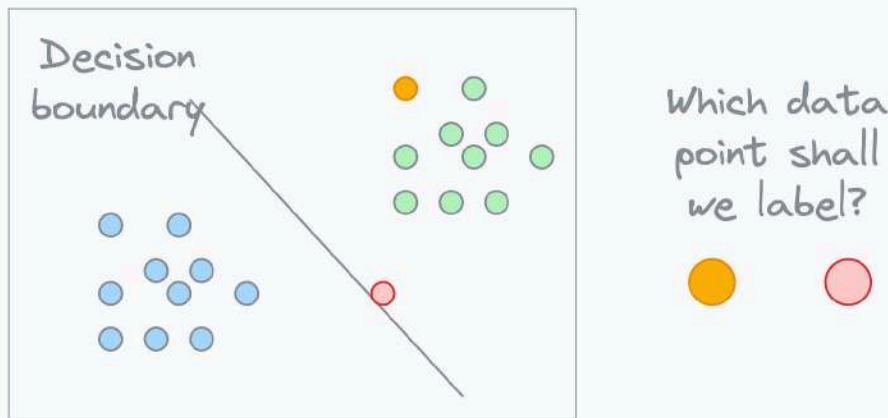
Now, go back to the predictions generated above and rank them in order of confidence:



In the above image:

- The model is already quite confident with the first two instances. There's no point checking those.
- Instead, it would be best if we (the human) annotate the instances with which it is least confident.

To get some more perspective, consider the image below. Logically speaking, which data point's human label will provide more information to the model? I know you already know the answer.



Thus, in the next step, we provide our human label to the low-confidence predictions and feed it back to the model with the previously labeled dataset:



Repeat this a few times and stop when you are satisfied with the performance.

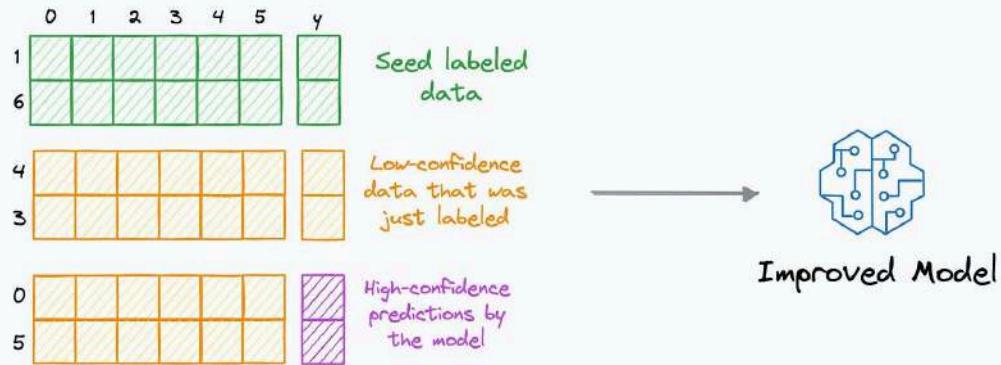
In my experience, active learning has always been an immensely time-saving approach to building supervised models on unlabeled datasets.

The only thing that you have to be careful about is generating confidence measures.

If you mess this up, it will affect every subsequent training step.

There's one more thing I like to do when using active learning.

While combining the low-confidence data with the seed data, we can also use the high-confidence data. The labels would be the model's predictions.

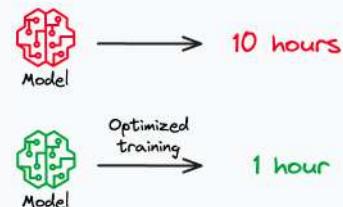


This variant of active learning is called cooperative learning.

Run-time and Memory Optimization

Momentum

As we progress towards building larger and larger models, every bit of possible optimization becomes crucial.



And, of course, there are various ways to speed up model training, like:

- Batch processing
- Leverage distributed training using frameworks like PySpark MLLib.
- Use better Hyperparameter Optimization, like Bayesian Optimization, which we will discuss in [this chapter](#).
- and many other techniques.

Momentum is another reliable and effective technique to speed up model training. While Momentum is pretty popular, many people struggle to intuitively understand how it works and why it is effective. Let's understand in this chapter.

Issues with Gradient Descent

In gradient descent, every parameter update solely depends on the current gradient. This is clear from the gradient weight update rule shown below:

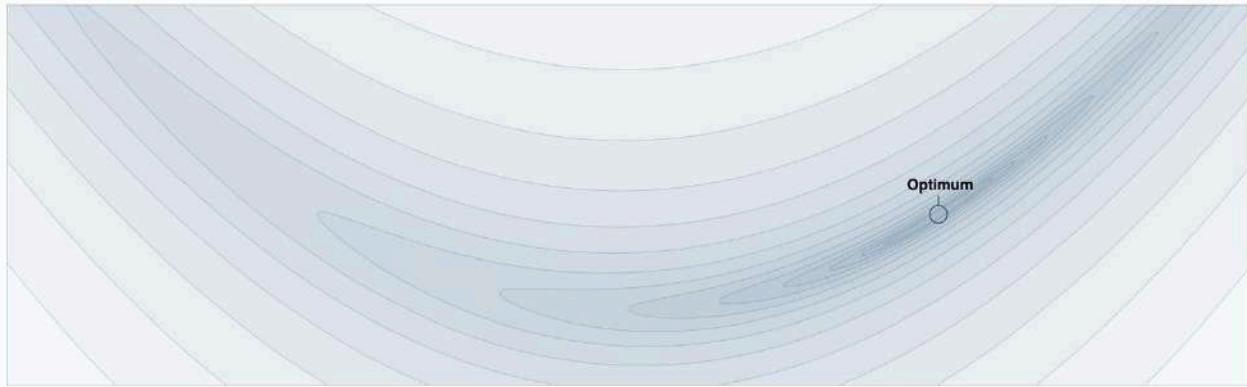
$$\theta^t = \theta^{t-1} - \alpha \frac{\delta J}{\delta \theta}$$

Only dependent on current gradient

As a result, we end up having many unwanted oscillations during the optimization process.

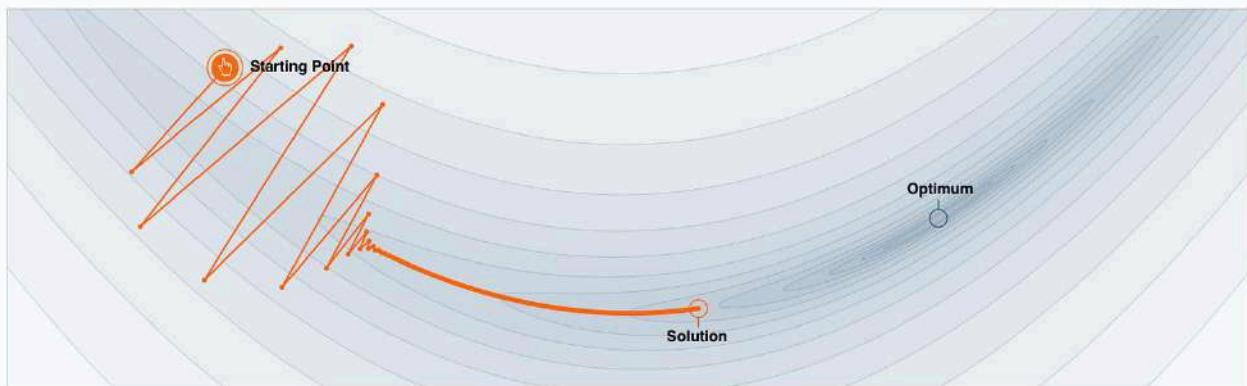
Let's understand this more visually.

Imagine this is the loss function contour plot, and the optimal location (parameter configuration where the loss function is minimum) is marked here:



Simply put, this plot illustrates how gradient descent moves towards the optimal solution. At each iteration, the algorithm calculates the gradient of the loss function at the current parameter values and updates the weights.

This is depicted below:



Notice two things here:

- It unnecessarily oscillates vertically.
- It ends up at the non-optimal solution after some epochs.

Ideally, we would have expected our weight updates to look this:



- It must have taken longer steps in the horizontal direction...
- ...and smaller vertical steps because a movement in this direction is unnecessary.

This idea is also depicted below:



How Momentum solves this problem?

Momentum-based optimization slightly modifies the update rule of gradient descent. More specifically, it also considers a moving average of past gradients:

$$\theta^t = \theta^{t-1} - \alpha \frac{\delta J}{\delta \theta} + \beta \cdot F(\theta^{t-1}, \theta^{t-2}, \dots)$$

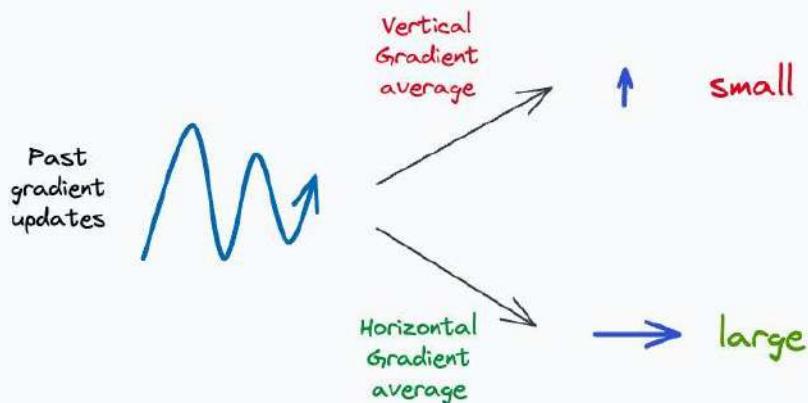
dependent on current gradient

dependent on previous updates

This helps us handle the unnecessary vertical oscillations we saw earlier.

How?

As Momentum considers a moving average of past gradients, so if the recent gradient update trajectory looks as shown in the following image, then it is clear that its average in the vertical direction will be very low while that in the horizontal direction will be large (which is precisely what we want):



As this moving average gets added to the gradient updates, it helps the optimization algorithm take larger steps in the desired direction.

$$\theta^t = \theta^{t-1} - \alpha \frac{\delta J}{\delta \theta} + \beta \cdot F(\theta^{t-1}, \theta^{t-2}, \dots)$$

↑ ↑

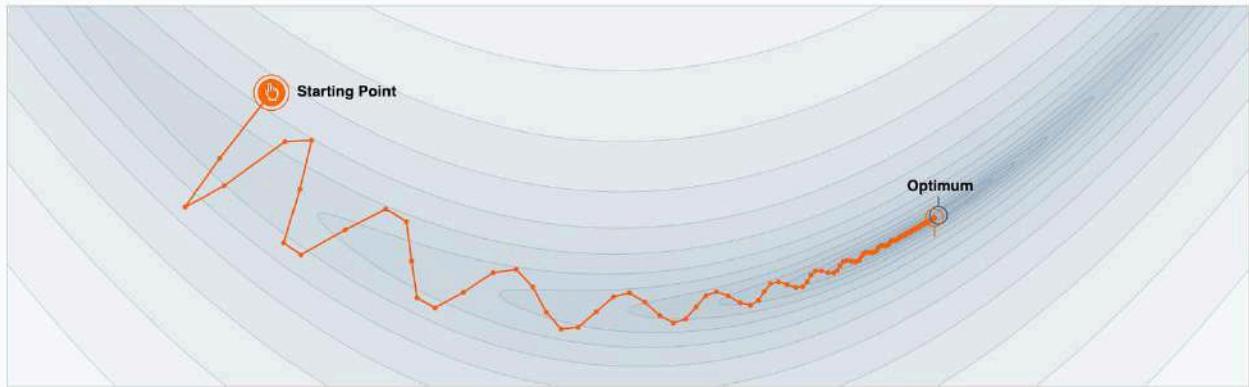
dependent on current gradient *dependent on previous updates*

This way, we can:

- Smoothen the optimization trajectory.

- Reduce unnecessary oscillations in parameter updates, which also speeds up training.

This is also evident from the image below:

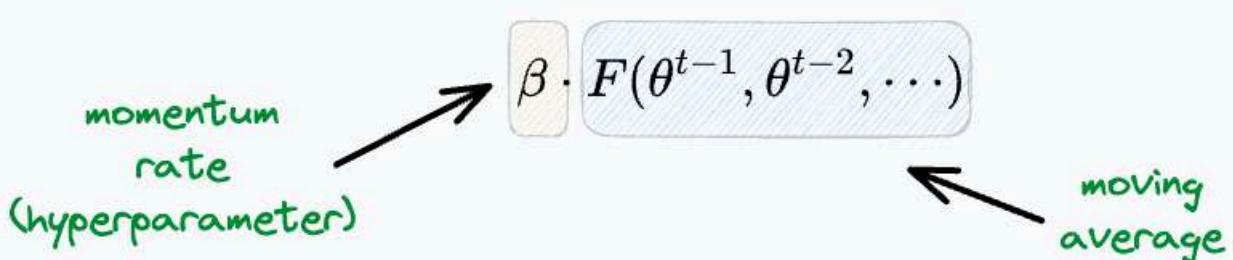


This time, the gradient update trajectory shows much smaller oscillations in the vertical direction, and it also manages to reach an optimum under the same number of epochs as earlier.

This is the core idea behind Momentum and how it works.

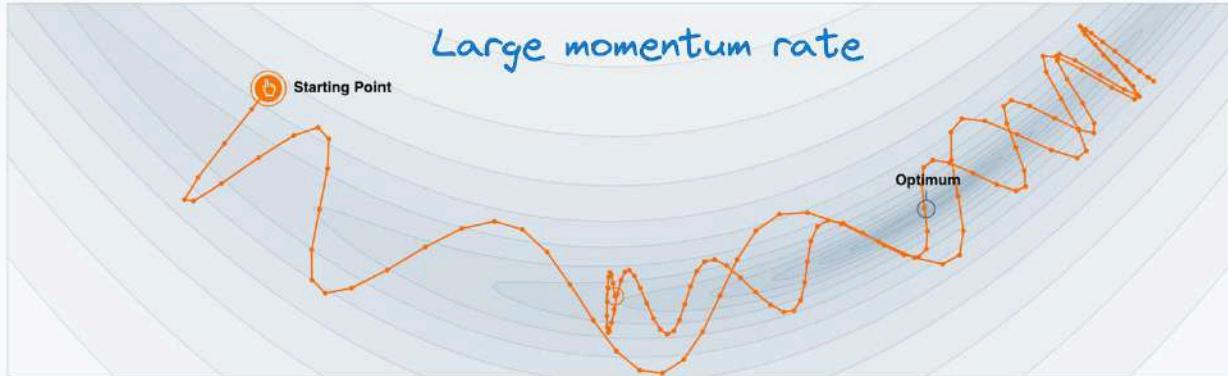
Of course, Momentum does introduce another hyperparameter (Momentum rate) in the model, which should be tuned appropriately like any other hyperparameter:

Momentum

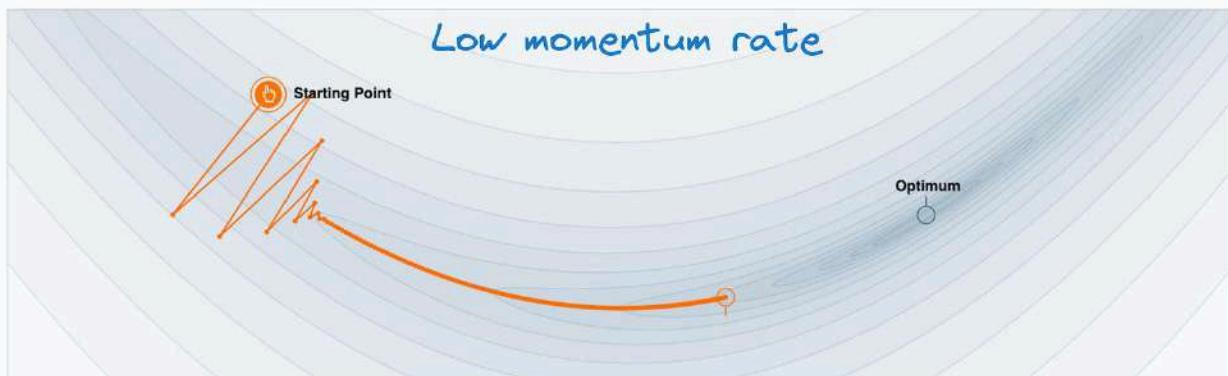


For instance, considering the 2D contours we discussed above:

- Setting an extremely large value of Momentum rate will significantly expedite gradient update in the horizontal direction. This may lead to overshooting the minima, as depicted below:



- What's more, setting an extremely small value of Momentum will slow down the optimal gradient update, defeating the whole purpose of Momentum.



If you want to have a more hands-on experience, check out this tool:

<https://bit.ly/4cOrJN1>.

Mixed Precision Training

Context

Typical deep learning libraries are really conservative when it comes to assigning data types.

The data type assigned by default is usually 64-bit or 32-bit, when there is also scope for 16-bit, for instance. This is also evident from the code below:

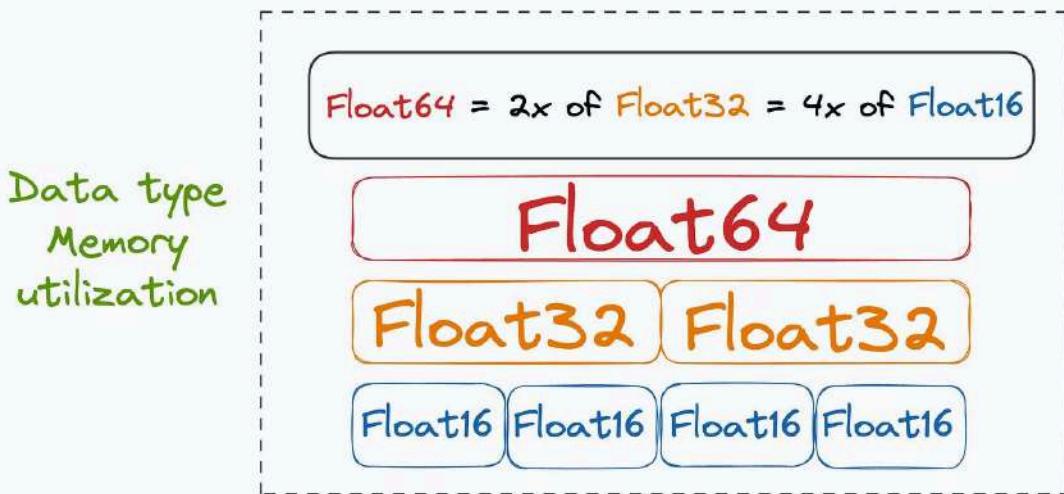
```
import torch
tensor = torch.nn.Parameter(torch.randn(5,5))
>>> tensor.dtype
torch.float32
```

PyTorch assigns a large data type

As a result, we are not entirely optimal at efficiently allocating memory. Of course, this is done to ensure better precision in representing information.

Lowest precision <i>→</i>	<pre>a = np.array([0.123456789123456789], dtype=np.float16) >>> print("16-bits: ", a[0]) "16-bits: 0.1235"</pre>
Intermediate precision <i>→</i>	<pre>b = np.array([0.123456789123456789], dtype=np.float32) >>> print("32-bits: ", b[0]) "32-bits: 0.12345679"</pre>
Highest precision <i>→</i>	<pre>c = np.array([0.123456789123456789], dtype=np.float64) >>> print("64bit: ", c[0]) "64-bit: 0.12345678912345678"</pre>

However, this precision always comes at the cost of additional memory utilization, which may not be desired in all situations.



In fact, it is also observed that many tensor operations, especially matrix multiplication, are much faster when we operate under smaller precision data types than larger ones, as demonstrated below:

Float 32 matrix multiplication

```
| matrixA = torch.randn(1000, 1000, dtype=torch.float32, device=device)
matrixB = torch.randn(1000, 1000, dtype=torch.float32, device=device)

| %timeit torch.matmul(matrixA, matrixB)
658 µs ± 12.4 µs per loop
```

Float 16 matrix multiplication

```
| matrixA = torch.randn(1000, 1000, dtype=torch.float16, device=device)
matrixB = torch.randn(1000, 1000, dtype=torch.float16, device=device)

| %timeit torch.matmul(matrixA, matrixB)
111 µs ± 4.07 µs per loop
```

6x Faster

Moreover, since `float16` is only half the size of `float32`, its usage reduces the memory required to train the network. This also allows us to train larger models, train on larger mini-batches (resulting in even more speedup), etc.

Mixed precision training is a pretty reliable and widely adopted technique in the industry to achieve this.

As the name suggests, the idea is to employ lower precision `float16` (wherever feasible, like in convolutions and matrix multiplications) along with `float32` — that is why the name “mixed precision.”

This is a list of some models I found that were trained using mixed precision:

Image Classification	Detection / Segmentation	Generative Models (Images)	Language Modeling
AlexNet	DeepLab	DLSS	BERT
DenseNet	Faster R-CNN	Vid2vid	GPT
Inception	Mask R-CNN	GauGAN	TrellisNet
MobileNet	SSD	Partial Image Inpainting	Gated Convolutions
EfficientNet	NVIDIA Automotive	Progress GAN	BigLSTM/mLSTM
ResNet	RetinaNet	Pix2Pix	RoBERTa
ResNeXt	UNET		Transformer XL
ShuffleNet	DETR		
SqueezeNet			
VGG			
Xception			
Dilated ResNet			
Stacked U-Net			
Recommendation		Speech	Translation
	DeepRecommender	Deep Speech 2	Convolutional Seq2Seq
	DLRM	Jasper	Dynamic Convolutions
	NCF	Tacotron	GNMT (RNN)
		Wave2vec	Levenshtein Transformer
		WaveNet	Transformer (Self-Attention)
		WaveGlow	

It's pretty clear that mixed precision training is much more popularly used, but we don't get to hear about it often.

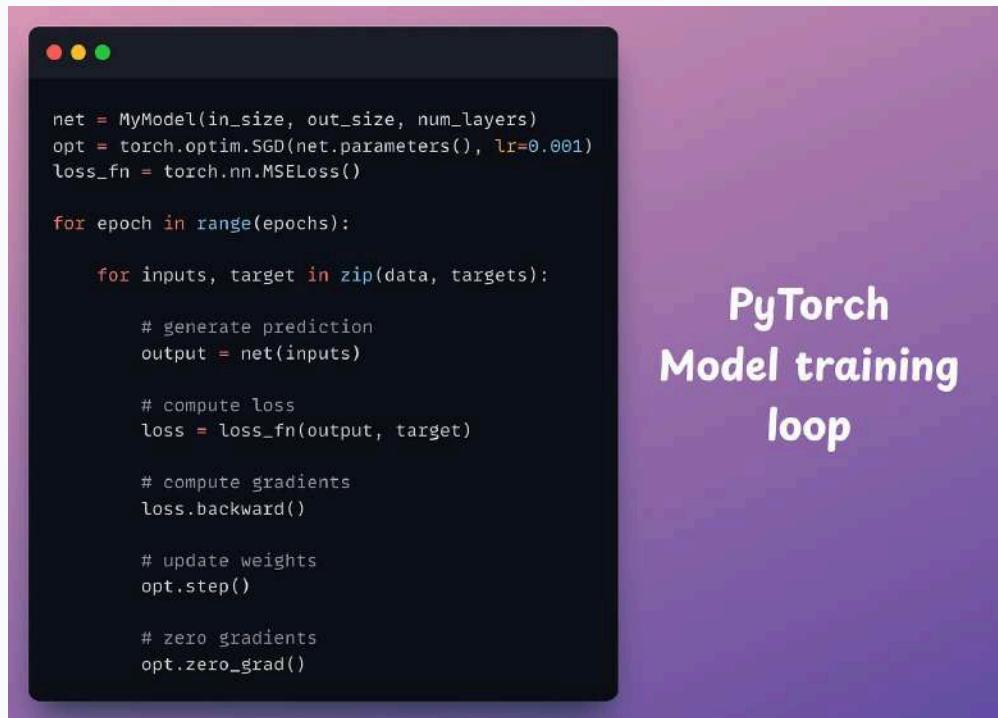
Before we get into the technical details...

From the above discussion, it must be clear that as we use a low-precision data type (`float16`), we might unknowingly introduce some numerical inconsistencies and inaccuracies.

To avoid them, there are some best practices for mixed precision training that I want to talk about next, along with the code.

Mixed precision training in PyTorch and Best Practices

Leveraging mixed precision training in PyTorch requires a few modifications in the existing network training implementation. Consider this is our current PyTorch model training implementation:



The diagram shows a dark rectangular window representing a terminal or code editor. Inside, Python code for a PyTorch training loop is displayed. To the right of the window, the text "PyTorch Model training loop" is written in white, with "PyTorch" and "loop" being the most prominent. A curved arrow points from the word "scaler" in the code towards the "Loss scaler" text.

```

net = MyModel(in_size, out_size, num_layers)
opt = torch.optim.SGD(net.parameters(), lr=0.001)
loss_fn = torch.nn.MSELoss()

for epoch in range(epochs):

    for inputs, target in zip(data, targets):

        # generate prediction
        output = net(inputs)

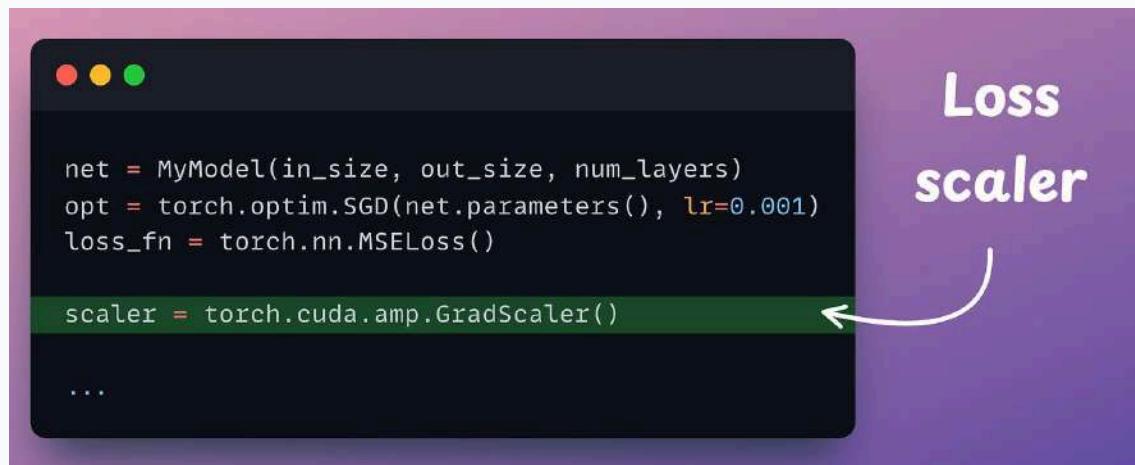
        # compute loss
        loss = loss_fn(output, target)

        # compute gradients
        loss.backward()

        # update weights
        opt.step()

        # zero gradients
        opt.zero_grad()
    
```

The first thing we introduce here is a scaler object that will scale the loss value:



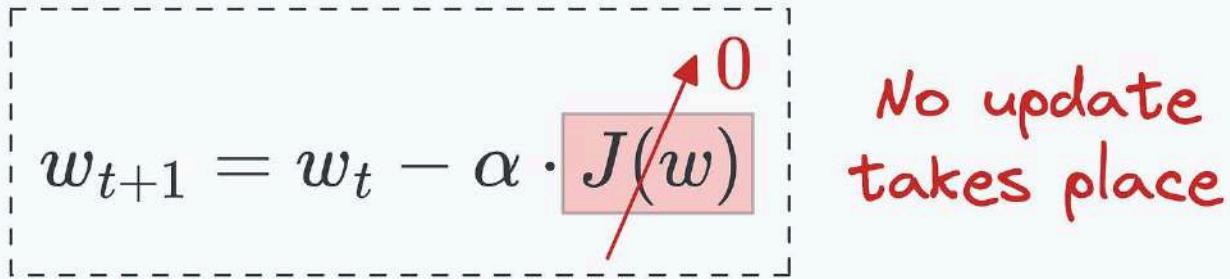
The diagram shows a dark rectangular window representing a terminal or code editor. Inside, Python code for a PyTorch training loop is displayed. A green horizontal bar highlights the line "scaler = torch.cuda.amp.GradScaler()". To the right of the window, the text "Loss scaler" is written in white, with "Loss" and "scaler" being the most prominent. A curved arrow points from the word "scaler" in the code towards the "Loss scaler" text.

```

net = MyModel(in_size, out_size, num_layers)
opt = torch.optim.SGD(net.parameters(), lr=0.001)
loss_fn = torch.nn.MSELoss()

scaler = torch.cuda.amp.GradScaler()
...
```

We do this because, at times, the original loss value can be so low, that we might not be able to compute gradients in `float16` with full precision. Such situations may not produce any update to the model's weights.



Scaling the loss to a higher numerical range ensures that even small gradients can contribute to the weight updates.

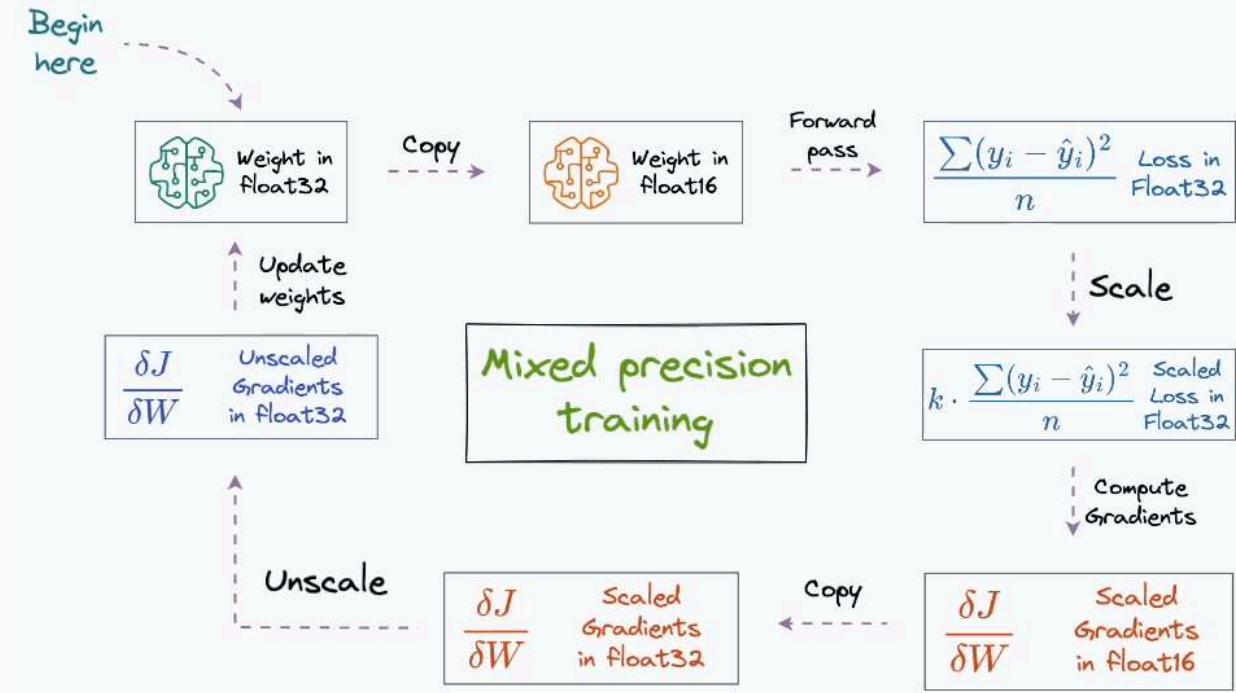
But these minute gradients can only be accommodated into the weight matrix when the weight matrix itself is represented in high precision, i.e., `float32`. Thus, as a conservative measure, we tend to keep the weights in `float32`.

That said, the loss scaling step is not entirely necessary because, in my experience, these little updates typically appear towards the end stages of the model training. Thus, it can be fair to assume that small updates may not drastically impact the model performance. But don't take this as a definite conclusion, so it's something that I want you to validate when you use mixed precision training.

Moving on, as the weights (which are matrices) are represented in `float32`, we can not expect the speedup from representing them in `float16`, if they remain this way:

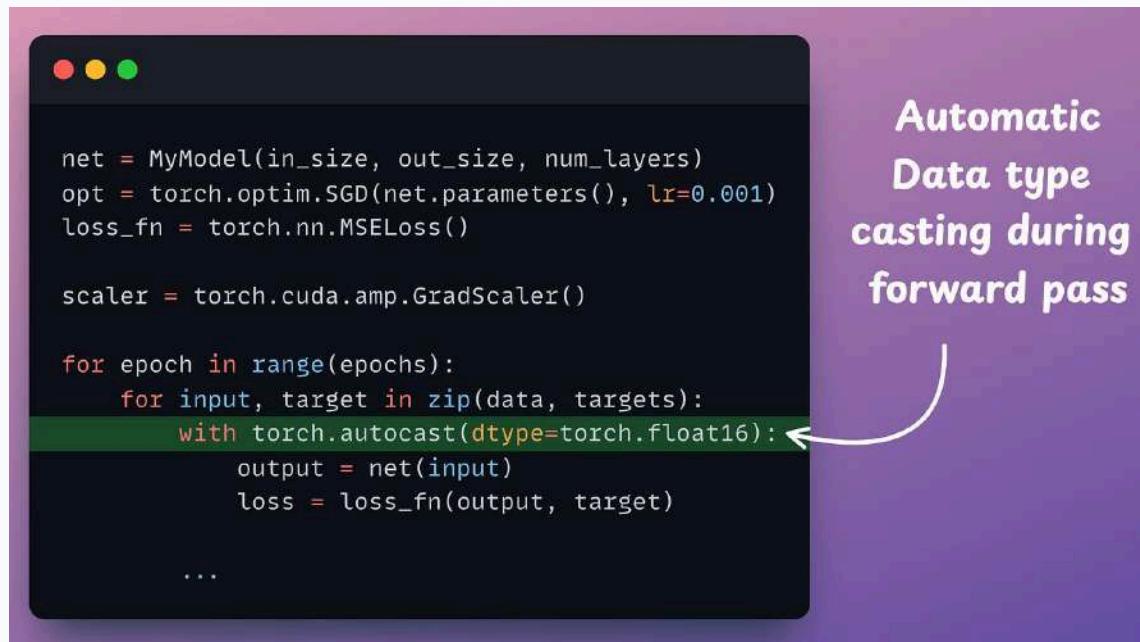
<pre>Float 32 matrix multiplication</pre> <pre>matrixA = torch.randn(1000, 1000, dtype=torch.float32, device=device) matrixB = torch.randn(1000, 1000, dtype=torch.float32, device=device) %timeit torch.matmul(matrixA, matrixB) 658 µs ± 12.4 µs per loop</pre> <pre>Float 16 matrix multiplication</pre> <pre>matrixA = torch.randn(1000, 1000, dtype=torch.float16, device=device) matrixB = torch.randn(1000, 1000, dtype=torch.float16, device=device) %timeit torch.matmul(matrixA, matrixB) 111 µs ± 4.07 µs per loop</pre>	6x Faster
--	------------------

To leverage these float16-based speedups, here are the steps we follow:



1. We make a **float16** copy of weights during the forward pass.
2. Next, we compute the loss value in **float32** and scale it to have more precision in gradients, which works in **float16**.
 - a. The reason we compute gradients in **float16** is because, like forward pass, gradient computations also involve matrix multiplications.
 - b. Thus, keeping them in **float16** can provide additional speedup.
3. Once we have computed the gradients in **float16**, the heavy matrix multiplication operations have been completed. Now, all we need to do is update the original weight matrix, which is in **float32**.
4. Thus, we make a **float32** copy of the above gradients, remove the scale we applied in Step 2, and update the **float32** weights.
5. Done!

The mixed-precision settings in the forward pass are carried out by the `torch.autocast()` context manager:



```

net = MyModel(in_size, out_size, num_layers)
opt = torch.optim.SGD(net.parameters(), lr=0.001)
loss_fn = torch.nn.MSELoss()

scaler = torch.cuda.amp.GradScaler()

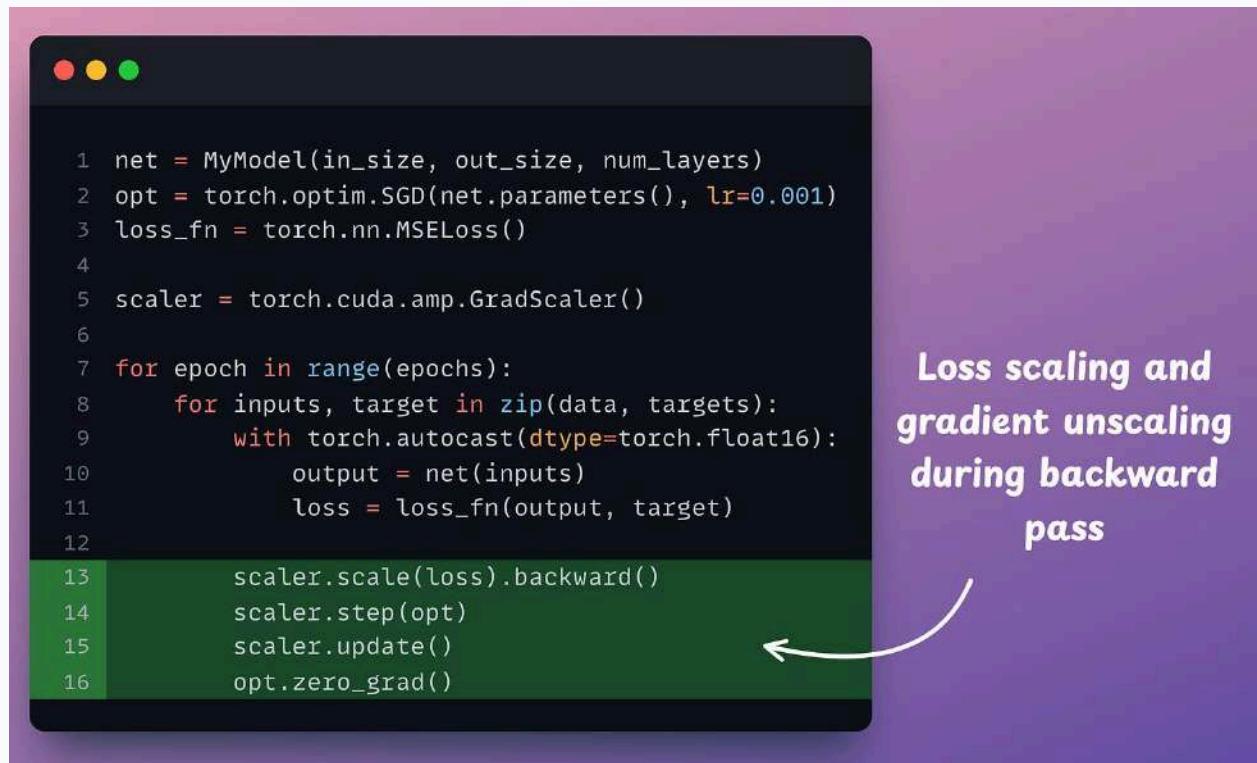
for epoch in range(epochs):
    for input, target in zip(data, targets):
        with torch.autocast(dtype=torch.float16): ←
            output = net(input)
            loss = loss_fn(output, target)

    ...

```

Automatic Data type casting during forward pass

Now, it's time to handle the backward pass.



```

1 net = MyModel(in_size, out_size, num_layers)
2 opt = torch.optim.SGD(net.parameters(), lr=0.001)
3 loss_fn = torch.nn.MSELoss()
4
5 scaler = torch.cuda.amp.GradScaler()
6
7 for epoch in range(epochs):
8     for inputs, target in zip(data, targets):
9         with torch.autocast(dtype=torch.float16): ←
10             output = net(inputs)
11             loss = loss_fn(output, target)
12
13             scaler.scale(loss).backward() ←
14             scaler.step(opt)
15             scaler.update()
16             opt.zero_grad()

```

Loss scaling and gradient unscaling during backward pass

- Line 13 → `scaler.scale(loss).backward()`: The scaler object scales the loss value and `backward()` is called to compute the gradients.
- Line 14 → `scaler.step(opt)`: Unscale gradients and update weights.
- Line 15 → `scaler.update()`: Update the scale for the next iteration.
- Line 16 → `opt.zero_grad()`: Zero gradients.

Done!

The efficacy of mixed precision scaling over traditional training is evident from the image below:

Usual training loop

```
[4] net = make_model(in_size, out_size, num_layers)
    opt = torch.optim.SGD(net.parameters(), lr=0.001)

    start_timer()
    for epoch in range(epochs):
        for input, target in zip(data, targets):
            output = net(input)
            loss = loss_fn(output, target)
            loss.backward()
            opt.step()
            opt.zero_grad()
    end_timer_and_print("Default precision:")

Default precision:
Total execution time = 5.277 sec
```

Execution Time:
5.2 seconds

Mixed precision training

```
[5] net = make_model(in_size, out_size, num_layers)
    opt = torch.optim.SGD(net.parameters(), lr=0.001)

    scaler = torch.cuda.amp.GradScaler()

    start_timer()
    for epoch in range(epochs):
        for input, target in zip(data, targets):
            with torch.autocast(device_type=device, dtype=torch.float16):
                output = net(input)
                loss = loss_fn(output, target)
            scaler.scale(loss).backward()
            scaler.step(opt)
            scaler.update()
            opt.zero_grad()
    end_timer_and_print("Mixed precision:")

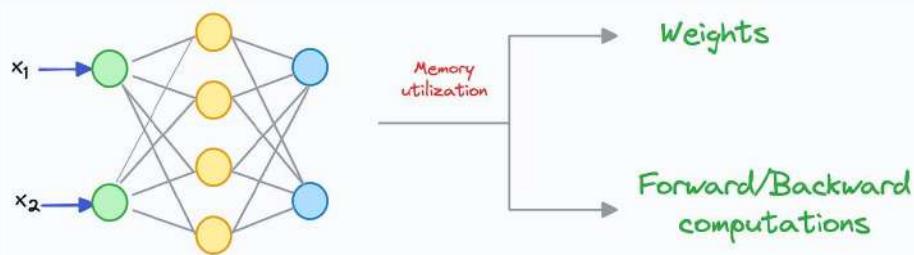
Mixed precision:
Total execution time = 2.076 sec
```

Execution Time:
2 seconds

Mixed precision training is over 2.5x faster than conventional training.

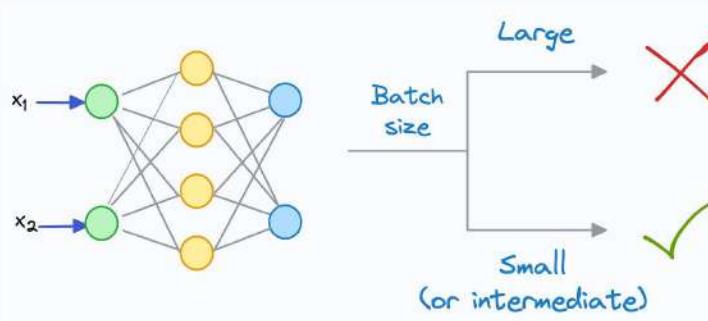
Gradient Checkpointing

Neural networks primarily utilize memory in two ways:



1. When they store model weights (this is fixed memory utilization).
2. When they are trained (this is dynamic). It happens in two ways:
 - a. During forward pass while computing and storing activations of all layers.
 - b. During backward pass while computing gradients at each layer.

The latter, i.e., dynamic memory utilization, often restricts us from training larger models with bigger batch sizes.



This is because memory utilization scales proportionately with the batch size.

That said, there's a pretty incredible technique that lets us increase the batch size while maintaining the overall memory utilization.

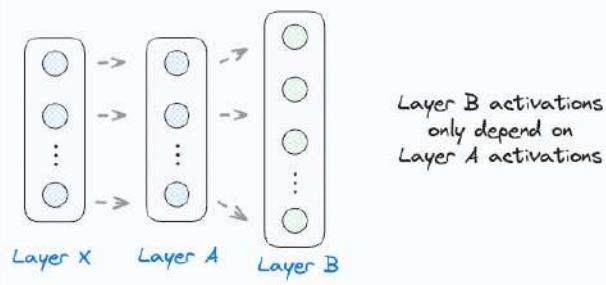
It is called Gradient checkpointing, and in my experience, it's a highly underrated technique to reduce the memory overheads of neural networks.

Let's understand this in more detail.

How gradient checkpointing works?

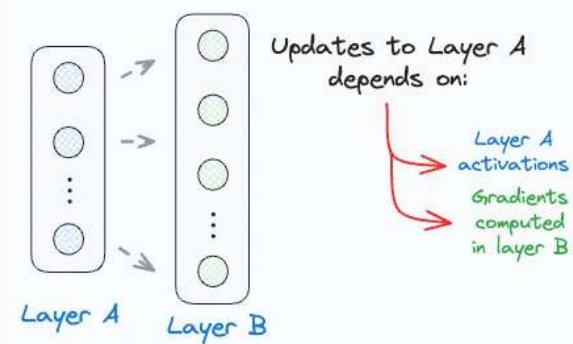
Gradient checkpointing is based on two key observations on how neural networks typically work:

- 1) The activations of a specific layer can be solely computed using the activations of the previous layer. For instance, in the image below, “Layer B” activations can be computed from “Layer A” activations only.



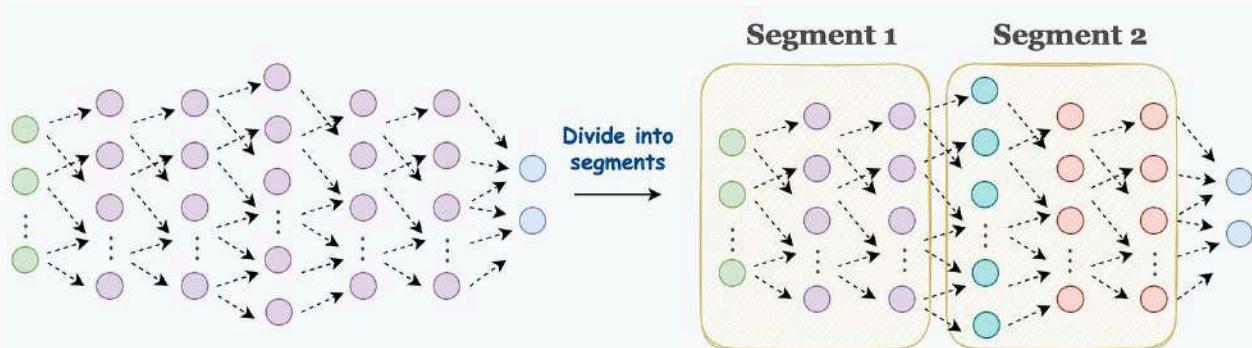
- 2) Updating the weights of a layer only depends on two things:

- The activations of that layer.
- The gradients computed in the next (right) layer (or rather, the running gradients).

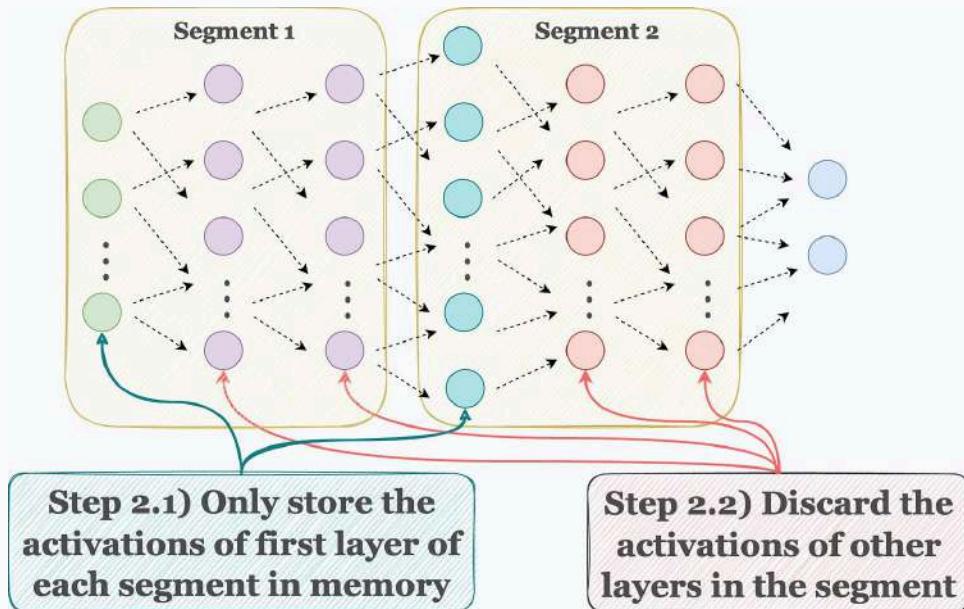


Gradient checkpointing exploits these two observations to optimize memory utilization. Here's how it works:

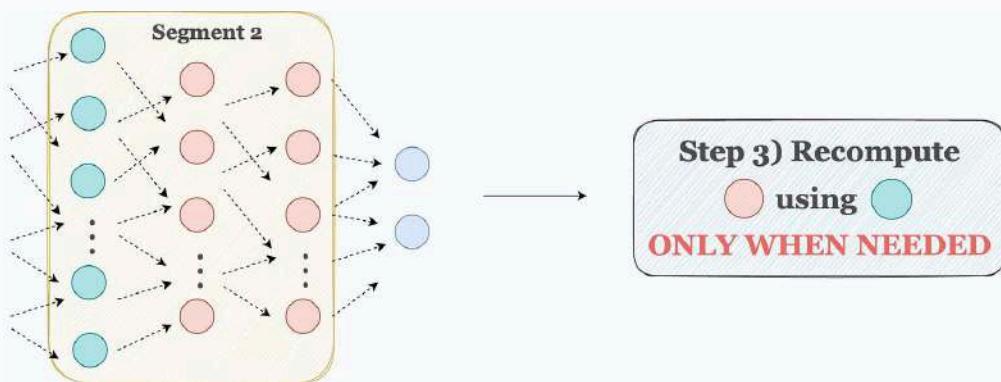
- Step 1) Divide the network into segments before the forward pass:



- Step 2) During the forward pass, only store the activations of the first layer in each segment. Discard the rest when they have been used to compute the activations of the next layer.



- Step 3) Now comes backpropagation. To update the weights of a layer, we need its activations. Thus, we recompute those activations using the first layer in that segment. For instance, as shown in the image below, to update the weights of the red layers, we recompute their activations using the activations of the cyan layer, which are already available in memory.

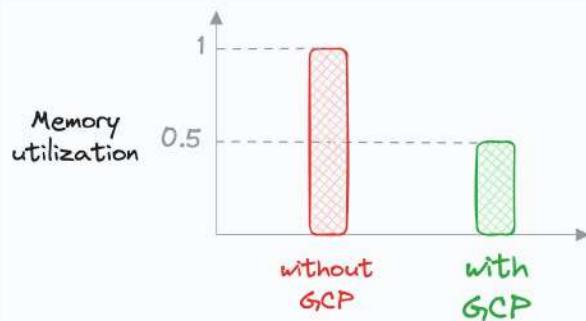


Done!

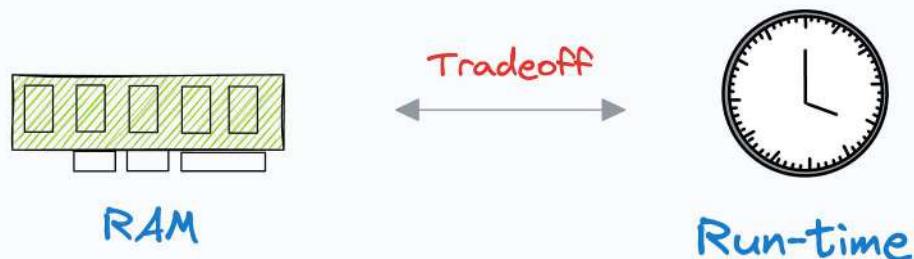
This is how gradient checkpointing works.

To summarize, the idea is that we don't need to store all the intermediate activations in memory. Instead, storing a few of them and recomputing the rest only when they are needed can significantly reduce the memory requirement. The whole idea makes intuitive sense as well.

This also allows us to train the network on larger batches of data. Typically, my observation has been that gradient checkpointing (GCP) can reduce memory usage by at least 50-60%, which is massive.



Of course, as we compute some activations twice, this does come at the cost of increased run-time, which can typically range between 15-25%. So there's always a tradeoff between memory and run-time.



That said, another advantage is that it allows us to use a larger batch size, which can slightly (not entirely though) counter the increased run-time.

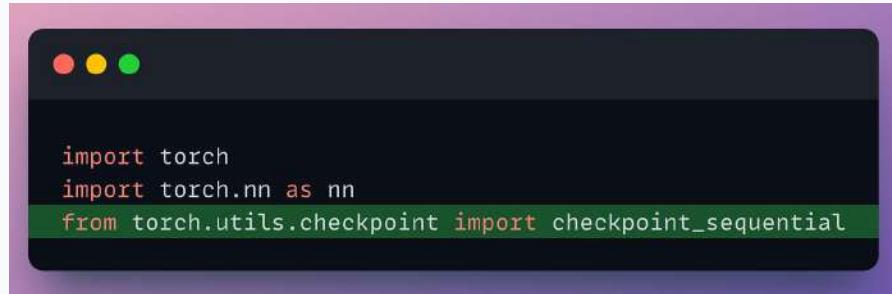
Nonetheless, gradient checkpointing is an extremely powerful technique to train larger models, which I have found to be pretty helpful at times, without resorting to more intensive techniques like distributed training, for instance.

Thankfully, gradient checkpointing is also implemented by many open-source deep learning frameworks like Pytorch, etc.

Here's a demo.

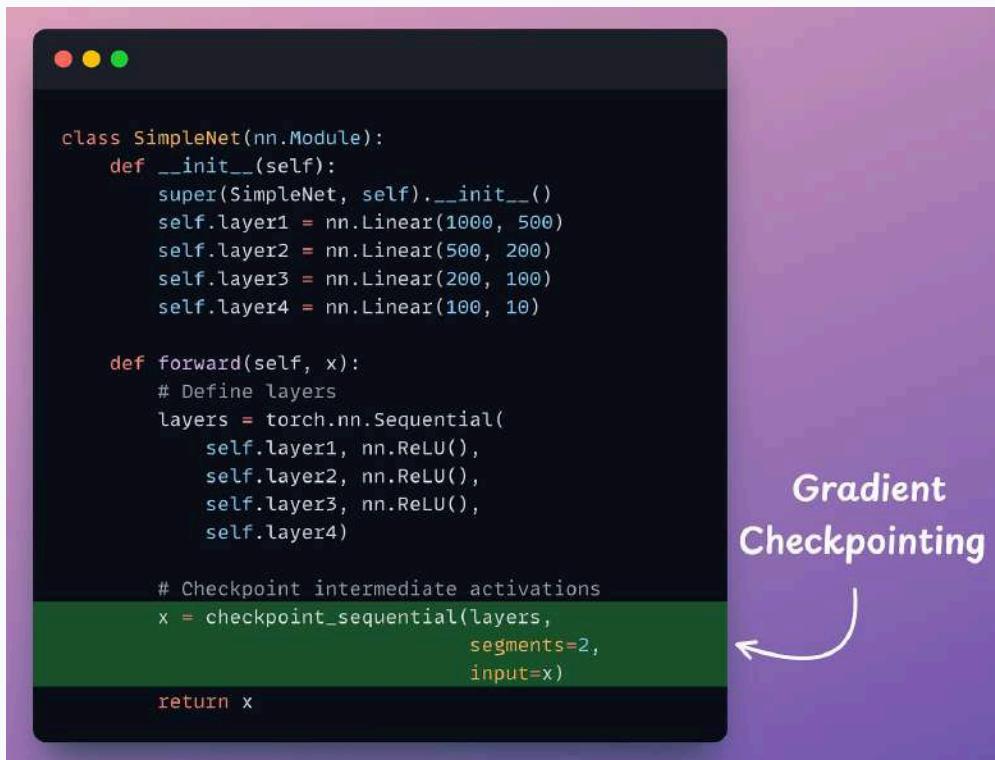
Gradient checkpointing in PyTorch

To utilize this, we begin by importing the necessary libraries and functions:



```
import torch
import torch.nn as nn
from torch.utils.checkpoint import checkpoint_sequential
```

Next, we define our neural network:



```
class SimpleNet(nn.Module):
    def __init__(self):
        super(SimpleNet, self).__init__()
        self.layer1 = nn.Linear(1000, 500)
        self.layer2 = nn.Linear(500, 200)
        self.layer3 = nn.Linear(200, 100)
        self.layer4 = nn.Linear(100, 10)

    def forward(self, x):
        # Define layers
        layers = torch.nn.Sequential(
            self.layer1, nn.ReLU(),
            self.layer2, nn.ReLU(),
            self.layer3, nn.ReLU(),
            self.layer4)

        # Checkpoint intermediate activations
        x = checkpoint_sequential(layers,
                                  segments=2,
                                  input=x)
        return x
```

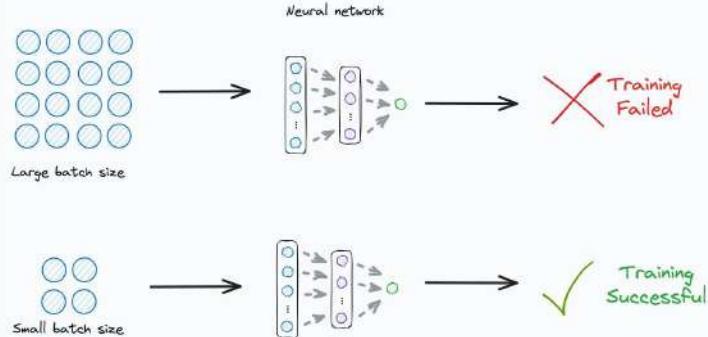
Gradient
Checkpointing

As demonstrated above, in the forward method, we use the `checkpoint_sequential` method to use gradient checkpointing and divide the network into two segments.

Next, we can proceed with network training as we usually would.

Gradient Accumulation

Under memory constraints, it is always recommended to train the neural network with a small batch size.



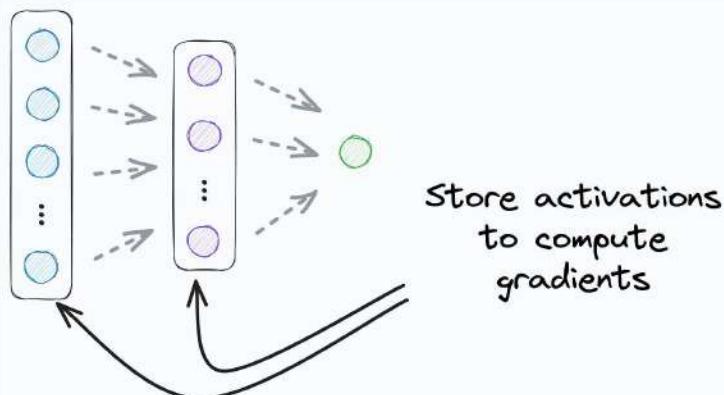
Despite that, there's a technique called **gradient accumulation**, which lets us (logically) increase batch size without explicitly increasing the batch size.

Confused?

Let's understand in this chapter. But before that, we must understand...

Why do neural networks typically explode during training?

The primary memory overhead in a neural network comes from backpropagation. This is because, during backpropagation, we must store the layer activations in memory. After all, they are used to compute the gradients.



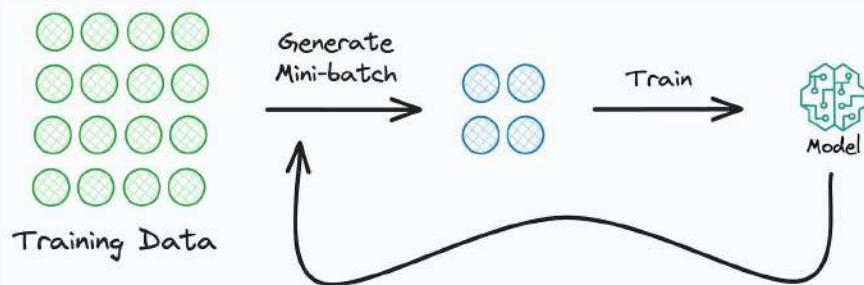
The bigger the network, the more activations a network must store in memory. Also, under memory constraints, having a large batch size will result in:

- storing many activations
- using those many activations to compute the gradients

This may lead to more resource consumption than available — resulting in training failure. But by reducing the batch size, we can limit the memory usage and train the network.

What is Gradient Accumulation and how does it help in increasing batch size in memory constraints?

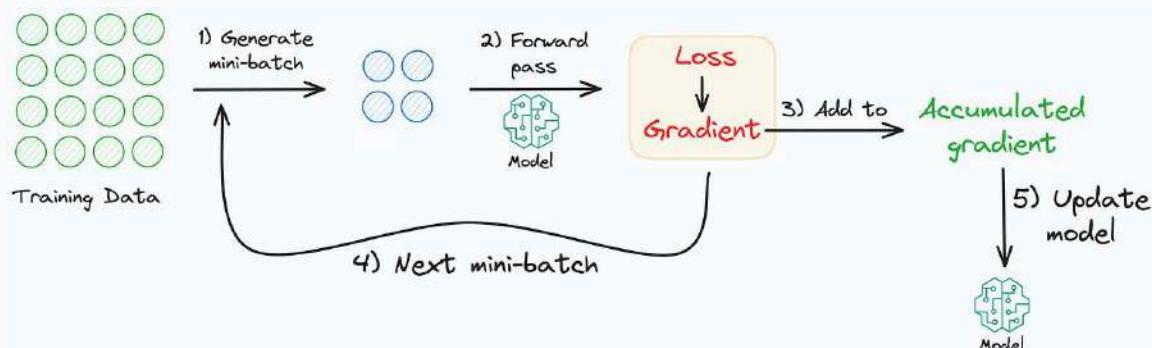
Consider we are training a neural network on mini-batches.



We train the network as follows:

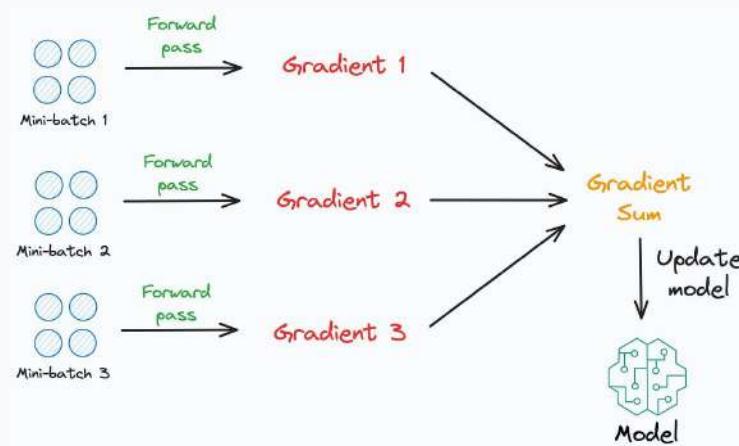
- On every mini-batch:
 - Run the forward pass while storing the activations.
 - During backward pass:
 - Compute the loss
 - Compute the gradients
 - Update the weights

Gradient accumulation modifies the last step of the backward pass, i.e., weight updates. More specifically, instead of updating the weights on every mini-batch, we can do this:



1. Run the forward pass on a mini-batch.
2. Compute the gradient values for weights in the network.
3. Don't update the weights yet.
4. Run the forward pass on the next mini-batch.
5. Compute the gradient values for weights and add them to the gradients obtained in step 2.
6. Repeat steps 3-5 for a few more mini-batches.
7. Update the weights only after processing a few mini-batches.

This technique works because accumulating the gradients across multiple mini-batches results in the same sum of gradients as if we were processing them together. Thus, logically speaking, using gradient accumulation, we can mimic a larger batch size without having to explicitly increase the batch size.



For instance, say we want to use a batch size of 64. However, current memory can only support a batch size of 16.

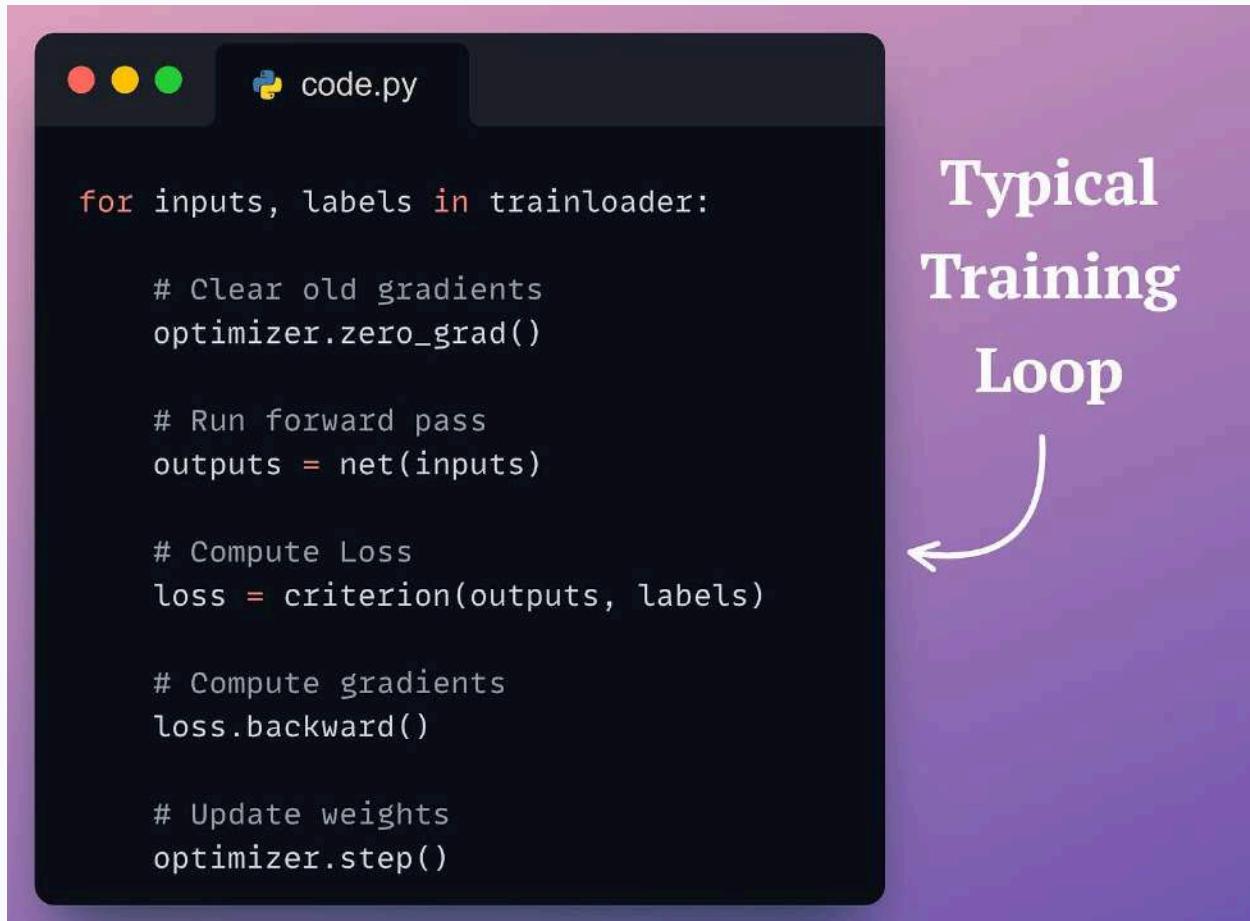
No worries!

- We can use a batch size of size 16.
- We can accumulate the gradients from every mini-batch.
- We can update the weights only once every 8 mini-batches.

Thus, effectively, we used a batch size of $16 \times 8 (=128)$ instead of what we originally intended — 64.

Implementation

Let's look at how we can implement this. In PyTorch, a typical training loop is implemented as follows:



```

for inputs, labels in trainloader:

    # Clear old gradients
    optimizer.zero_grad()

    # Run forward pass
    outputs = net(inputs)

    # Compute Loss
    loss = criterion(outputs, labels)

    # Compute gradients
    loss.backward()

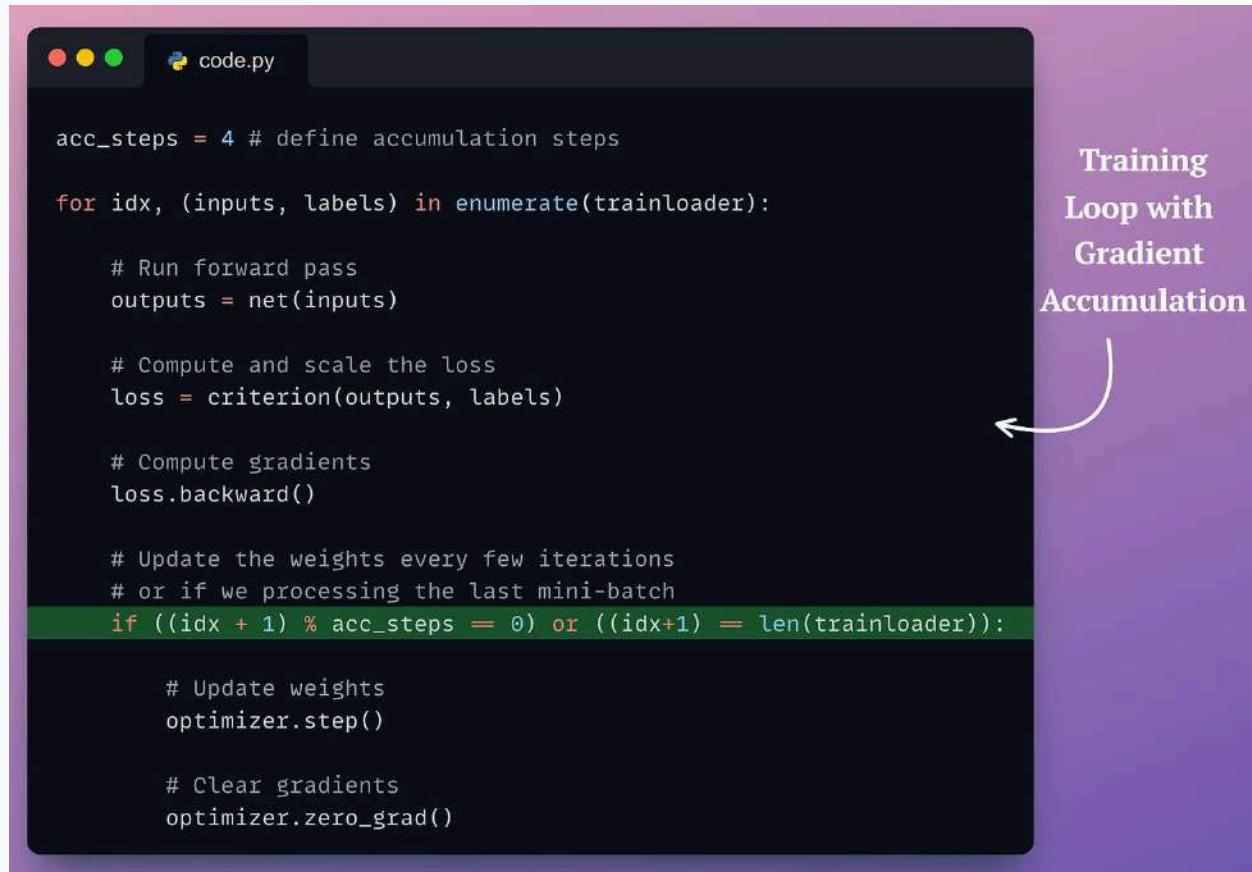
    # Update weights
    optimizer.step()

```

- We clear the gradients
- Run the forward pass
- Compute the loss
- Compute the gradients
- Update the weights

However, as discussed earlier, if needed, we can only update the weights after a few iterations. Thus, we must continue to accumulate the gradients, which is precisely what `loss.backward()` does.

Also, as `optimizer.zero_grad()` clears the gradients, we must only execute it after updating the weights. This idea is implemented below:



```

acc_steps = 4 # define accumulation steps

for idx, (inputs, labels) in enumerate(trainloader):

    # Run forward pass
    outputs = net(inputs)

    # Compute and scale the loss
    loss = criterion(outputs, labels)

    # Compute gradients
    loss.backward()

    # Update the weights every few iterations
    # or if we processing the last mini-batch
    if ((idx + 1) % acc_steps == 0) or ((idx+1) == len(trainloader)):

        # Update weights
        optimizer.step()

        # Clear gradients
        optimizer.zero_grad()

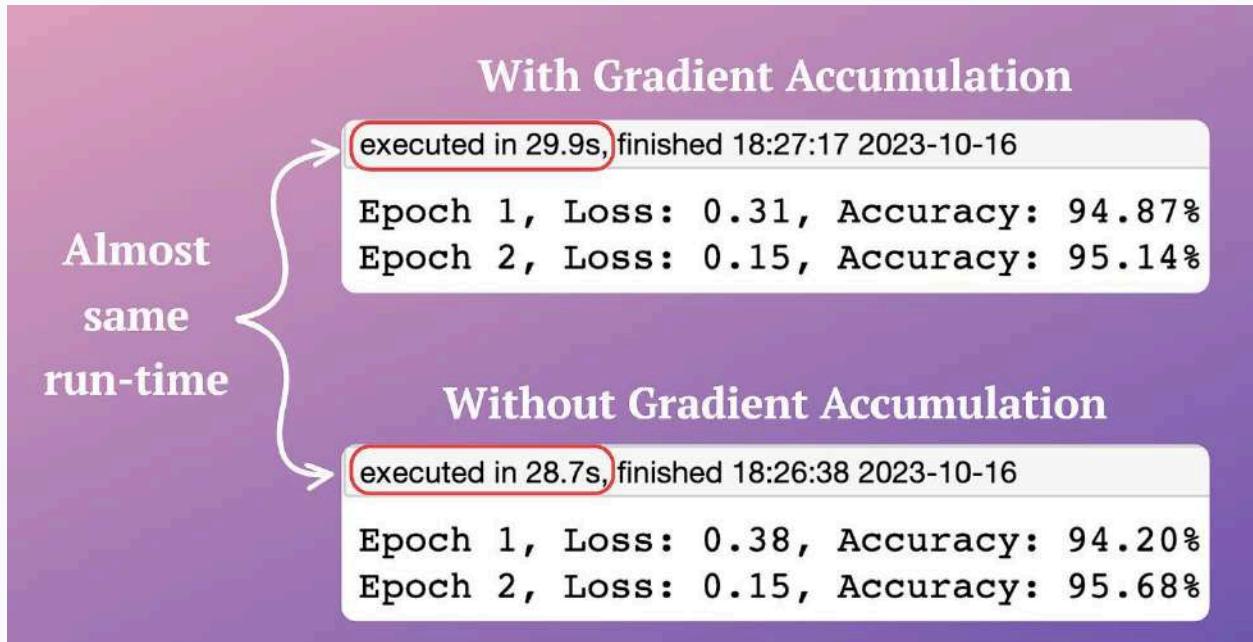
```

- First, we define `acc_steps` — the number of mini-batches after which we want to update the weights.
- Next, we run the forward pass.
- Moving on, we compute the loss and the gradients.
- As discussed earlier, we will not update the weights yet and instead let the gradients accumulate for a few more mini-batches.
- We only update the weights when the if condition is true.
- After updating, we clear the accumulated gradients.

This way, we can optimize neural network training in memory-constrained settings.

Departing note

Before we end, it is essential to note that gradient accumulation is NOT a remedy to improve run-time in memory-constrained situations. In fact, we can also verify this from my experiment:



Instead, its objective is to reduce overall memory usage.

Of course, it's true that we are updating the weights only after a few iterations. So, it will be a bit faster than updating on every iteration. Yet, we are still processing and computing gradients on small mini-batches, which is the core operation here.

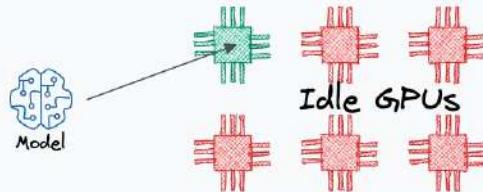
Nonetheless, the good thing is that even if you are not under memory constraints, you can still use gradient accumulation.

- Specify your typical batch size.
- Run forward pass.
- Compute loss and gradients.
- Update only after a few iterations.

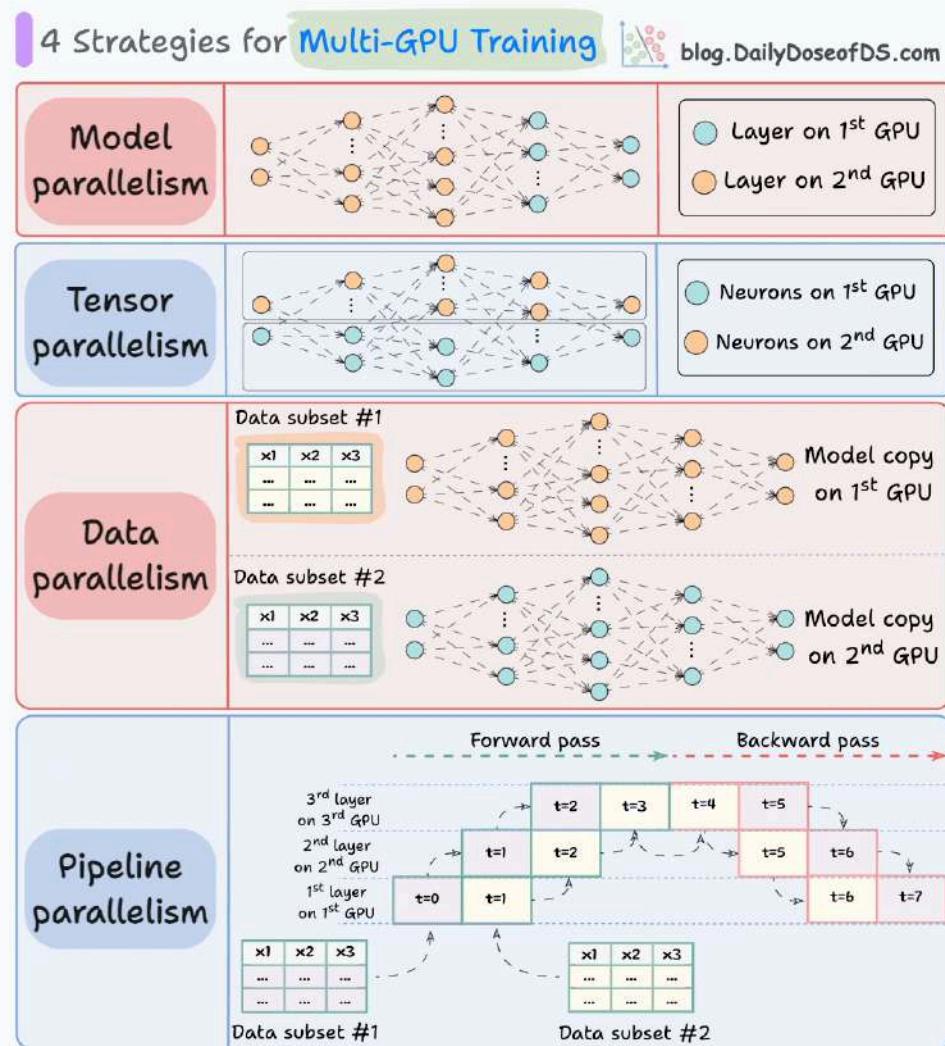
You can download the notebook here: <https://bit.ly/3xNCfFt>.

4 Strategies for Multi-GPU Training

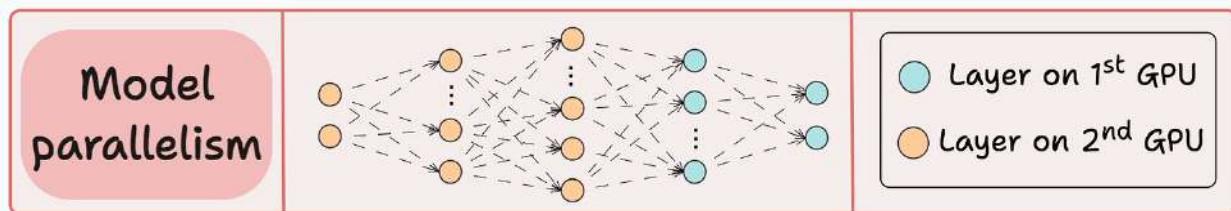
By default, deep learning models only utilize a single GPU for training, even if multiple GPUs are available.



An ideal way to proceed (especially in big-data settings) is to distribute the training workload across multiple GPUs. The graphic below depicts four common strategies for multi-GPU training:

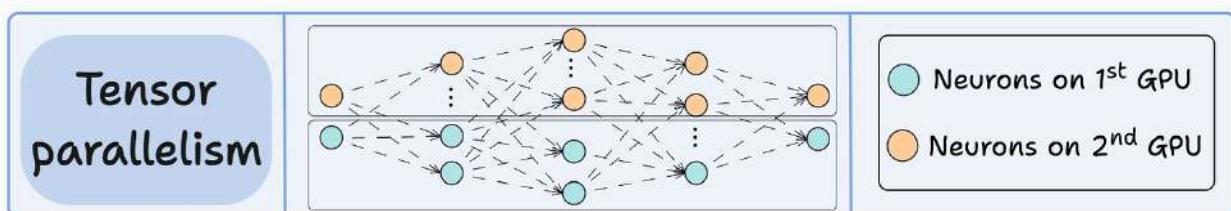


#1) Model parallelism

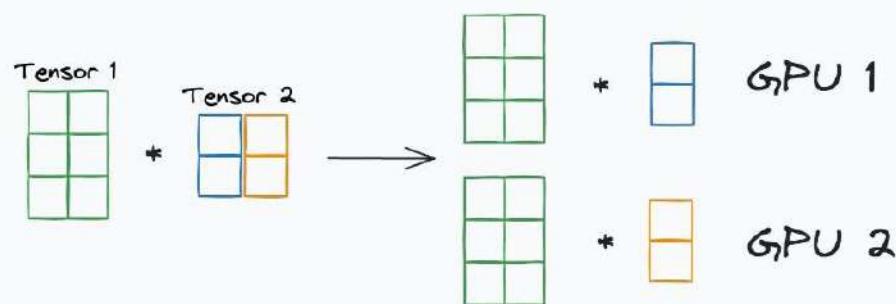


- Different parts (or layers) of the model are placed on different GPUs.
- Useful for huge models that do not fit on a single GPU.
- However, model parallelism also introduces severe bottlenecks as it requires data flow between GPUs when activations from one GPU are transferred to another GPU.

#2) Tensor parallelism

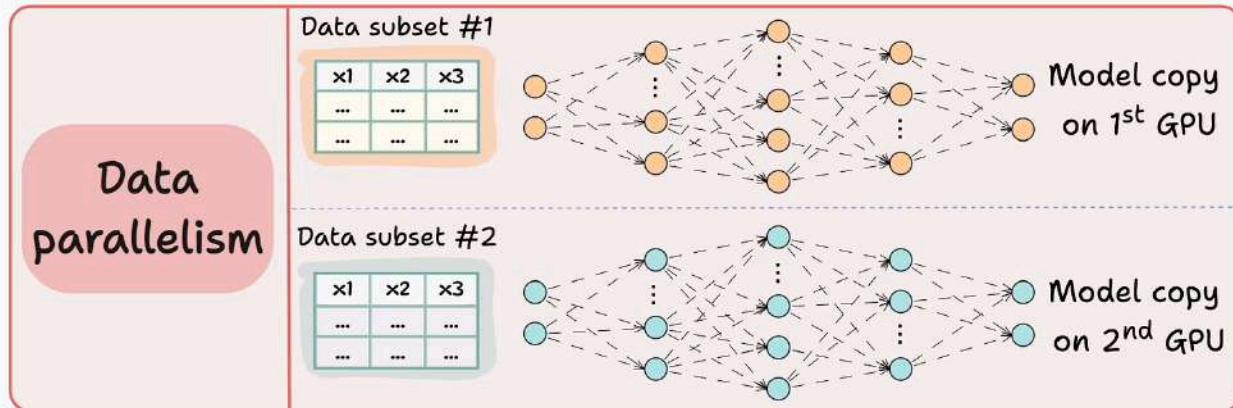


- Distributes and processes individual tensor operations across multiple devices or processors.
- It is based on the idea that a large tensor operation, such as matrix multiplication, can be divided into smaller tensor operations, and each smaller operation can be executed on a separate device or processor.



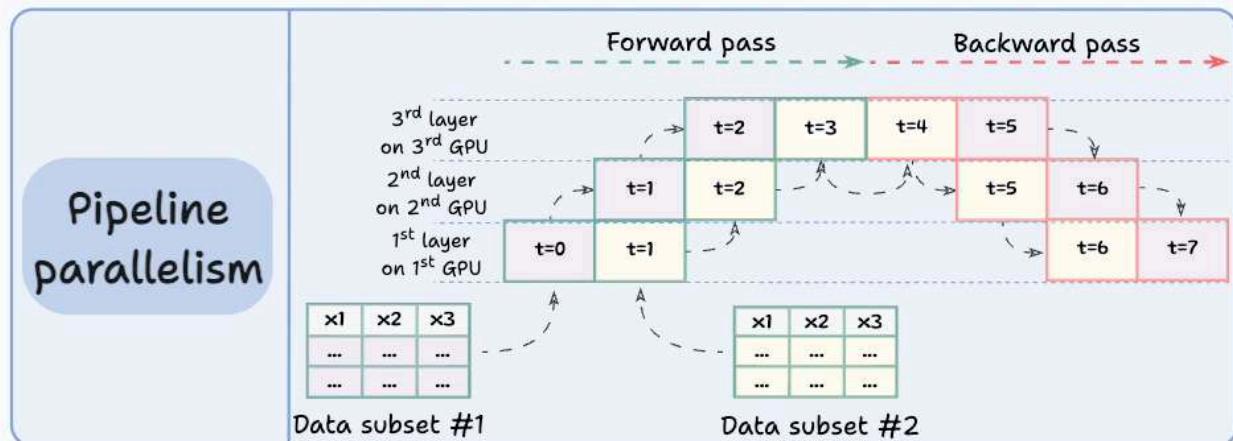
- Such parallelization strategies are inherently built into standard implementations of PyTorch and other deep learning frameworks, but they become much more pronounced in a distributed setting.

#3) Data parallelism



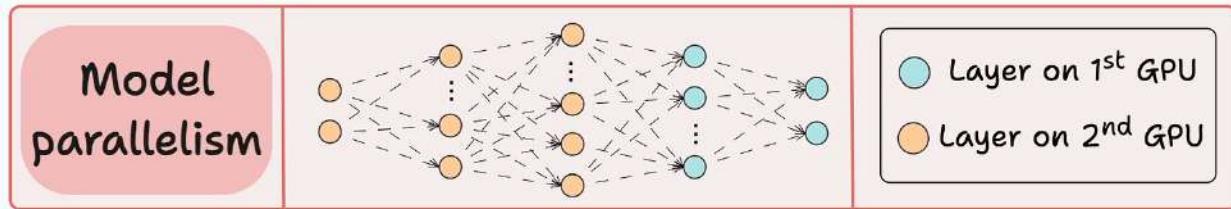
- Replicate the model across all GPUs.
- Divide the available data into smaller batches, and each batch is processed by a separate GPU.
- The updates (or gradients) from each GPU are then aggregated and used to update the model parameters on every GPU.

#4) Pipeline parallelism

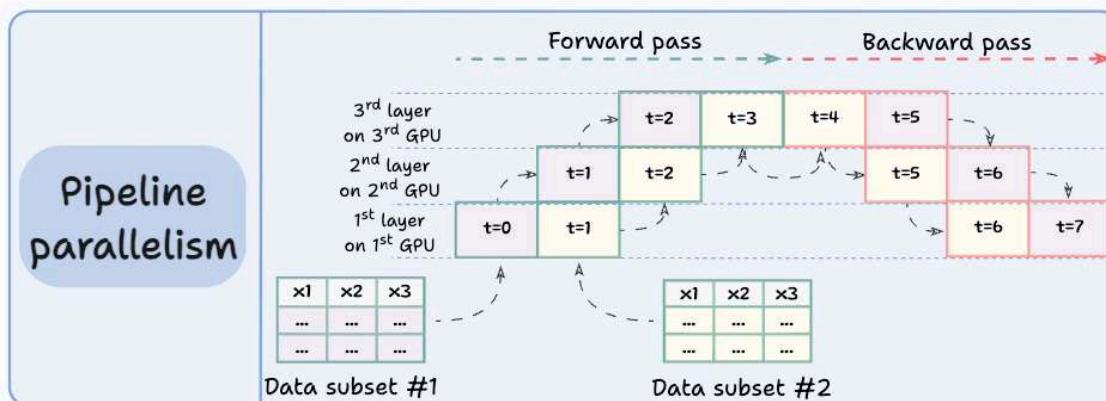


- This is often considered a combination of data parallelism and model parallelism.

- So the issue with standard model parallelism is that 1st GPU remains idle when data is being propagated through layers available in 2nd GPU:



- Pipeline parallelism addresses this by loading the next micro-batch of data once the 1st GPU has finished the computations on the 1st micro-batch and transferred activations to layers available in the 2nd GPU. The process looks like this:
 - 1st micro-batch passes through the layers on 1st GPU.
 - 2nd GPU receives activations on 1st micro-batch from 1st GPU.
 - While the 2nd GPU passes the data through the layers, another micro-batch is loaded on the 1st GPU.
 - And the process continues.
- GPU utilization drastically improves this way. This is evident from the animation below where multi-GPUs are being utilized at the same timestamp (look at t=1, t=2, t=5, and t=6):

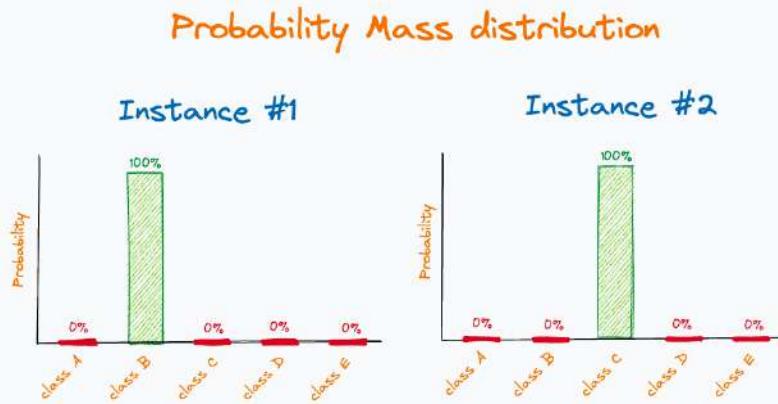


Those were four common strategies for multi-GPU training.

Miscellaneous

Label Smoothing

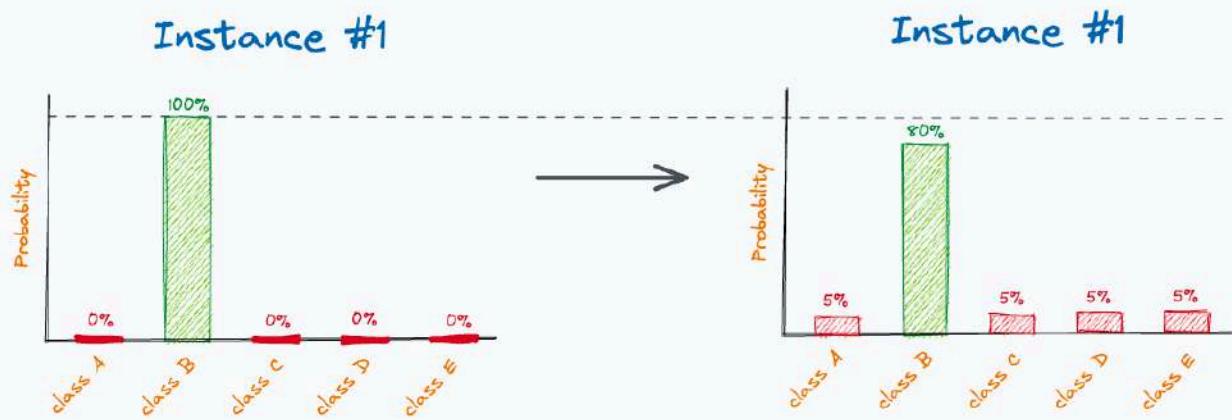
For every instance in single-label classification datasets, the entire probability mass belongs to a single class, and the rest are zero. This is depicted below:



The issue is that, at times, such label distributions excessively motivate the model to learn the true class for every sample with pretty high confidence. This can impact its generalization capabilities.

Label smoothing is a lesser-talked regularisation technique that elegantly addresses this issue.

Label Smoothing

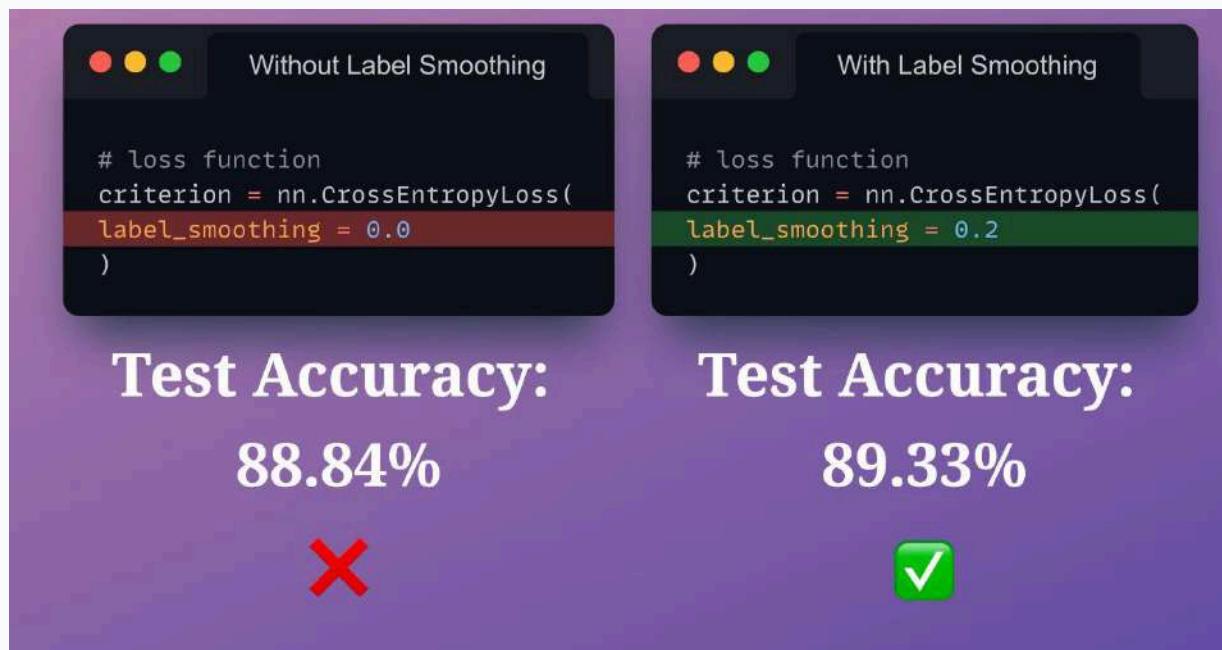


As depicted above, with label smoothing:

- We intentionally reduce the probability mass of the true class slightly.
- The reduced probability mass is uniformly distributed to all other classes.

Simply put, this can be thought of as asking the model to be “less overconfident” during training and prediction while still attempting to make accurate predictions.

The efficacy of this technique is evident from the image below:



In this experiment, I trained two neural networks on the Fashion MNIST dataset with the exact same weight initialization.

- One without label smoothing.
- Another with label smoothing.

The model with label smoothing resulted in a better test accuracy, i.e., better generalization.

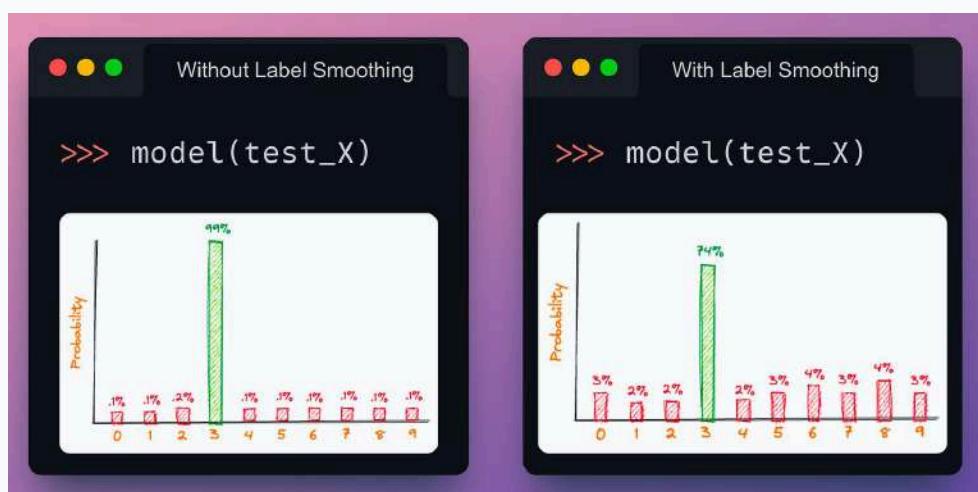
When not to use label smoothing?

After using label smoothing for many of my projects, I have also realized that it is not well suited for all use cases. So it's important to know when you should not use it.

See, if you only care about getting the final prediction correct and improving generalization, label smoothing will be a pretty handy technique. However, I wouldn't recommend utilizing it if you care about:

- Getting the prediction correct.
- **And** understanding the model's confidence in generating a prediction.

This is because as we discussed above, label smoothing guides the model to become “less overconfident” about its predictions. Thus, we typically notice a drop in the confidence values for every prediction, as depicted below:



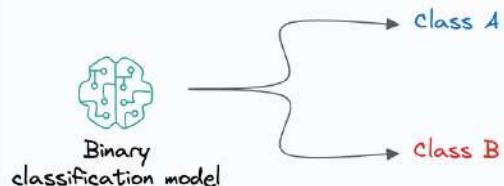
On a specific test instance:

- The model without label smoothing outputs 99% probability for class 3.
- With label smoothing, although the prediction is still correct, the confidence drops to 74%.

This is something to keep in mind when using label smoothing. Nonetheless, the technique is indeed pretty promising for regularizing deep learning models. You can download the code notebook for this chapter here: <https://bit.ly/4ePt08d>.

Focal loss

Binary classification tasks are typically trained using the binary cross entropy (BCE) loss function:



$$CE(p, y) = \begin{cases} -\log(p) & : y = 1 \\ -\log(1-p) & : y = 0 \end{cases}$$

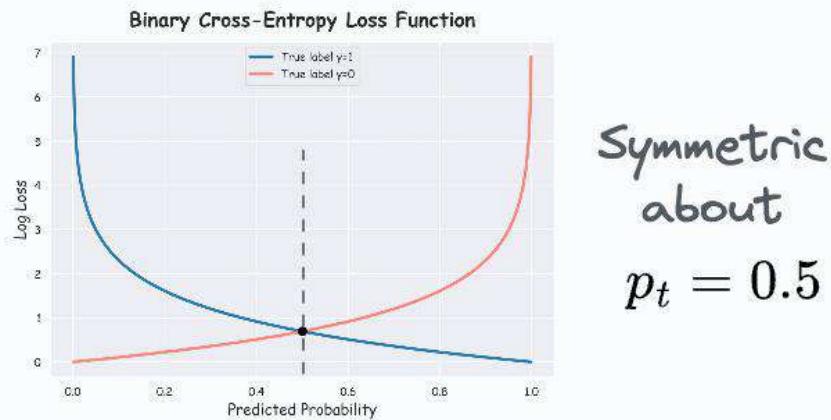
For notational convenience, if we define p_t as the following:

$$p_t = \begin{cases} p & : y = 1 \\ 1 - p & : y = 0 \end{cases}$$

...then we can also write the cross-entropy loss function as:

$$CE(p, y) = -\log(p_t)$$

That said, one limitation of BCE loss is that it weighs probability predictions for both classes equally, which is evident from its symmetry:



For more clarity, consider the table below, which depicts two instances, one from the minority class and another from the majority class, both with the same loss:

True label	$y = 1$ minority class	$y = 0$ majority class
Probability output	$\hat{y} = 0.3$	$\hat{y} = 0.7$
log loss	$-\log(0.3)$	$-\log(0.3)$

This causes problems when we use BCE for imbalanced datasets, wherein most instances from the dominating class are “easily classifiable.” Thus, a loss value of, say, $-\log(0.3)$ from the majority class instance should (ideally) be weighed LESS than the same loss value from the minority class.

$$\begin{array}{ccc} -\log(0.3) & \gg & -\log(0.3) \\ \text{from minority} & & \text{from majority} \\ \text{class} & & \text{class} \end{array}$$

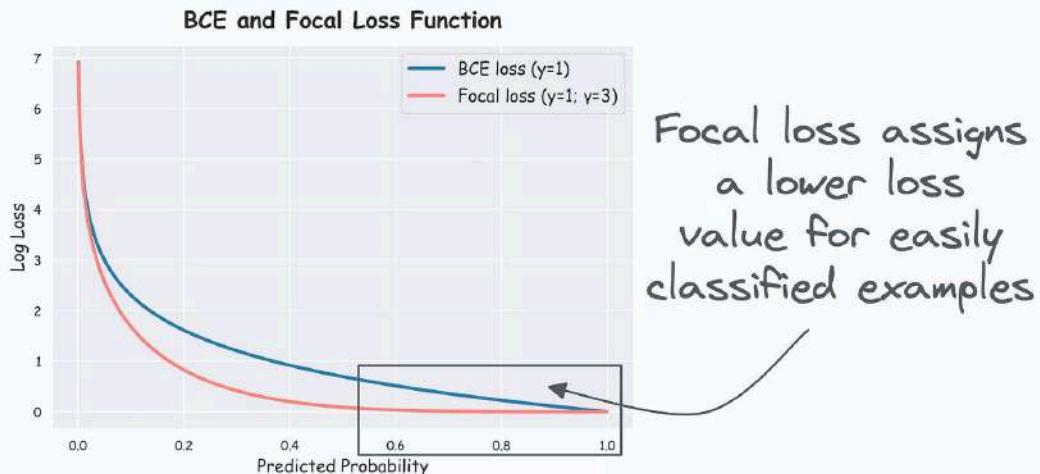
Focal loss is a pretty handy and useful alternative to address this issue. It is defined as follows:

$$CE(p, y) = -\log(p_t)$$

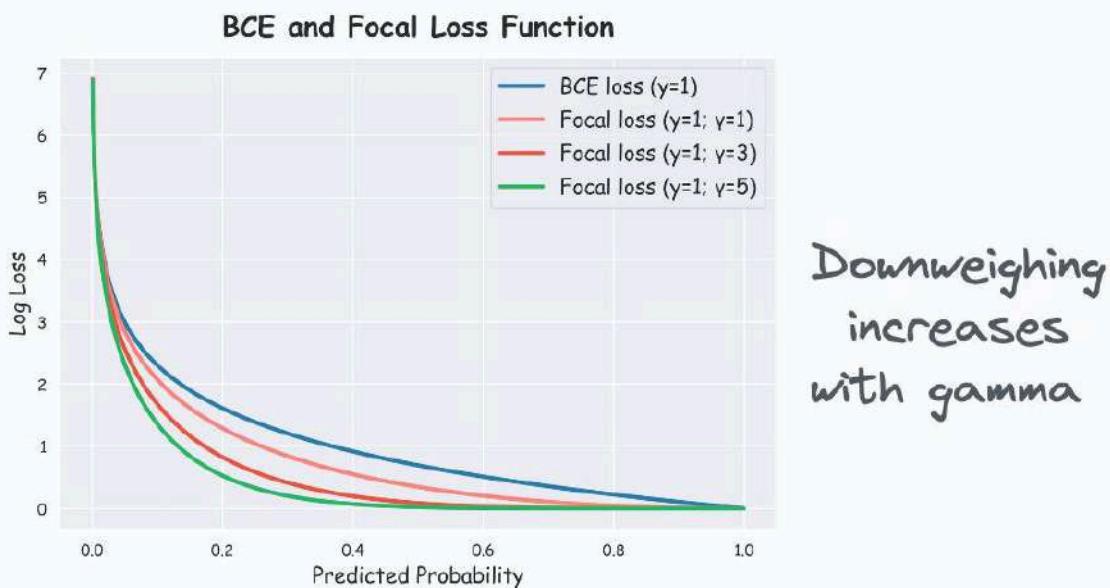
$$Focal(p, y) = -(1 - p_t)^\gamma \cdot \log(p_t)$$

As depicted above, it introduces an additional multiplicative factor called downweighing, and the parameter γ (*Gamma*) is a hyperparameter.

Plotting BCE (class $y=1$) and Focal loss (for class $y=1$ and $\gamma=3$), we get the following curve:

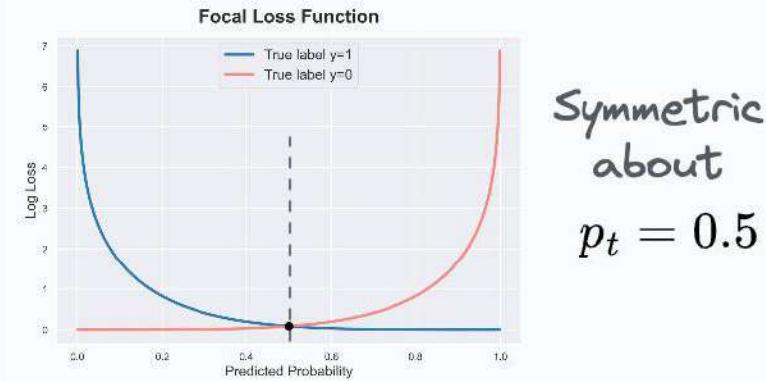


As shown in the figure above, focal loss reduces the contribution of the predictions the model is pretty confident about. Also, the higher the value of γ (Gamma), the more downweighting takes place, as shown in this plot below:

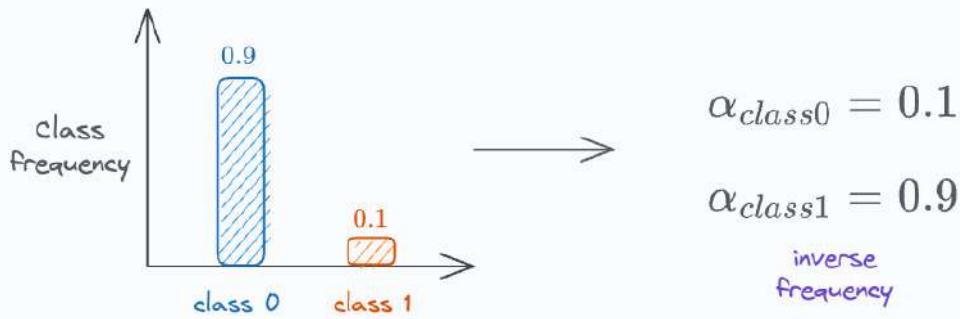


Moving on, while the Focal loss function reduces the contribution of confident predictions, we aren't done yet.

The focal loss function now is still symmetric like BCE:



To address this, we must add another weighing parameter (α), which is the inverse of the class frequency, as depicted below:



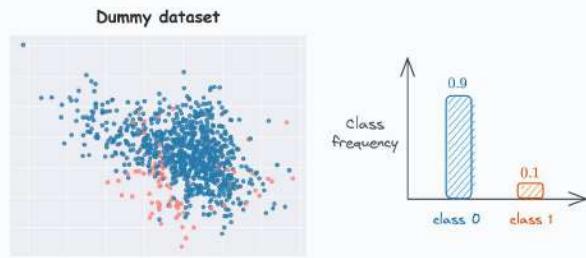
The α parameter is the inverse of the class frequency. Thus, the final loss function comes out to be the following:

$$Focal(p, y) = - \left[\alpha_t \cdot (1 - p_t)^\gamma \cdot \log(p_t) \right]$$

class inverse weighing Probability downweighting

By using both downweighting and inverse weighting, the model gradually learns patterns specific to the hard examples instead of always being overly confident in predicting easy instances.

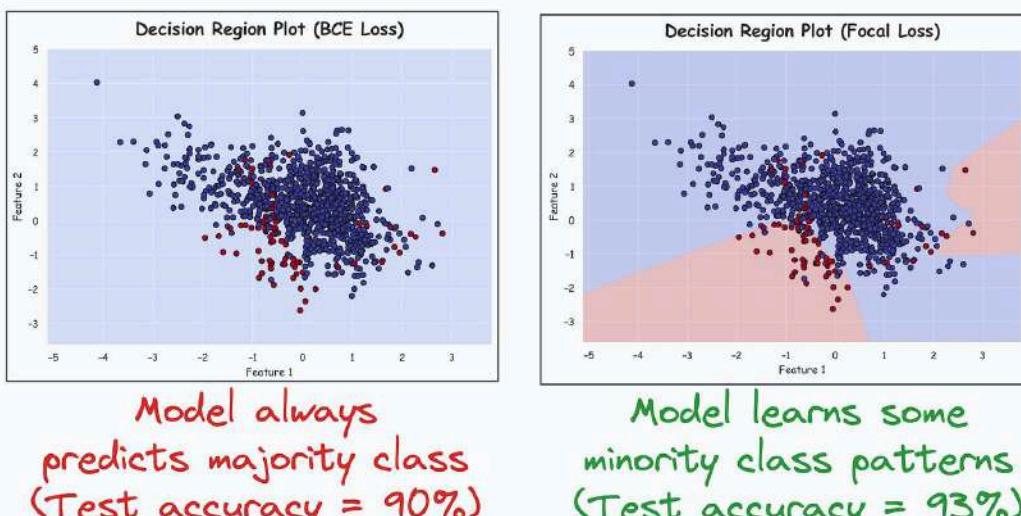
To test the efficacy of focal loss in a class imbalance setting, I created a dummy classification dataset with a 90:10 imbalance ratio:



Next, I trained two neural network models (with the same architecture of 2 hidden layers):

- One with BCE loss
- Another with Focal loss

The decision region plot and test accuracy for these two models is depicted below:



It is clear that:

- The model trained with BCE loss (left) always predicts the majority class.
- The model trained with focal loss (right) focuses relatively more on minority class patterns. As a result, it performs better.

Download this Jupyter notebook to get started with Focal loss:

<https://bit.ly/45XzNZC>.

How Dropout Actually Works?

Some time back, I was invited by a tech startup to conduct their ML interviews. I interviewed 12 candidates and mostly asked practical ML questions.

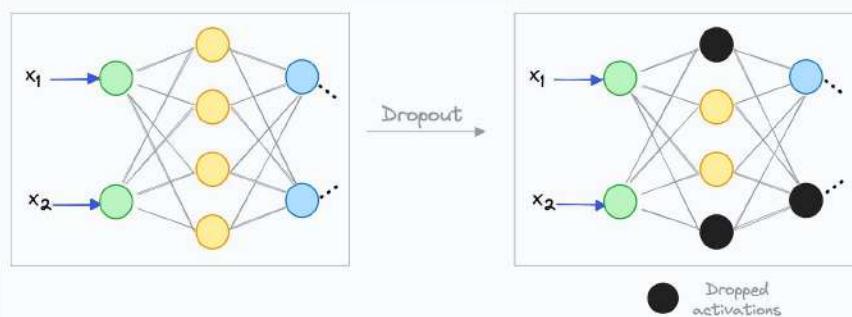
However, there were some conceptual questions as well, like the one below, which I intentionally asked every candidate:

How does Dropout work?

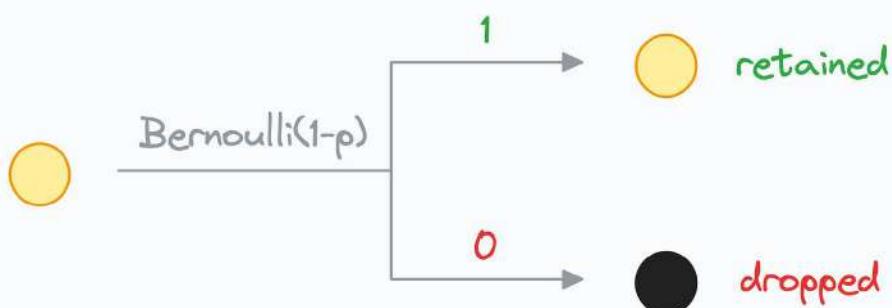
Pretty simple, right? Apparently, every candidate gave me an incomplete answer, which I have mentioned below:

Candidates' Answer

In a gist, the idea is to zero out neurons randomly in a neural network. This is done to regularize the network.



Dropout is only applied during training, and which neuron activations to zero out (or drop) is decided using a Bernoulli distribution:



“p” is the dropout probability specified in, say, PyTorch → nn.Dropout(p).

My follow-up question: Is there anything else that we do in Dropout?

Candidates: No, that's it. We only zero out neurons and train the network as we usually would.

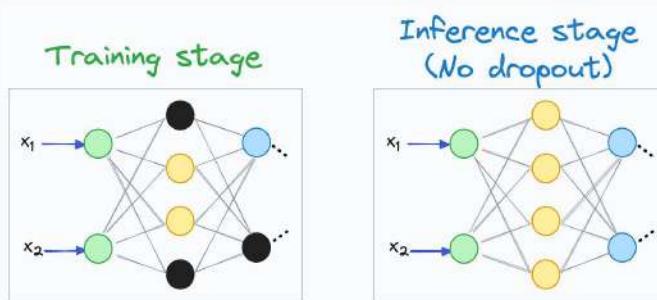
Now, coming back to the topic...

Of course, I am not saying that the above details are incorrect. They are correct.

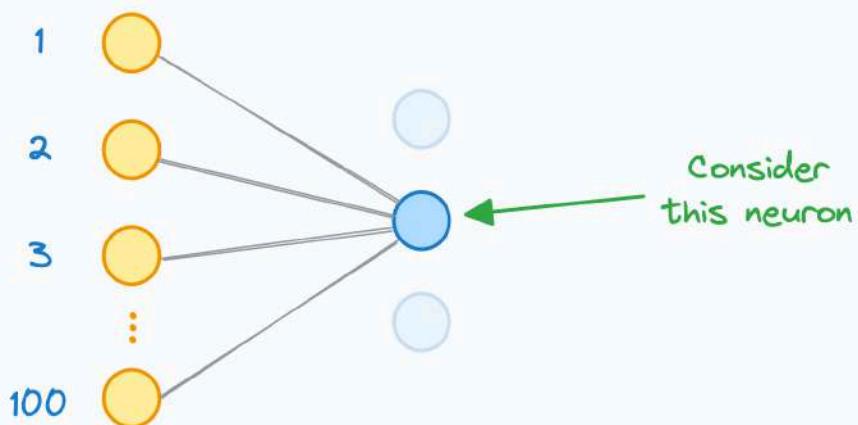
However, this is just 50% of how Dropout works, and disappointingly, most resources don't cover the remaining 50%. If you too are only aware of the 50% details I mentioned above, continue reading.

How Dropout actually works?

To begin, we must note that Dropout is only applied during training, but not during the inference/evaluation stage:

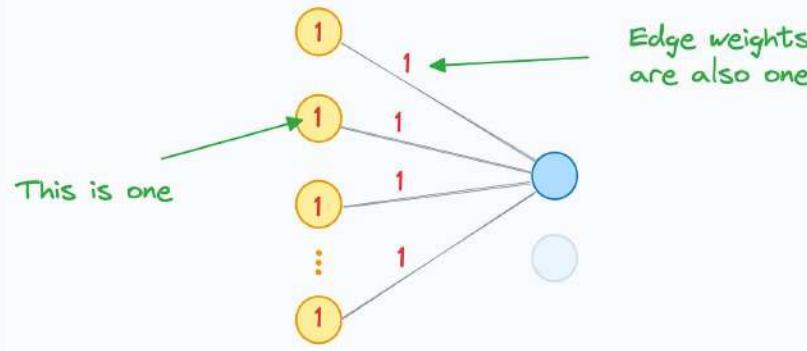


Now, consider that a neuron's input is computed using 100 neurons in the previous hidden layer:

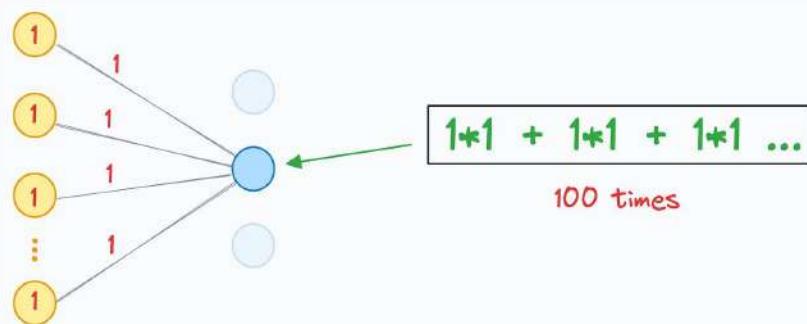


For simplicity, let's assume a couple of things here:

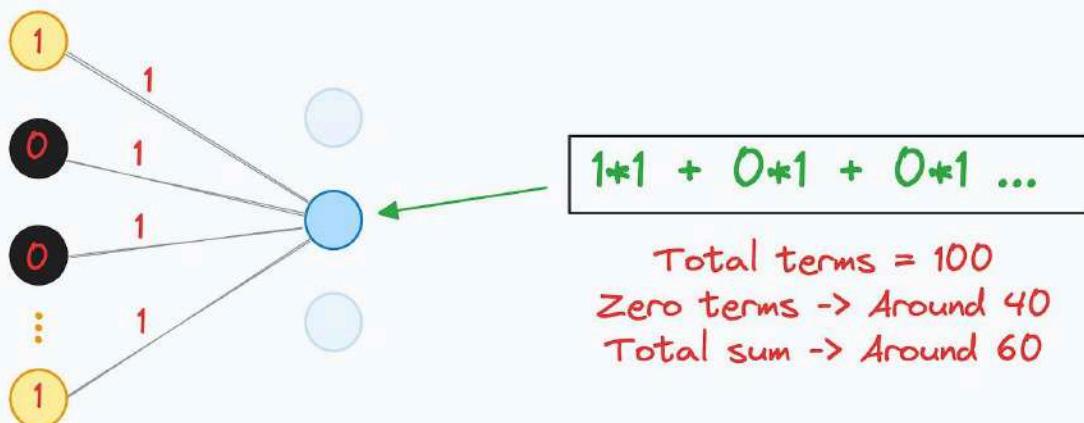
- The activation of every yellow neuron is 1.
- The edge weight from the yellow neurons to the blue neuron is also 1.



As a result, the input received by the blue neuron will be 100, as depicted below:



Now, during training, if we were using Dropout with, say, a 40% dropout rate, then roughly 40% of the yellow neuron activations would have been zeroed out. As a result, the input received by the blue neuron would have been around 60:



However, the above point is only valid for the training stage.

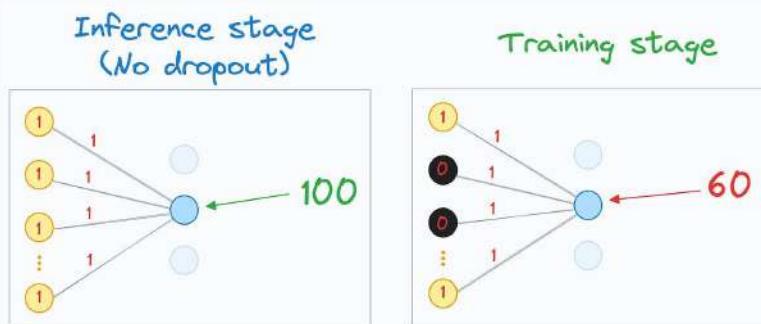
If the same scenario had existed during the inference stage instead, then the input received by the blue neuron would have been 100.

Thus, under similar conditions:

- The input received during training → 60.
- The input received during inference → 100.

Do you see any problem here?

During training, the average neuron inputs are significantly lower than those received during inference.



More formally, using Dropout significantly affects the scale of the activations. However, it is desired that the neurons throughout the model must receive the roughly same mean (or expected value) of activations during training and inference. To address this, Dropout performs one additional step.

This idea is to scale the remaining active inputs during training. The simplest way to do this is by scaling all activations during training by a factor of $1/(1-p)$, where p is the dropout rate. For instance, using this technique on the neuron input of 60, we get the following (recall that we set $p=40\%$):

$$60 \rightarrow 60 * \frac{1}{1 - 0.4} = \boxed{100}$$

← output rescaled

As depicted above, scaling the neuron input brings it to the desired range, which makes training and inference stages coherent for the network.

Verifying experimentally

In fact, we can verify that typical implementations of Dropout, from PyTorch, for instance, do carry out this step. Let's define a dropout layer as follows:

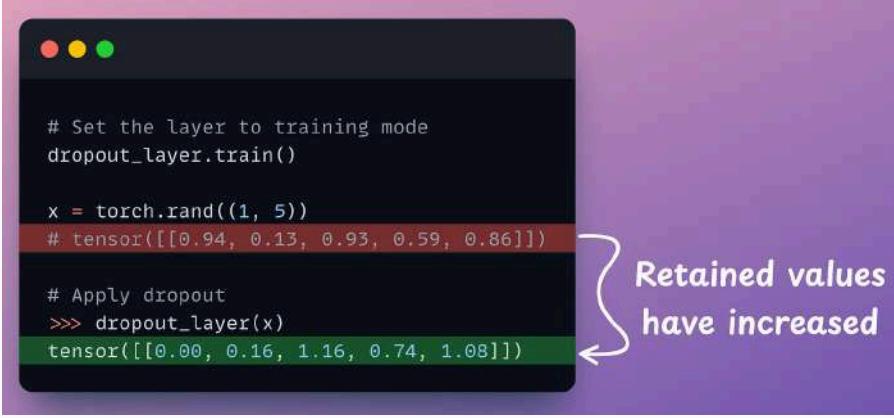


```
import torch
import torch.nn as nn

dropout_layer = nn.Dropout(0.2)
```

Dropout layer with $p=0.2$

Now, let's consider a random tensor and apply this dropout layer to it:



```
# Set the layer to training mode
dropout_layer.train()

x = torch.rand((1, 5))
# tensor([[0.94, 0.13, 0.93, 0.59, 0.86]])

# Apply dropout
>>> dropout_layer(x)
tensor([[0.00, 0.16, 1.16, 0.74, 1.08]])
```

Retained values have increased

As depicted above, the retained values have increased.

- The second value goes from $0.13 \rightarrow 0.16$.
- The third value goes from $0.93 \rightarrow 1.16$.
- and so on...

What's more, the retained values are precisely the same as we would have obtained by explicitly scaling the input tensor with $1/(1-p)$:

```

>>> dropout_layer(x)
tensor([0.00, 0.16, 1.16, 0.74, 1.08])

>>> x/(1-p)
tensor([0.94, 0.16, 1.16, 0.74, 1.08])

```

Same values

If we were to do the same thing in evaluation mode instead, we notice that no value is dropped and no scaling takes place either, which makes sense as Dropout is only used during training:

```

# Set the layer to evaluation mode
dropout_layer.eval()

>>> x
tensor([0.94, 0.13, 0.93, 0.59, 0.86])

# Apply dropout during evaluation
>>> dropout_layer(x)
tensor([0.94, 0.13, 0.93, 0.59, 0.86])

```

No change

This is the remaining 50% details, which, in my experience, most resources do not cover, and as a result, most people aren't aware of.

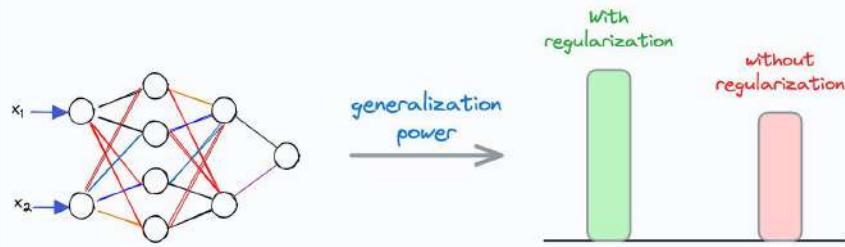
But it is a highly important step in Dropout, which maintains numerical coherence between training and inference stages.

With that, now you know 100% of how Dropout works.

Next, let's discuss an issue with Dropout in case of CNNs.

Issues with Dropout in CNNs

When it comes to training neural networks, it is always recommended to use Dropout to improve its generalization power.

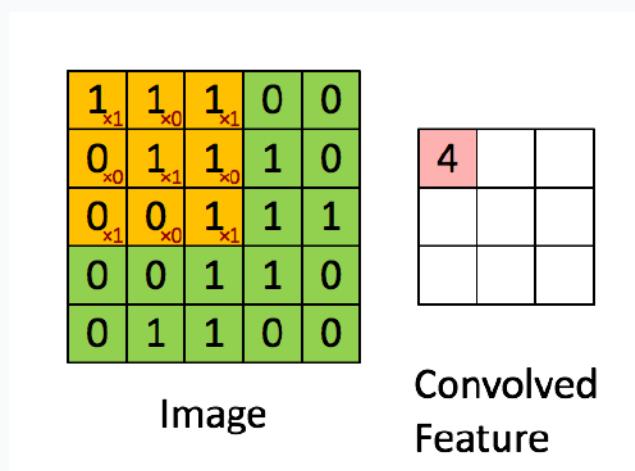


This applies not just to CNNs but to all other neural networks. And I am sure you already know the above details, so let's get into the interesting part.

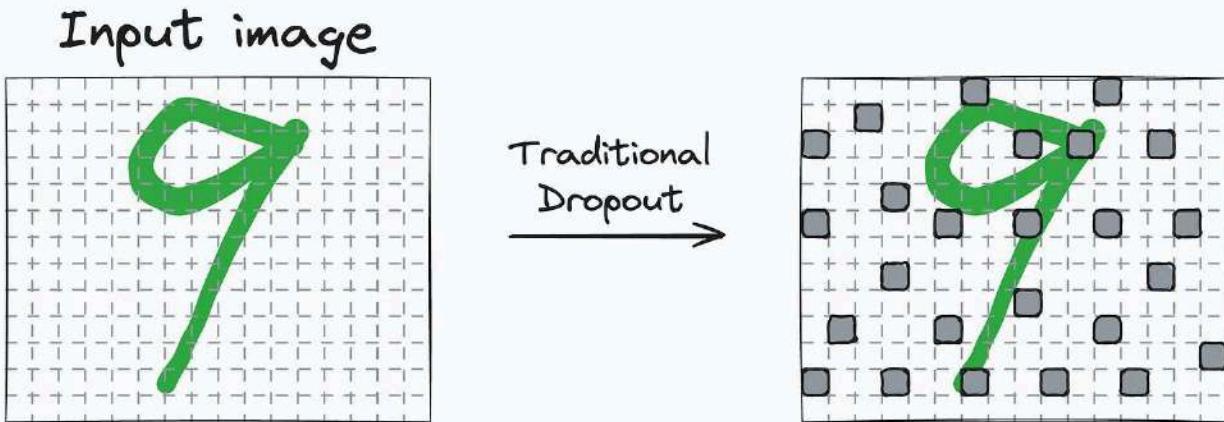
The problem of using Dropout in CNNs

The core operation that makes CNNs so powerful is convolution, which allows them to capture local patterns, such as edges and textures, and helps extract relevant information from the input.

From a purely mathematical perspective, we slide a filter (shown in yellow below) over the input (shown in green below) and take the element-wise sum between the filter and the overlapped input to get the convolution output:



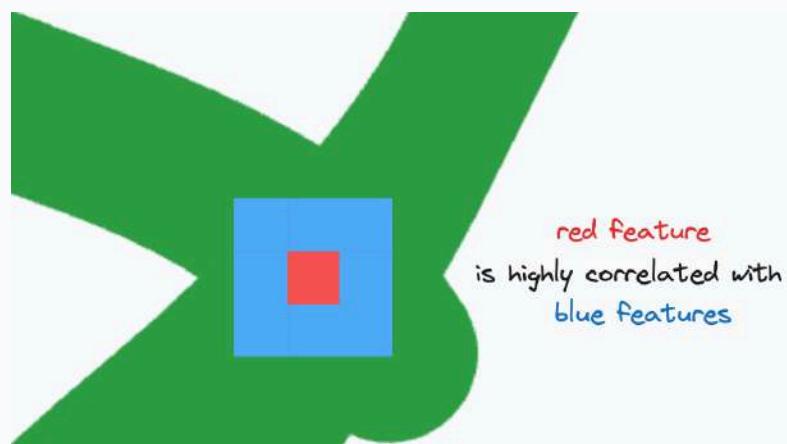
Here, if were to apply the traditional Dropout, the input features would look something like this:



In fully connected layers, we zero out neurons. In CNNs, however, we randomly zero out the pixel values before convolution, as depicted above.

But this isn't found to be that effective specifically for convolution layers. To understand this, consider we have some image data. In every image, we would find that nearby features (or pixels) are highly correlated spatially.

For instance, imagine zooming in on the pixel level of the digit '9'. Here, we would notice that the red pixel (or feature) is highly correlated with other features in its vicinity:

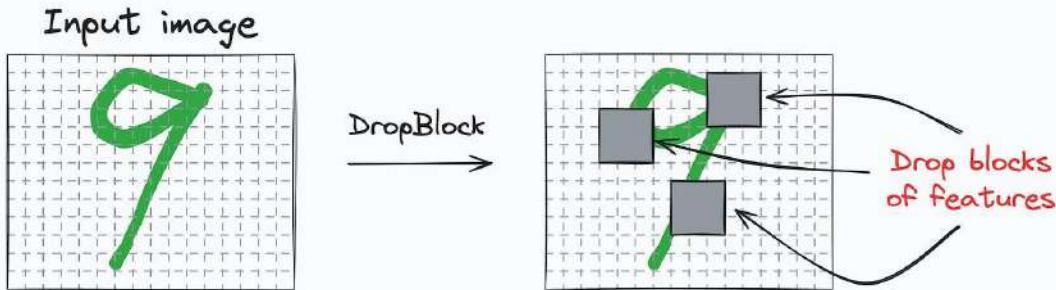


Thus, dropping the red feature using Dropout will likely have no effect and its information can still be sent to the next layer.

Simply put, the nature of the convolution operation defeats the entire purpose of the traditional Dropout procedure.

The solution

DropBlock is a much better, effective, and intuitive way to regularize CNNs. The core idea in DropBlock is to drop a contiguous region of features (or pixels) rather than individual pixels. This is depicted below:

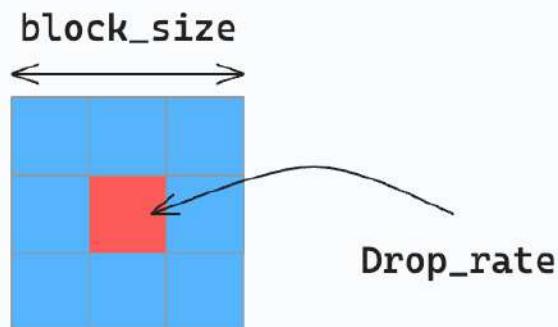


Similar to Dropout in fully connected layers, wherein the network tries to generate more robust ways to fit the data in the absence of some activations, in the case of DropBlock, the convolution layers get more robust to fit the data despite the absence of a block of features.

Moreover, the idea of DropBlock also makes intuitive sense — if a contiguous region of a feature is dropped, the problem of using Dropout with convolution operation can be avoided.

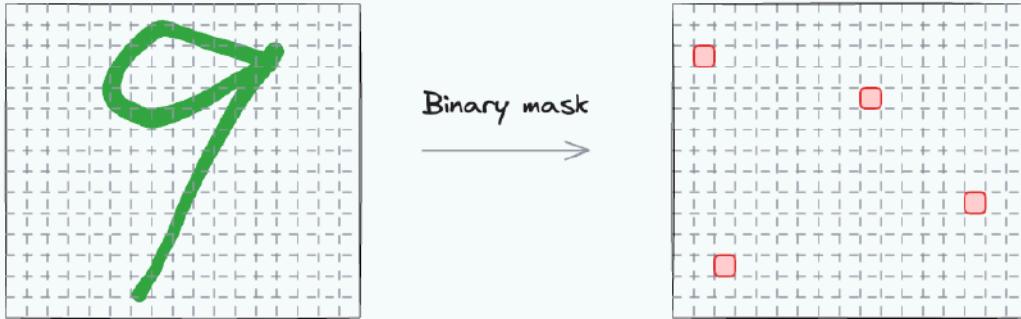
DropBlock parameters

DropBlock has two main parameters:

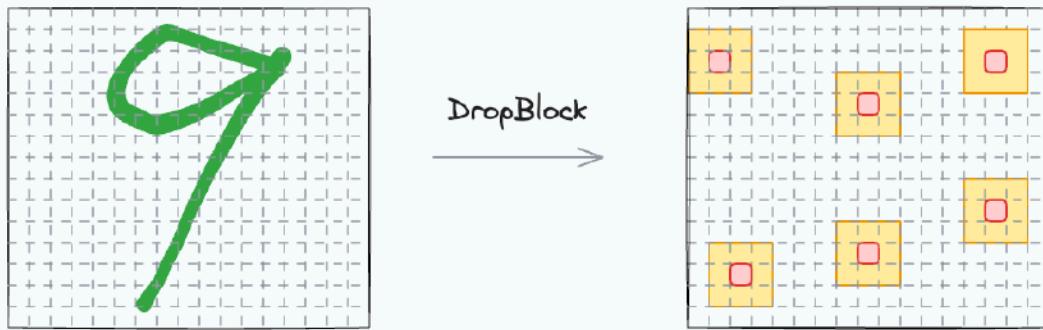


- **Block_size:** The size of the box to be dropped.
- **Drop_rate:** The drop probability of the central pixel.

To apply DropBlock, first, we create a binary mask on the input sampled from the Bernoulli distribution:



Next, we create a block of size `block_size*block_size` which has the sampled pixels at the center:



The efficacy of DropBlock over Dropout is evident from the results table below:

Model	top-1(%)	top-5(%)
ResNet-50	76.51 ± 0.07	93.20 ± 0.05
ResNet-50 + dropout ($kp=0.7$) [1]	76.80 ± 0.04	93.41 ± 0.04
ResNet-50 + DropPath ($kp=0.9$) [17]	77.10 ± 0.08	93.50 ± 0.05
ResNet-50 + SpatialDropout ($kp=0.9$) [20]	77.41 ± 0.04	93.74 ± 0.02
ResNet-50 + Cutout [23]	76.52 ± 0.07	93.21 ± 0.04
ResNet-50 + AutoAugment [27]	77.63	93.82
ResNet-50 + label smoothing (0.1) [28]	77.17 ± 0.05	93.45 ± 0.03
ResNet-50 + DropBlock, ($kp=0.9$)	78.13 ± 0.05	94.02 ± 0.02
ResNet-50 + DropBlock ($kp=0.9$) + label smoothing (0.1)	78.35 ± 0.05	94.15 ± 0.03

Annotations: A yellow arrow points from the 'ResNet-50 + DropBlock, ($kp=0.9$)' row to the '1.33% gain' label. A green arrow points from the 'ResNet-50 + DropBlock ($kp=0.9$) + label smoothing (0.1)' row to the '1.55% gain' label.

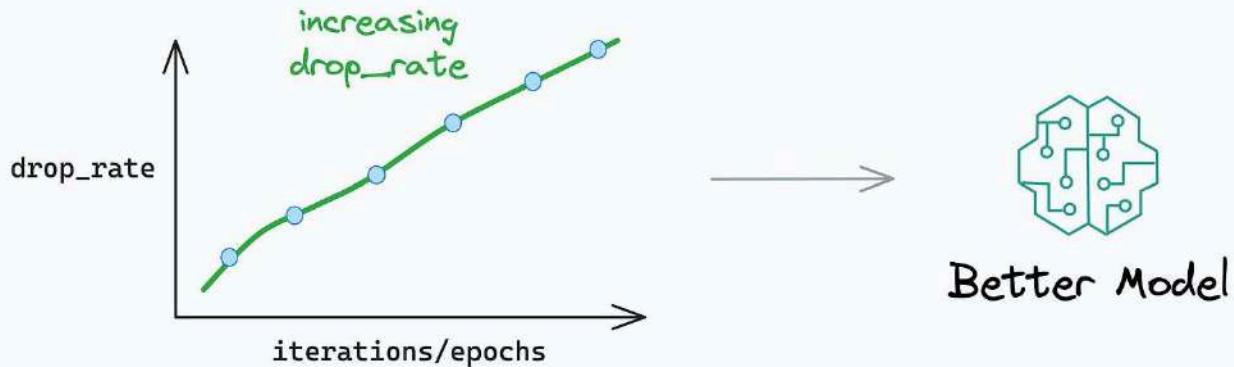
On the ImageNet classification dataset:

- DropBlock provides a 1.33% gain over Dropout.
- DropBlock with Label smoothing (discussed in the last chapter) provides a 1.55% gain over Dropout.

Thankfully, DropBlock is also integrated with PyTorch.

There's also a library for DropBlock, called “`dropblock`,” which also provides the linear scheduler for `drop_rate`.

So the thing is that the researchers who proposed DropBlock found the technique to be more effective when the `drop_rate` was increased gradually.



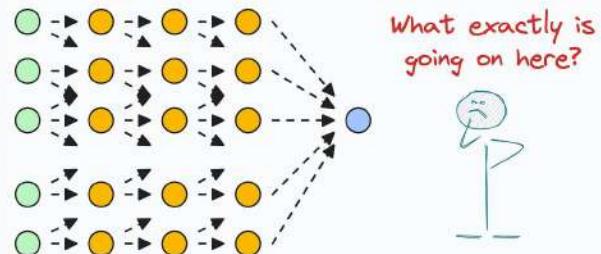
The DropBlock library implements the scheduler. But of course, there are ways to do this in PyTorch as well. So it's entirely up to you which implementation you want to use:

- DropBlock PyTorch: <https://bit.ly/3xZfT3E>.
- DropBlock library: <https://github.com/miguelvr/dropblock>.

What Hidden Layers and Activation Functions Actually Do

Everyone knows the objective of an activation function in a neural network. They let the network learn non-linear patterns. There is nothing new here, and I am sure you are aware of that too.

However, one thing I have often realized is that most people struggle to build an intuitive understanding of what exactly a neural network consistently tries to achieve during its layer-after-layer transformations.

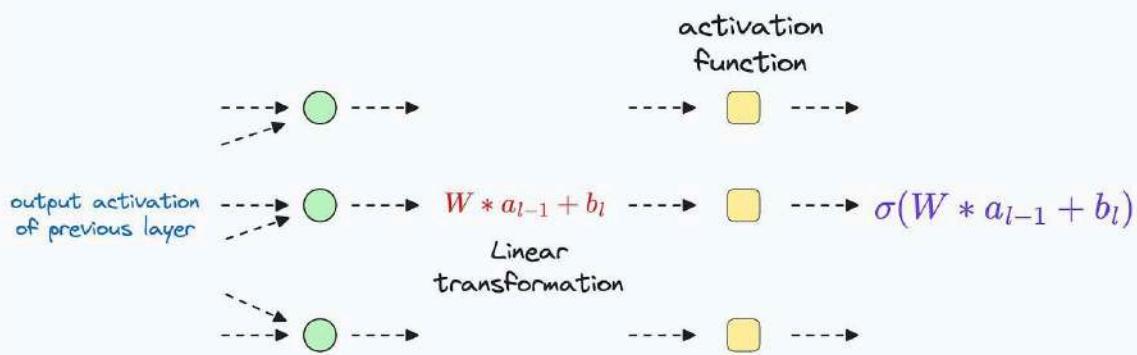


In this chapter, let me share a unique perspective on this, which would really help you understand the internal workings of a neural network.

I have supported this chapter with plenty of visuals for better understanding. Also, for simplicity, we shall consider a binary classification use case.

Background

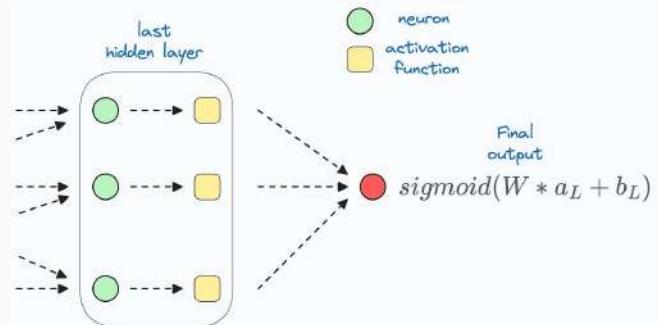
The data undergoes a series of transformations at each hidden layer:



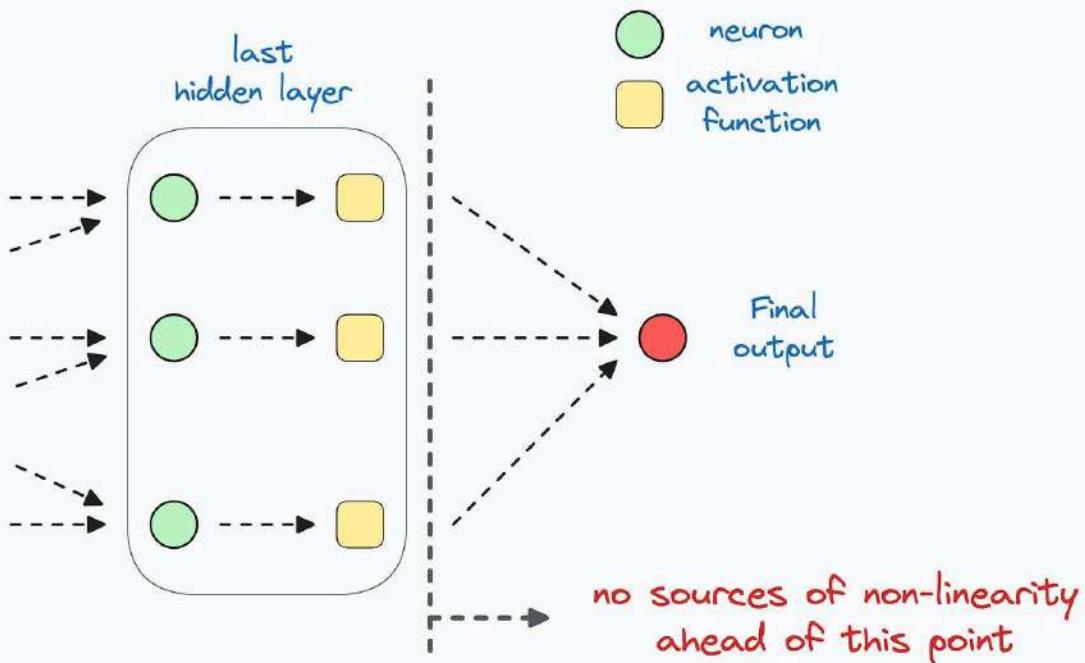
- Linear transformation of the data obtained from the previous layer
- ...followed by a non-linearity using an activation function — ReLU, Sigmoid, Tanh, etc.

The above transformations are performed on every hidden layer of a neural network. Now, notice something here.

Assume that we just applied the above data transformation on the **very last hidden layer** of the neural network. Once we do that, the activations progress toward the output layer of the network for one final transformation, which is entirely linear.



The above transformation is entirely linear because all sources of non-linearity (activations functions) exist on or before the last hidden layer. And during the forward pass, once the data leaves the last hidden layer, there is no further scope for non-linearity.



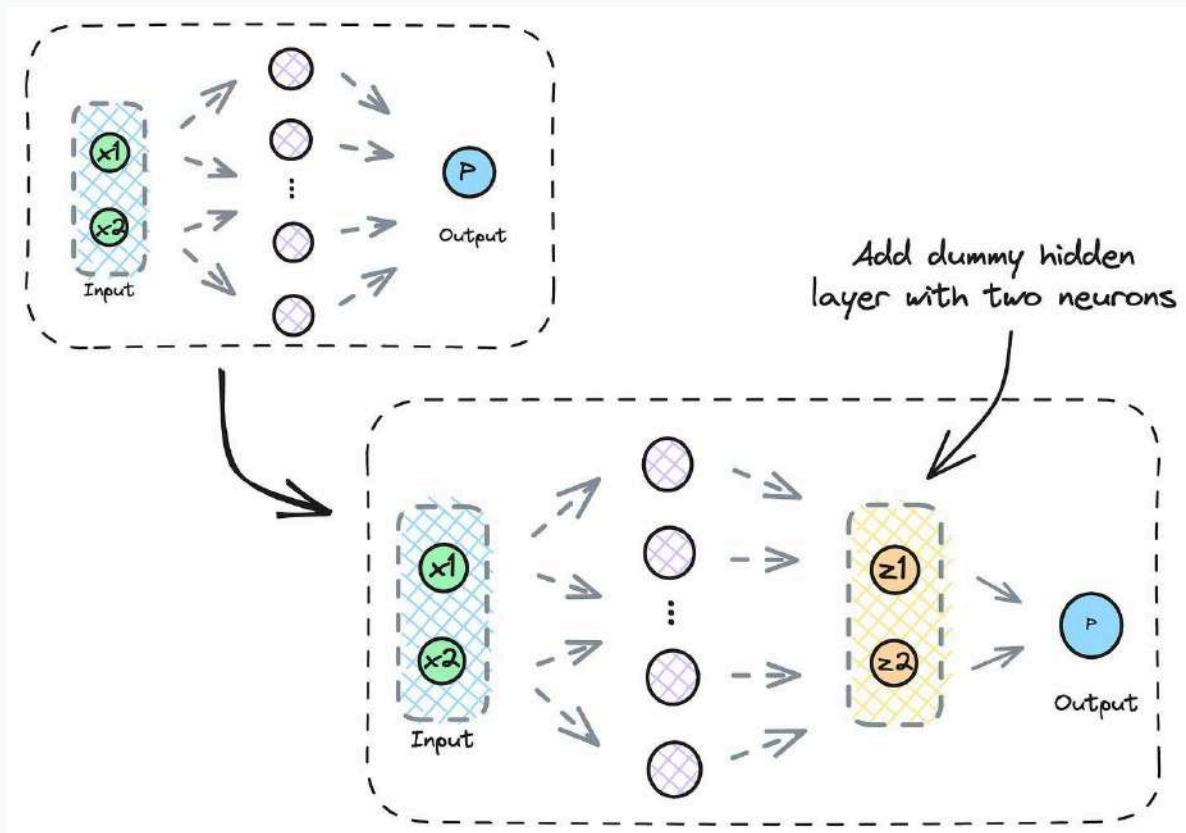
Thus, to make accurate predictions, the data received by the output layer from the last hidden layer **MUST BE** linearly separable.

To summarize....

While transforming the data through all its hidden layers and just before reaching the output layer, a neural network is constantly hustling to project the data to a space where it somehow becomes linearly separable. If it does, the output layer becomes analogous to a logistic regression model, which can easily handle this linearly separable data.

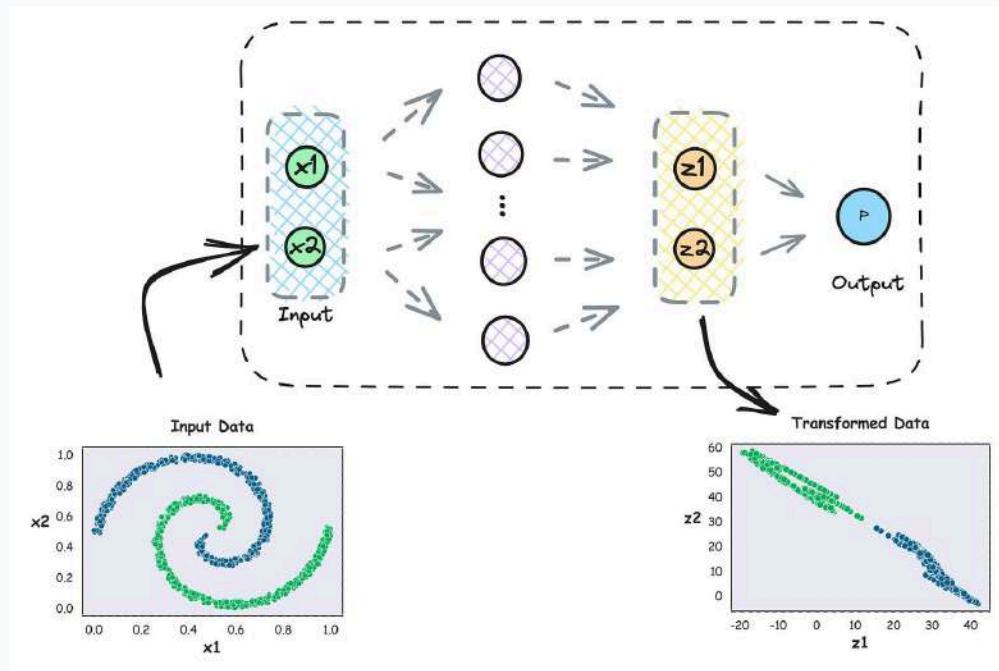
In fact, we can also verify this experimentally.

To visualize the input transformation, we can add a dummy hidden layer with just two neurons right before the output layer and train the neural network again.



Why do we add a layer with just two neurons?

This way, we can easily visualize the transformation. We expect that if we plot the activations of this 2D dummy hidden layer, they must be linearly separable. The below visual precisely depicts this.



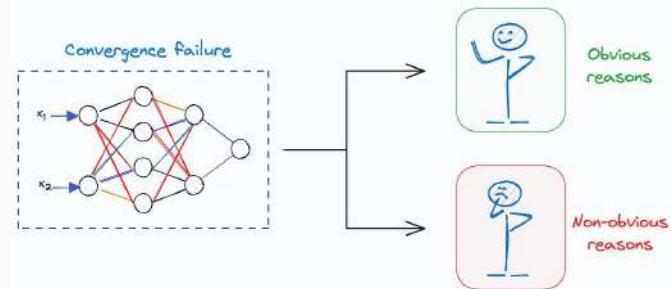
As we notice above, while the input data was linearly inseparable, the input received by the output layer is indeed linearly separable.

This transformed data can be easily handled by the output classification layer.

And this shows that all a neural network is trying to do is transform the data into a linearly separable form before reaching the output layer.

Shuffle Data Before Training

Deep learning models may fail to converge due to various reasons. Some causes are obvious and common, and therefore, quickly rectifiable, like too high/low learning rate, no data normalization, no batch normalization, etc.

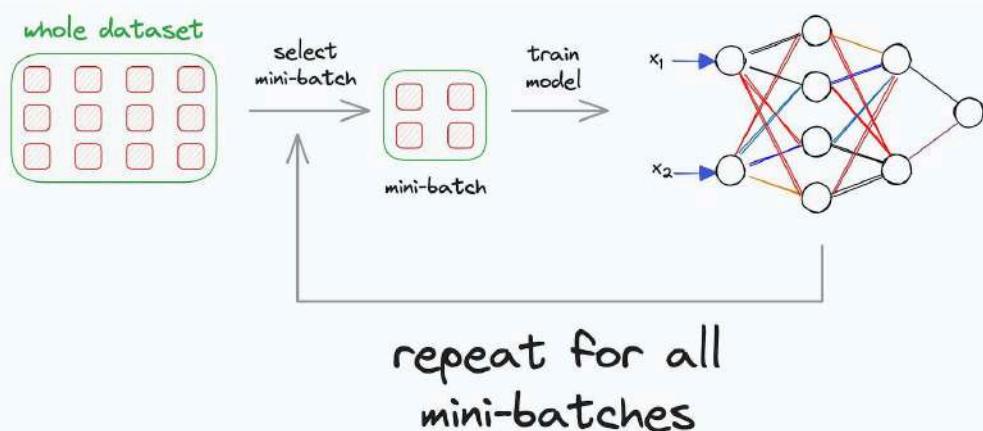


But the problem arises when the cause isn't that apparent. Therefore, it may take some serious time to debug if you are unaware of them. In this chapter, I want to talk about one such data-related mistake, which I once committed during my early days in machine learning. Admittedly, it took me quite some time to figure it out back then because I had no idea about the issue.

Experiment

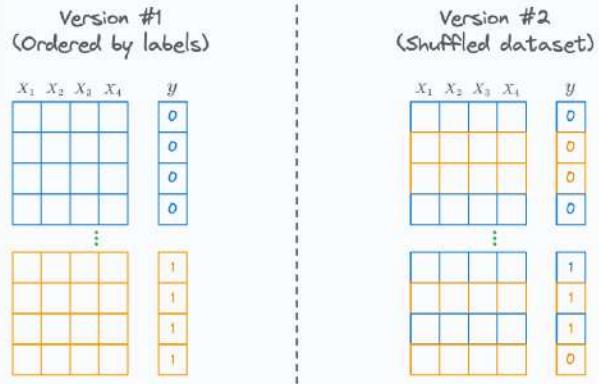
Consider a classification neural network trained using mini-batch gradient descent.

Mini-batch gradient descent: Update network weights using a few data points at a time.



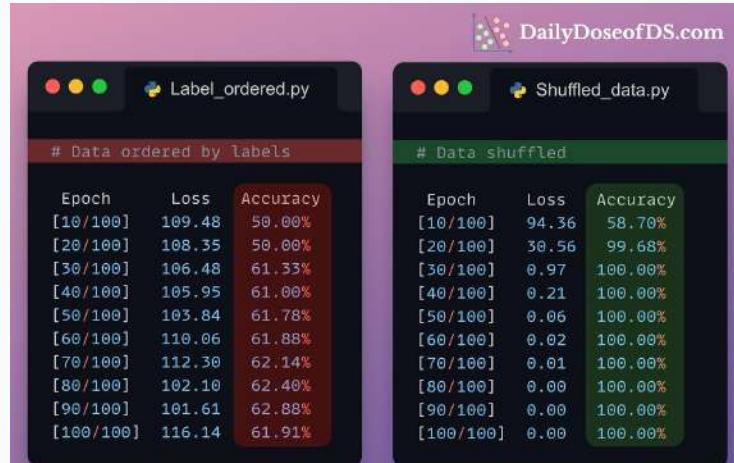
Here, we train two different neural networks:

- Version 1: The dataset is ordered by labels.
- Version 2: The dataset is properly shuffled by labels.



And, of course, before training, we ensure that both networks had the same initial weights, learning rate, and other settings.

The image depicts the epoch-by-epoch performance of the two models. On the left, we have the model trained on label-ordered data, and the one on the right was trained on the shuffled dataset.



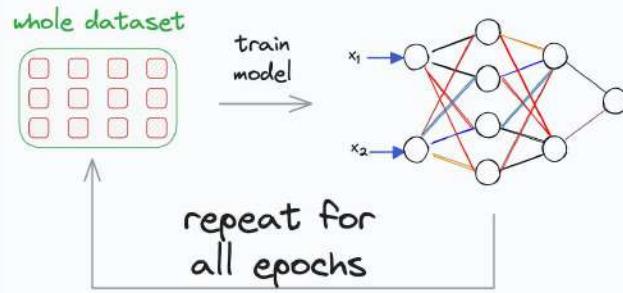
It is clear that the model receiving a label-ordered dataset miserably fails to converge while the other model, although overfits, shows that model has been learned effectively.

Why does that happen?

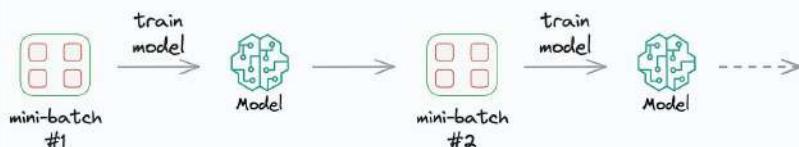
Now, if you think about it for a second, overall, both models received the same data, didn't they? Yet, the order in which the data was fed to these models totally determined their performance. I vividly remember that when I faced this issue, I knew that my data was ordered by labels.

Yet, it never occurred to me that ordering may influence the model performance because the data will always be the same regardless of the ordering.

But later, I realized that this point will only be valid when the model sees the entire data and updates the model weights in one go, i.e., in batch gradient descent, as depicted in this image.



But in the case of mini-batch gradient descent, the weights are updated after every mini-batch. Thus, the prediction and weight update on a subsequent mini-batch is influenced by the previous mini-batches.



In the context of label-ordered data, where samples of the same class are grouped together, mini-batch gradient descent will lead the model to learn patterns specific to the class it excessively saw early on in training. In contrast, randomly ordered data ensures that each mini-batch contains a balanced representation of classes. This allows the model to learn a more comprehensive set of features throughout the training process.

Of course, the idea of shuffling is not valid for time-series datasets as their temporal structure is important. The good thing is that if you happen to use, say, PyTorch DataLoader, you are safe. This is because it already implements shuffling. But if you have a custom implementation, ensure that you are not making any such error.

Before I end, one thing that you must ALWAYS remember when training neural networks is that these models can proficiently learn entirely non-existing patterns about your dataset. So never give them any chance to do so.

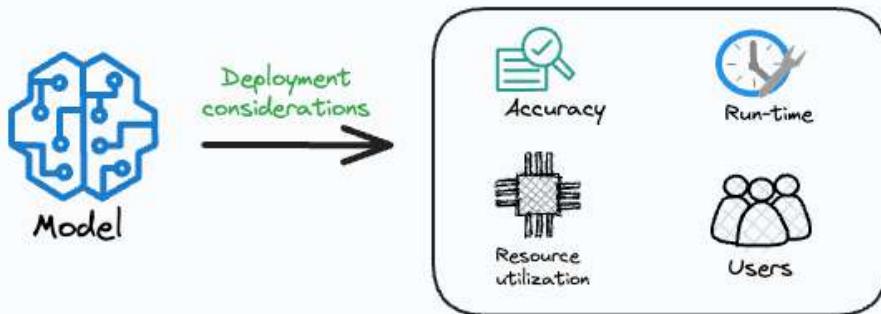
Model compression

Knowledge Distillation for Model Compression

Model accuracy alone (or an equivalent performance metric) rarely determines which model will be deployed.

This is because we also consider several operational metrics, such as:

- Inference Latency: Time taken by the model to return a prediction.
- Model size: The memory occupied by the model.
- Ease of scalability, etc.



In this chapter, let me share a technique (with a demo) called knowledge distillation, which is commonly used to compress ML models and contribute to the above operational metrics.

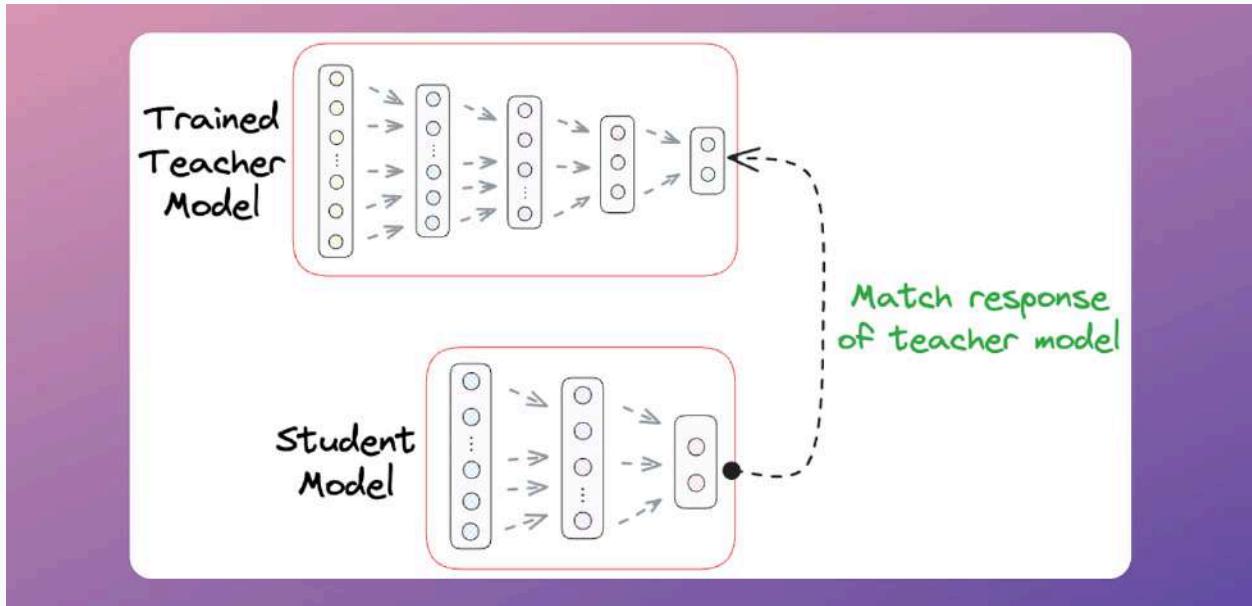
What is knowledge distillation?

In a gist, the idea is to train a smaller/simpler model (called the “student” model) that mimics the behavior of a larger/complex model (called the “teacher” model).



This involves two steps:

- Train the teacher model as we typically would.
- Train a student model that matches the output of the teacher model.



If we compare it to an academic teacher-student scenario, the student may not be as performant as the teacher.

But with consistent training, a smaller model may get (almost) as good as the larger one.

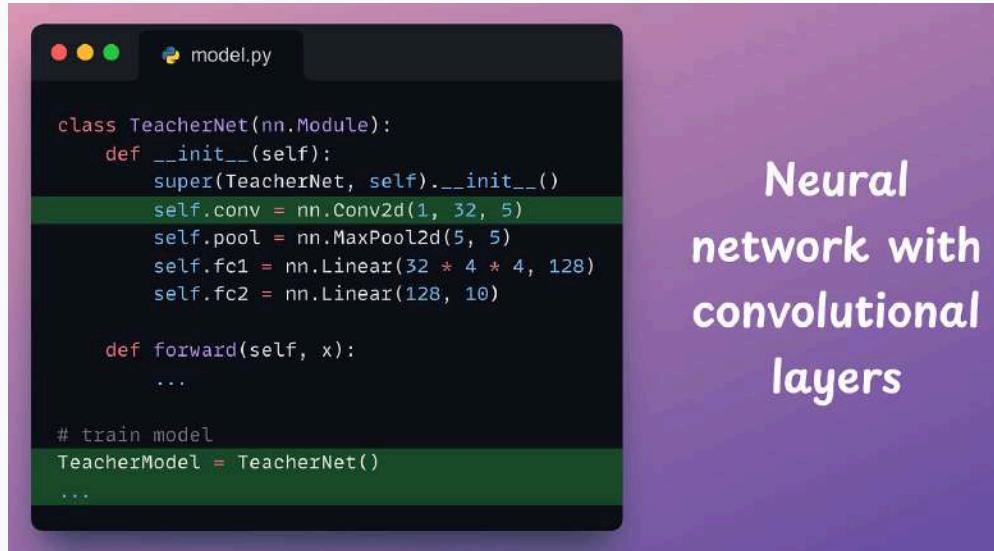
A classic example of a model developed in this way is DistillBERT. It is a student model of BERT.

- DistilBERT is approximately 40% smaller than BERT, which is a massive difference in size.
- Still, it retains approximately 97% of the BERT's capabilities.

Next, let's look at a demo.

Knowledge distillation demo

In the interest of time, let's say we have already trained the following CNN model on the MNIST dataset (I have provided the full Jupyter notebook towards the end, don't worry):



```

model.py

class TeacherNet(nn.Module):
    def __init__(self):
        super(TeacherNet, self).__init__()
        self.conv = nn.Conv2d(1, 32, 5)
        self.pool = nn.MaxPool2d(5, 5)
        self.fc1 = nn.Linear(32 * 4 * 4, 128)
        self.fc2 = nn.Linear(128, 10)

    def forward(self, x):
        ...
        # train model
TeacherModel = TeacherNet()
...

```

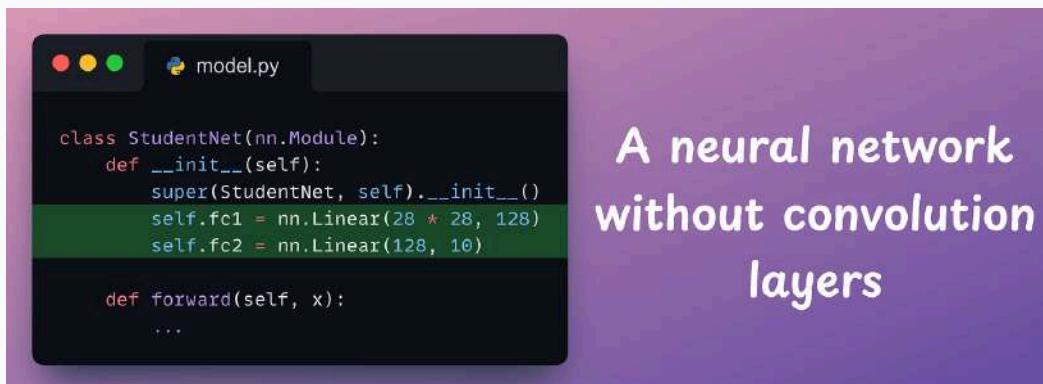
**Neural
network with
convolutional
layers**

The epoch-by-epoch training loss and validation accuracy is depicted below:

Epoch	Loss	Acc
1	0.20	98.79%
2	0.17	98.91%
3	0.15	98.78%
4	0.13	98.89%
5	0.11	98.63%

**Epoch-wise
loss and
accuracy**

Next, let's define a simpler model without any convolutional layers:



```

model.py

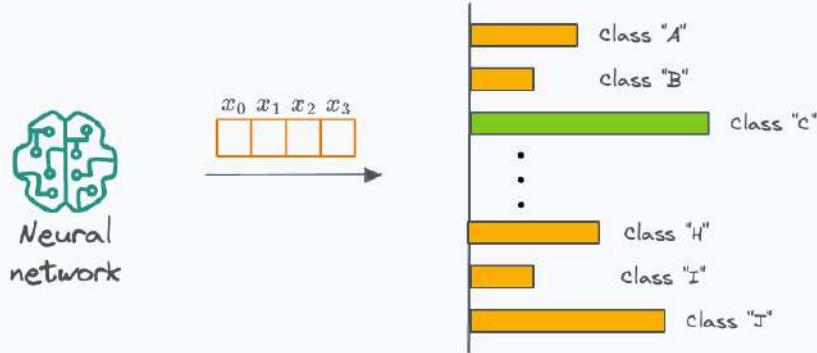
class StudentNet(nn.Module):
    def __init__(self):
        super(StudentNet, self).__init__()
        self.fc1 = nn.Linear(28 * 28, 128)
        self.fc2 = nn.Linear(128, 10)

    def forward(self, x):
        ...

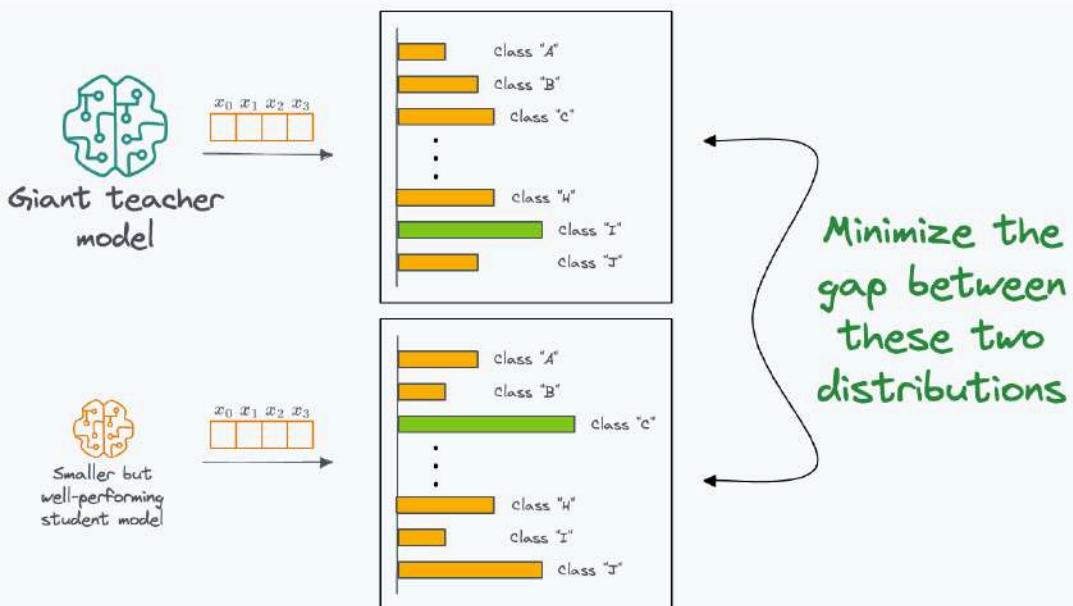
```

**A neural network
without convolution
layers**

Being a classification model, the output will be a probability distribution over the $\langle N \rangle$ classes:



Thus, we can train the student model such that its probability distribution matches that of the teacher model.



One way to do this is to use KL divergence as a loss function.

$$D_{KL}(P \parallel Q) = \sum_x P(x) \cdot \log \left(\frac{P(x)}{Q(x)} \right)$$

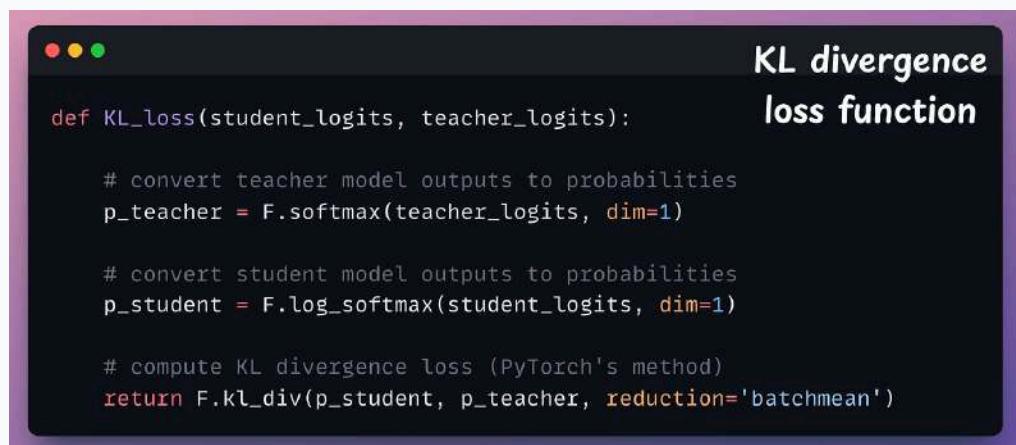
It measures how much information is lost when we use distribution Q to approximate distribution P.

A question for you: What will be the KL divergence if $P=Q$?

Thus, in our case:

- P will be the probability distribution from the teacher model.
- Q will be the probability distribution from the student model.

The loss function is implemented below:



```

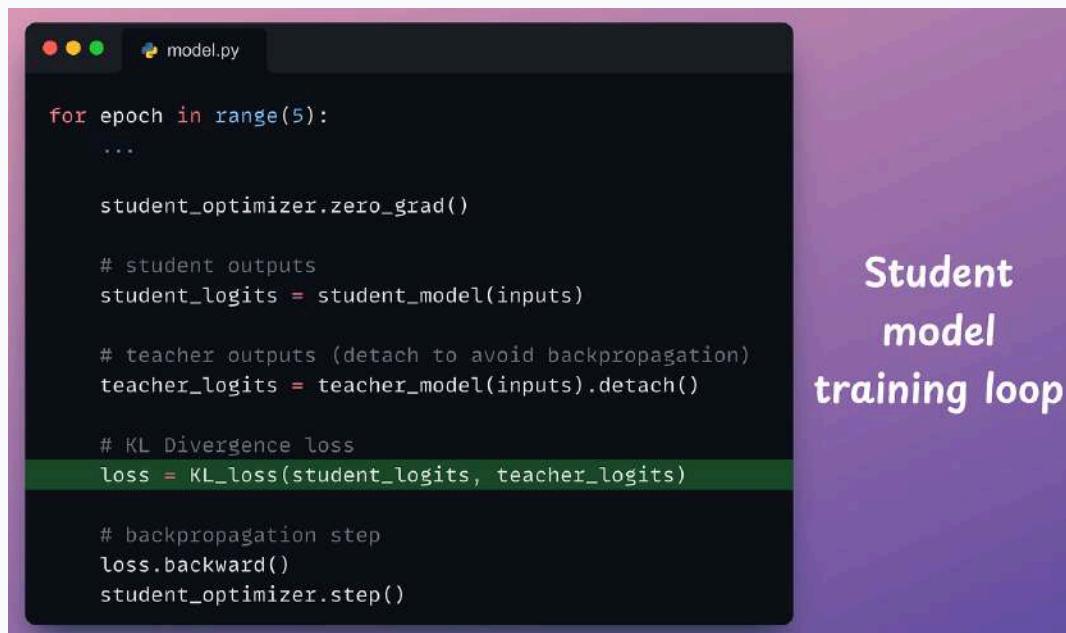
def KL_loss(student_logits, teacher_logits):
    # convert teacher model outputs to probabilities
    p_teacher = F.softmax(teacher_logits, dim=1)

    # convert student model outputs to probabilities
    p_student = F.log_softmax(student_logits, dim=1)

    # compute KL divergence loss (PyTorch's method)
    return F.kl_div(p_student, p_teacher, reduction='batchmean')

```

Finally, we train the student model as follows:



```

for epoch in range(5):
    ...
    student_optimizer.zero_grad()

    # student outputs
    student_logits = student_model(inputs)

    # teacher outputs (detach to avoid backpropagation)
    teacher_logits = teacher_model(inputs).detach()

    # KL Divergence loss
    loss = KL_loss(student_logits, teacher_logits)

    # backpropagation step
    loss.backward()
    student_optimizer.step()

```

Done!

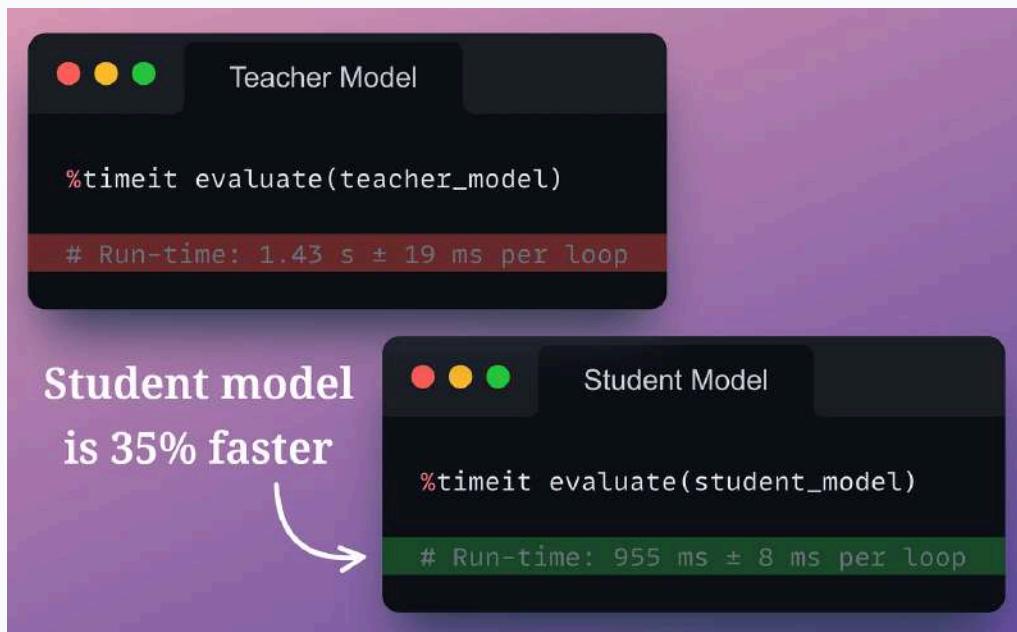
The following image compares the training loss and validation accuracy of the two models:



Of course, as shown in the highlighted lines above, the performance of the student model is not as good as the teacher model, which is expected.

However, it is still pretty promising, given that it was only composed of simple feed-forward layers.

Also, as depicted below, the student model is approximately 35% faster than the teacher model, which is a significant increase in the inference run-time of the model for about a 1-2% drop in the test accuracy.



That said, one of the biggest downsides of knowledge distillation is that one must still train a larger teacher model first to train the student model.

But in a resource-constrained environment, it may not be feasible to train a large teacher model.

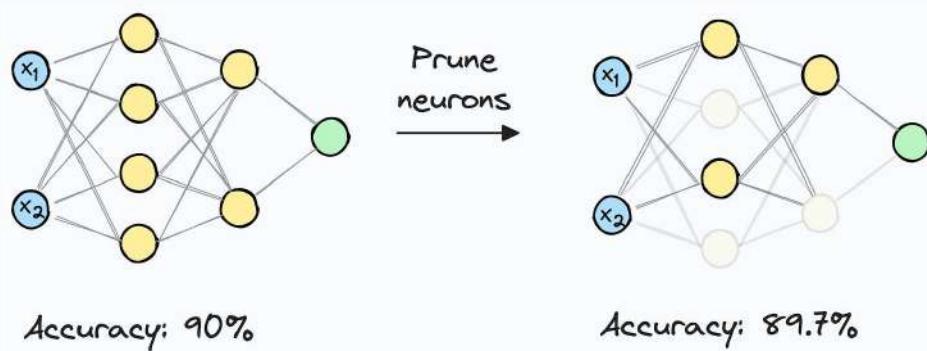
So this technique assumes that we are not resource-constrained at least in the development environment.

In the next chapter, let's discuss one more technique to compress ML models and reduce their memory footprint.

Activation Pruning

Once we complete network training, we are almost always left with plenty of useless neurons — ones that make nearly zero contribution to the network's performance, but they still consume memory.

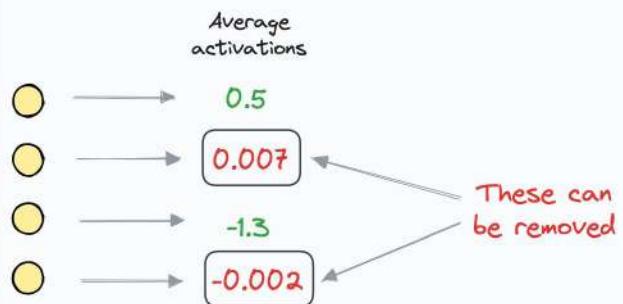
In other words, there is a high percentage of neurons, which, if removed from the trained network, will not affect the performance remarkably:



And, of course, I am not saying this as a random and uninformed thought. I have experimentally verified this over and over across my projects.

Here's the core idea.

After training is complete, we run the dataset through the model (no backpropagation this time) and analyze the average activation of individual neurons. Here, we often observe that many neuron activations are always close to near-zero values.



Thus, they can be pruned from the network, as they will have very little impact on the model's output.

For pruning, we can decide on a pruning threshold (λ) and prune all neurons whose activations are less than this threshold.

This makes intuitive sense as well.

More specifically, if a neuron rarely possesses a high activation value, then it is fair to assume that it isn't contributing to the model's output, and we can safely prune it.

The following table compares the accuracy of the pruned model with the original (full) model across a range of pruning thresholds (λ):

Original Model	Pruning Threshold (λ)	Validation Accuracy	Parameter Reduction
Pruned Models	-	96.7	-
	0.1	96.62	62%
	0.2	96.41	66%
	0.3	96.31	69%
	0.4	96.08	72%
	0.5	94.88	74%
	0.6	94.49	76%
	0.7	92.19	79%
	0.8	90.65	81%
	0.9	88.39	83%
	1	87.87	84%

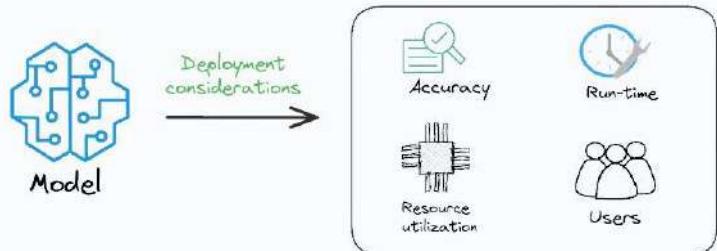


At a pruning threshold $\lambda=0.4$, the validation accuracy of the model drops by just 0.62%, but the number of parameters drops by 72%.

That is a huge reduction, while both models being almost equally good! Of course, there is a trade-off because we are not doing as well as the original model. But in many cases, especially when deploying ML models, accuracy is not the only primary metric that decides these.

Instead, several operational metrics like efficiency, speed, memory consumption, etc., are also a key deciding factor.

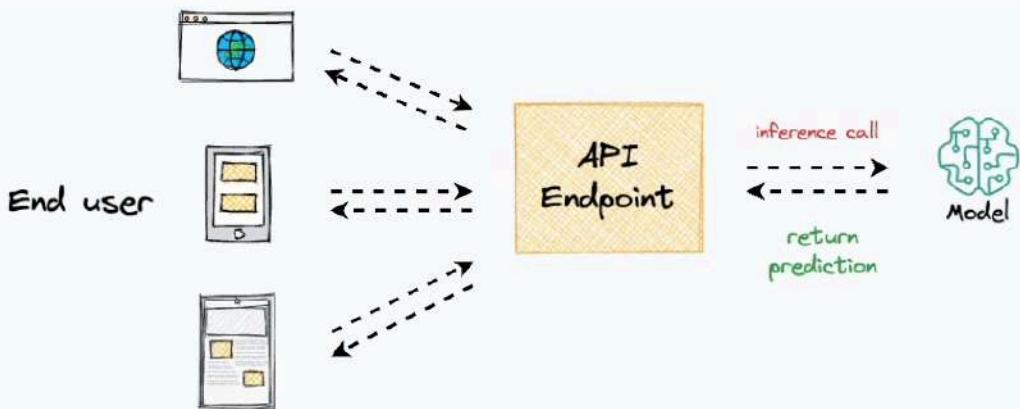
That is why model compression techniques are so crucial in such cases.



Deployment

Deploy ML Models from Jupyter Notebook

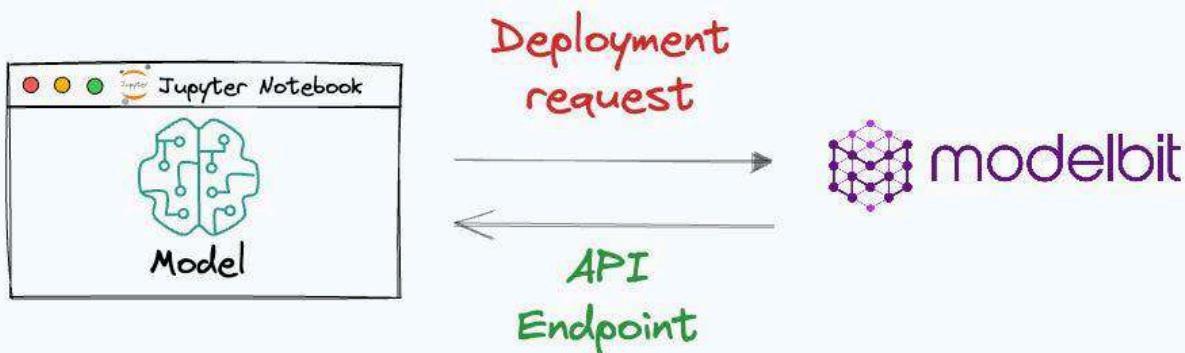
The core objective of model deployment is to obtain an API endpoint that can be used for inference purposes:



While this sounds simple, deployment is typically quite a tedious and time-consuming process. One must maintain environment files, configure various settings, ensure all dependencies are correctly installed, and many more.

So, in this chapter, I want to help you simplify this process. More specifically, we shall learn how to deploy any ML model right from a Jupyter Notebook in just three simple steps using the Modelbit API.

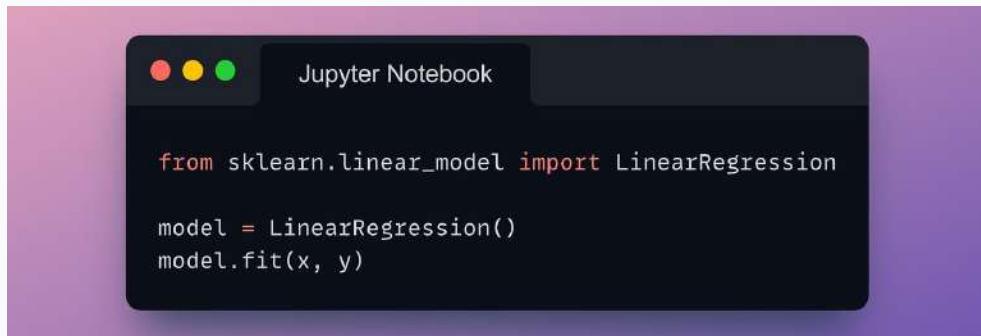
Modelbit lets us seamlessly deploy ML models directly from our Python notebooks (or git) to Snowflake, Redshift, and REST.



Deployment with Modelbit

Assume we have already trained our model.

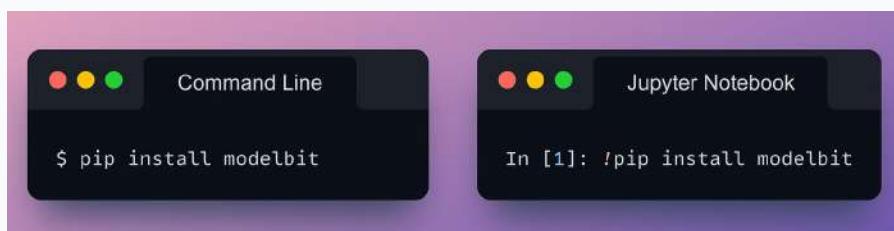
For simplicity, let's assume it to be a linear regression model trained using sklearn, but it can be any other model as well:



```
from sklearn.linear_model import LinearRegression
model = LinearRegression()
model.fit(x, y)
```

Let's see how we can deploy this model with Modelbit!

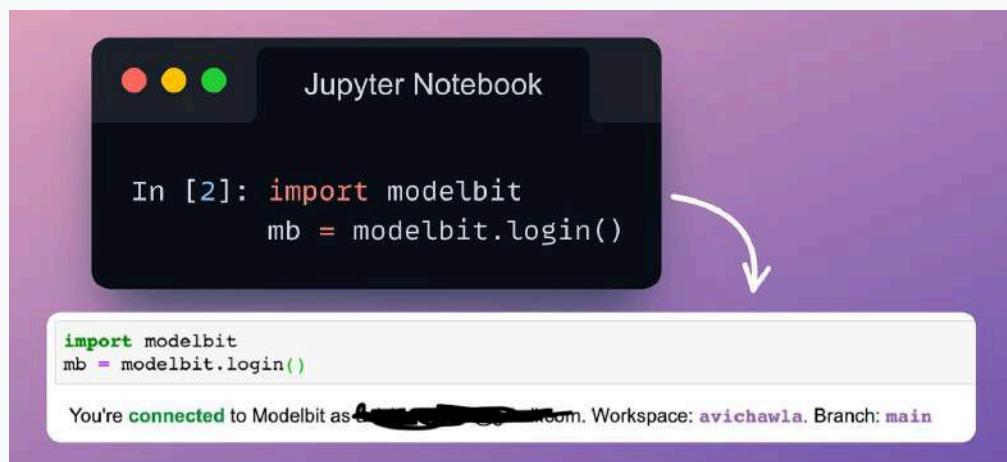
- First, we install the Modelbit package via pip:



```
$ pip install modelbit
```

In [1]: !pip install modelbit

- Next, we log in to Modelbit from our Jupyter Notebook (make sure you have created an account here: [Modelbit](#))



```
In [2]: import modelbit
mb = modelbit.login()
```

```
import modelbit
mb = modelbit.login()
```

You're connected to Modelbit as @92

- Finally, we deploy it, but here's an important point to note:

To deploy a model using Modelbit, we must define an inference function.

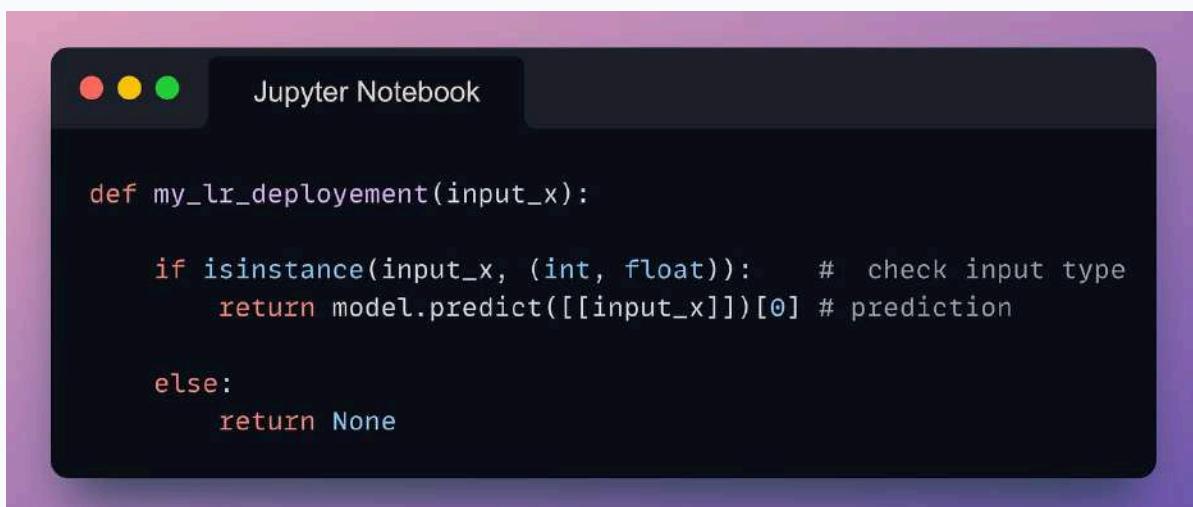
Simply put, this function contains the code that will be executed at inference. Thus, it will be responsible for returning the prediction.

```
def inference_function():
    # prediction code
```

A function that
returns model
predictions

We must specify the input parameters required by the model in this method. Also, we can name it anything we want.

For our linear regression case, the inference function can be as follows:



```
def my_lr_deployment(input_x):
    if isinstance(input_x, (int, float)):      # check input type
        return model.predict([[input_x]])[0] # prediction
    else:
        return None
```

- We define a function `my_lr_deployment()`.
- Next, we specify the input of the model as a parameter of this method.
- We validate the input for its data type.
- Finally, we return the prediction.

One good thing about Modelbit is that every dependency of the function (the model object in this case) is pickled and sent to production automatically along

with the function. Thus, we can reference any object in this method. Once we have defined the function, we can proceed with deployment as follows:

We have successfully deployed the model in three simple steps, that too, right from the Jupyter Notebook!

Once our model has been successfully deployed, it will appear in our Modelbit dashboard.

The Modelbit dashboard shows a deployment named "my_lr_deployment" owned by "Avi Chawla". It was last deployed 3 minutes ago and has 1 version.

As shown above, Modelbit provides an API endpoint. We can use it for inference purposes as follows:

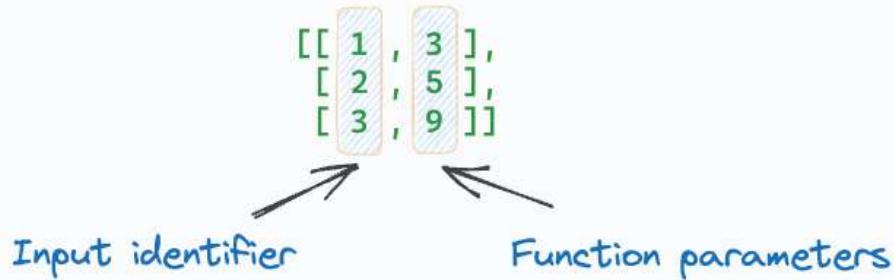
```
import json, requests

# invoke the API endpoint
requests.post("https://avichawla.app.modelbit.com/v1/my_lr_deployment/latest",
              headers={"Content-Type": "application/json"},
              data=json.dumps({"data": [[1, 3], [2, 5], [3, 9]]})).json()
```

Response

```
{"data": [[1, 12.41], # [Input ID, Output]
          [2, 19.33], # [Input ID, Output]
          [3, 33.16]] # [Input ID, Output]
 }
```

In the above request, data passed to the endpoint is a list of lists.



The first number in the list is the input ID. All entries following the ID in a list are the function parameters.

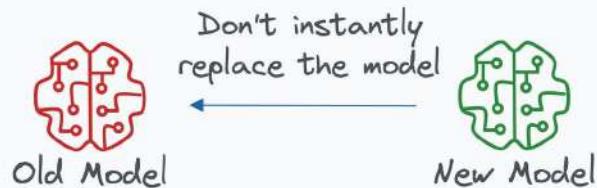
Lastly, we can also specify specific versions of the libraries or Python used while deploying our model. This is depicted below:

```
mb.deploy(my_lr_deployment,
          python_packages = ["scikit-learn=1.1.2", "pandas=1.5.0"],
          python_version   = "3.9")
```

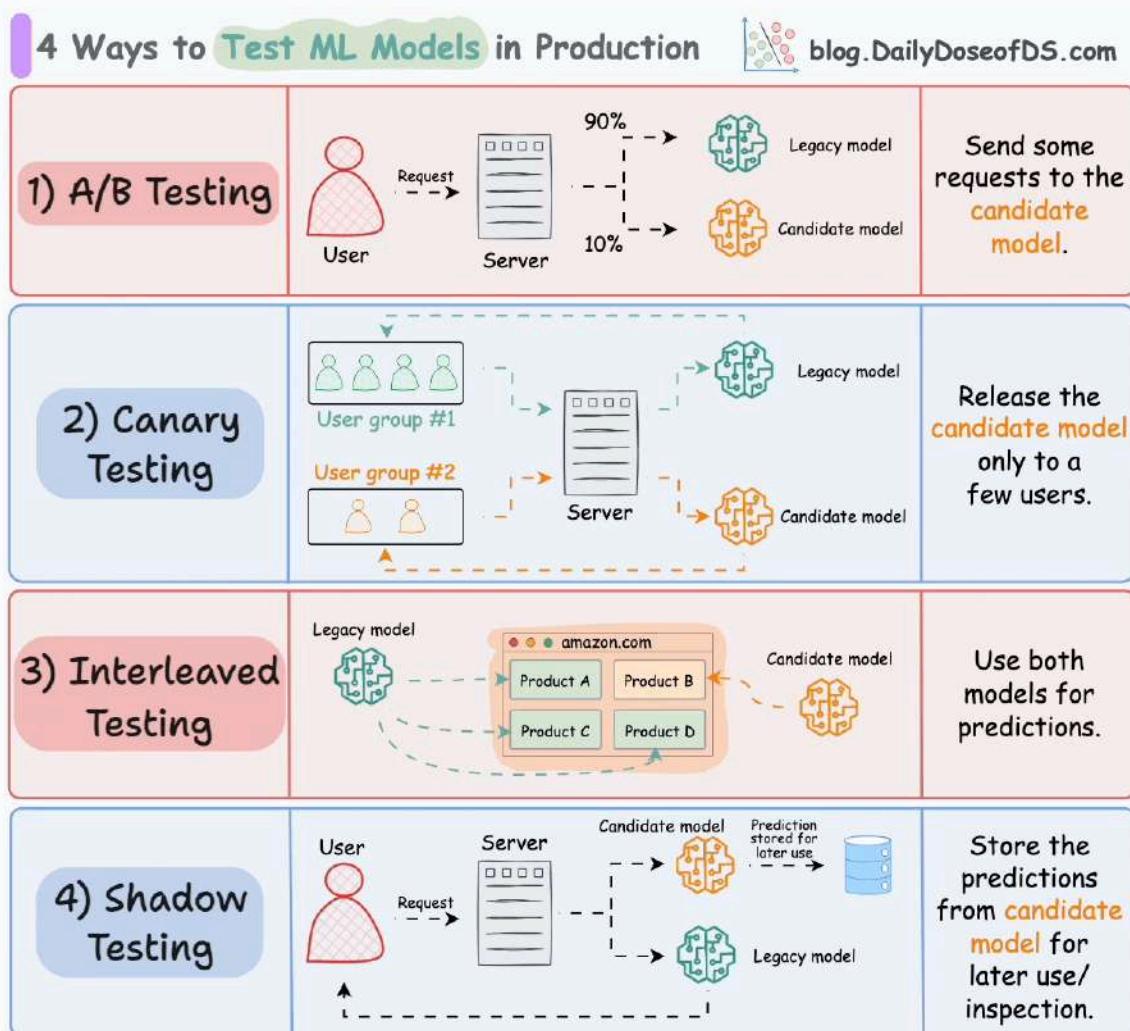
Isn't that cool, simple, and elegant over traditional deployment approaches?

4 Ways to Test ML Models in Production

Despite rigorously testing an ML model locally (on validation and test sets), it could be a terrible idea to instantly replace the previous model with a new model.

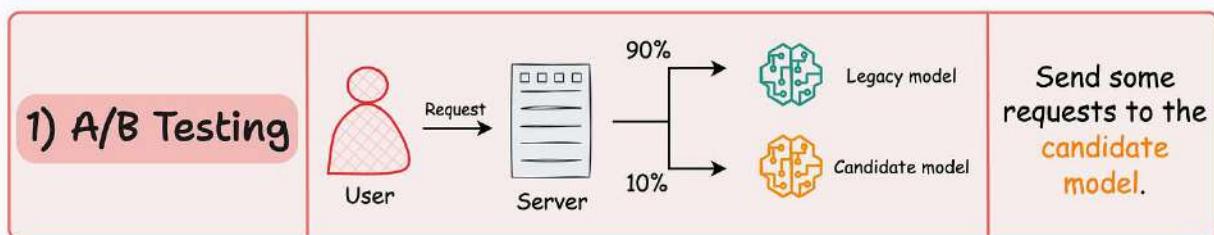


A more reliable strategy is to test the model in production (yes, on real-world incoming data). While this might sound risky, ML teams do it all the time, and it isn't that complicated. The following visual depicts 4 common strategies to do so:



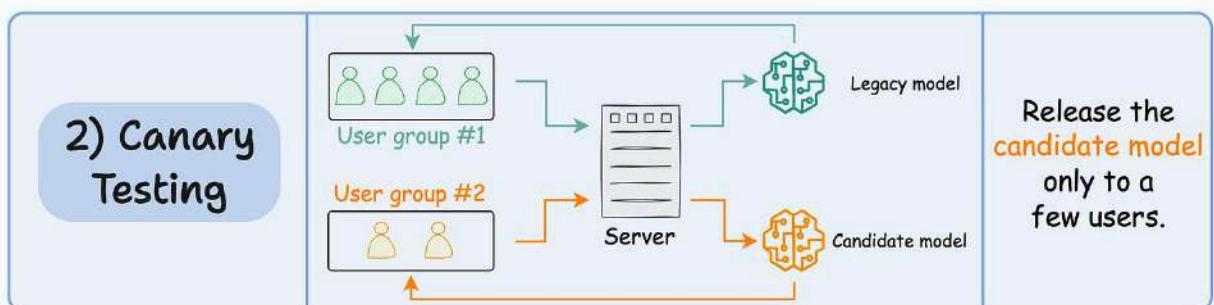
- The current model is called the legacy model.
- The new model is called the candidate model.

#1) A/B testing



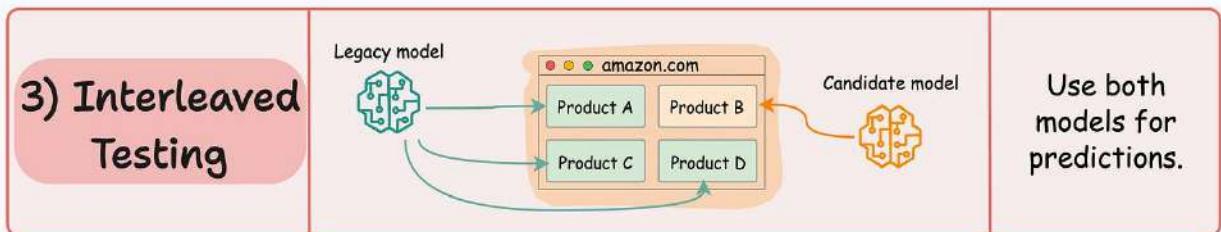
- Distribute the incoming requests non-uniformly between the legacy model and the candidate model.
- Intentionally limit the exposure of the candidate model to avoid any potential risks. Thus, the number of requests sent to the candidate model must be low.

#2) Canary testing



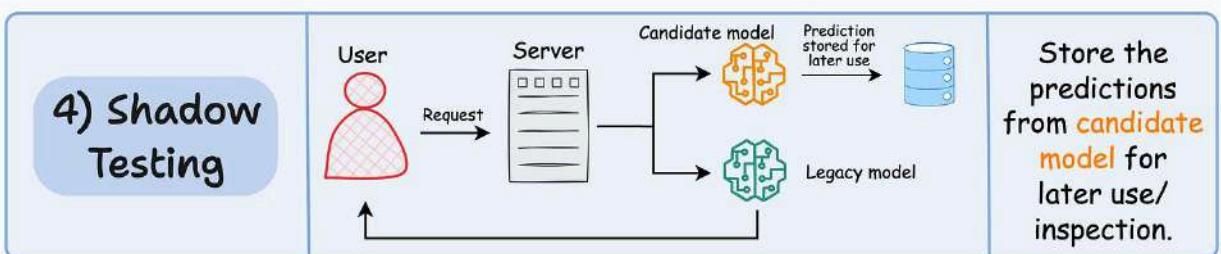
- In A/B testing, since traffic is randomly redirected to either model irrespective of the user, it can potentially affect all users.
- In canary testing, the candidate model is released to a small subset of users in production and gradually rolled out to more users.

#3) Interleaved testing



- This involves mixing the predictions of multiple models in the response.
- Consider Amazon's recommendation engine. In interleaved deployments, some product recommendations displayed on the homepage can come from the legacy model, while some can come from the candidate model.

#4) Shadow testing



- All of the above techniques affect some (or all) users.
- Shadow testing (or dark launches) lets us test a new model in a production environment without affecting the user experience.
- The candidate model is deployed alongside the existing legacy model and serves requests like the legacy model. However, the output is not sent back to the user. Instead, the output is logged for later use to benchmark its performance against the legacy model.
- We explicitly deploy the candidate model instead of testing offline because the production environment is difficult to replicate offline.

Shadow testing offers risk-free testing of the candidate model in a production environment.

Version Controlling and Model Registry

Real-world ML deployment is never just about “deployment” — host the model somewhere, obtain an API endpoint, integrate it into the application, and you are done!

This is because, in reality, plenty of things must be done post-deployment to ensure the model’s reliability and performance.

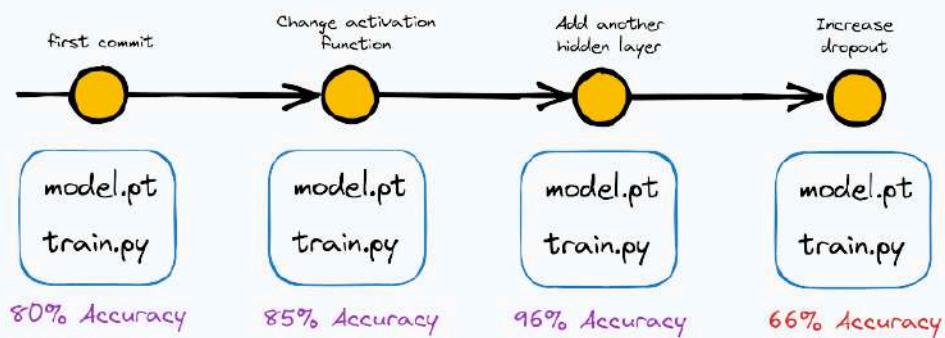
#1) Version control

To begin, it is immensely crucial to version control ML deployments. You may have noticed this while using ChatGPT, for instance.



But updating does not simply mean overwriting the previous version.

Instead, ML models are always version-controlled (using git tools).



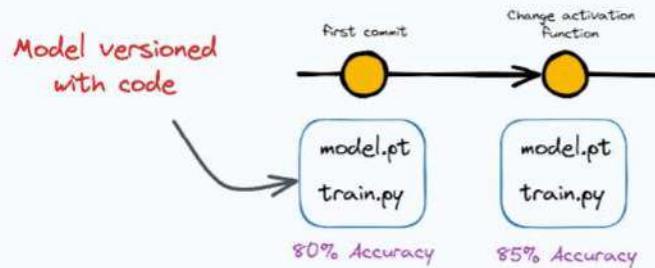
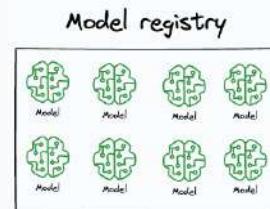
The advantages of version-controlling ML deployments are pretty obvious:

- In case of sudden mishaps post-deployment, we can instantly roll back to an older version.
- We can facilitate parallel development with branching, and many more.

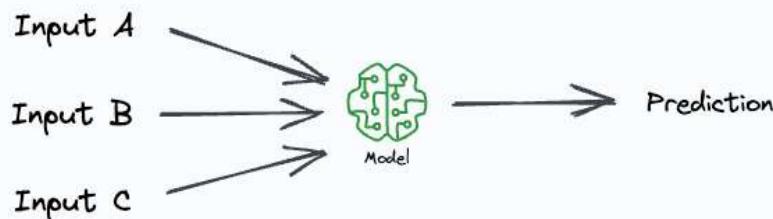
#2) Model registry

Another practical idea is to maintain a model registry for deployments. Let's understand what it is.

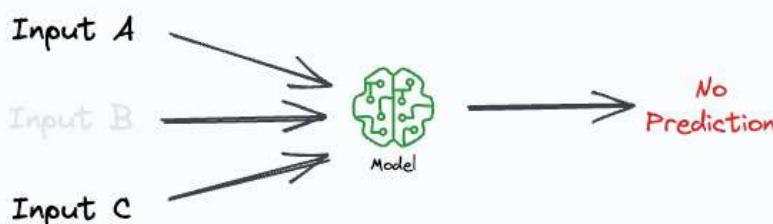
Simply put, a model registry can be considered repository of models. See, typically, we might be inclined to version our code and the ML model together:



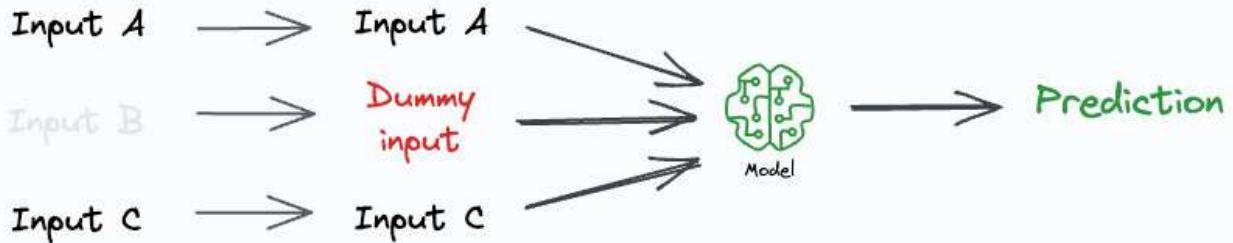
However, when we use a model registry, we version models separately from the code. Let me give you an intuitive example to understand this better. Imagine our deployed model takes three inputs to generate a prediction:



While writing the inference code, we overlooked that, at times, one of the inputs might be missing. We realized this by analyzing the model's logs.



We may want to fix this quickly (at least for a while) before we decide on the next steps more concretely. Thus, we may decide to update the inference code by assigning a dummy value for the missing input.

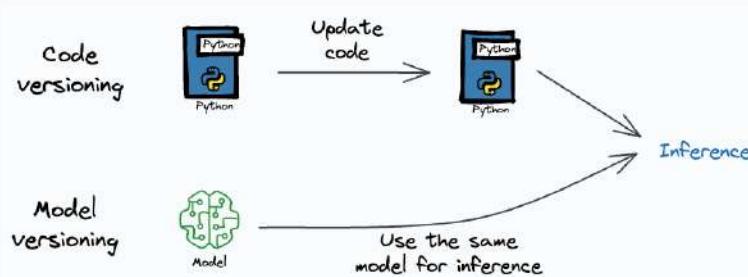


This will allow the model to still process the incoming request.

Let me ask you a question: “Did we update the model?”

No, right?

Here, we only need to update the inference code. The model will remain the same.



But if we were to version the model and code together, it would lead to a redundant model and take up extra space.

However, by maintaining a model registry:

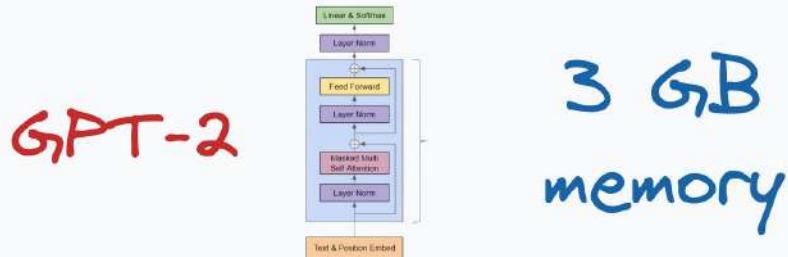
- We can only update the inference code.
- Avoid pushing a new (yet unwanted) model to deployment.

This makes intuitive sense as well, doesn't it?

LLMs

Where Did the GPU Memory Go?

GPT-2 (XL) has 1.5 Billion parameters, and its parameters consume ~3GB of memory in 16-bit precision.



Under 16-bit precision, one parameter takes up 2 bytes of memory, so 1.5B parameters will consume 3GB of memory.

What's your estimate for the minimum memory needed to train GPT-2 on a single GPU?

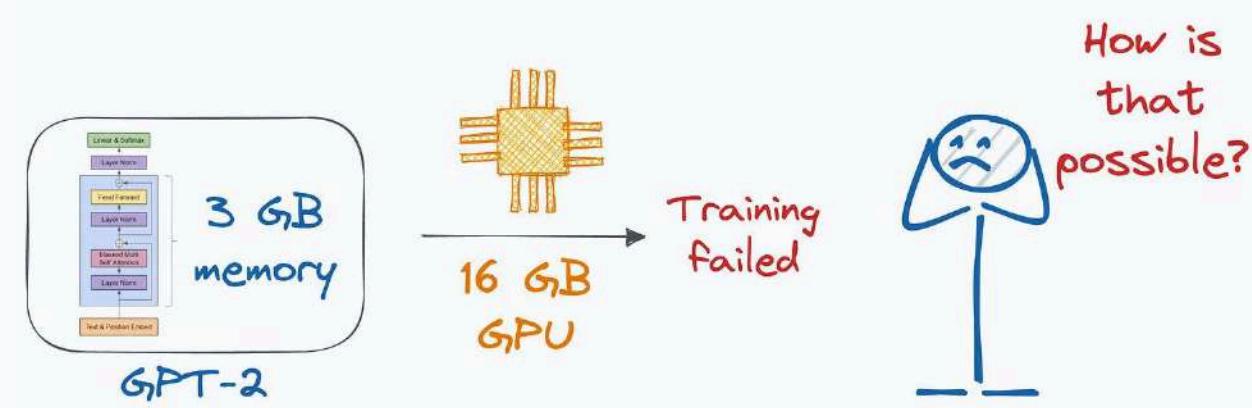
- Optimizer → Adam
- Batch size → 32
- Number of transformer layers → 48
- Sequence length → 1000

What's your estimate for the minimum memory needed to train a GPT-2 model (size 3GBs) on a single GPU?

- 4-6 GB
- 8-10 GB
- 12-15 GB
- 32-35 GB
- 50+ GB

The answer might surprise you.

One can barely train a 3GB GPT-2 model on a single GPU with 32GB of memory.



But how could that be even possible? Where does all the memory go?

Let's understand.

There are so many fronts on which the model consistently takes up memory during training.

#1) Optimizer states, gradients, and parameter memory

Mixed precision training is widely used to speed up model training.

As the name suggests, the idea is to utilize lower-precision `float16` (whenever feasible, like in convolutions and matrix multiplications) along with `float32` — that is why the name “mixed precision.”

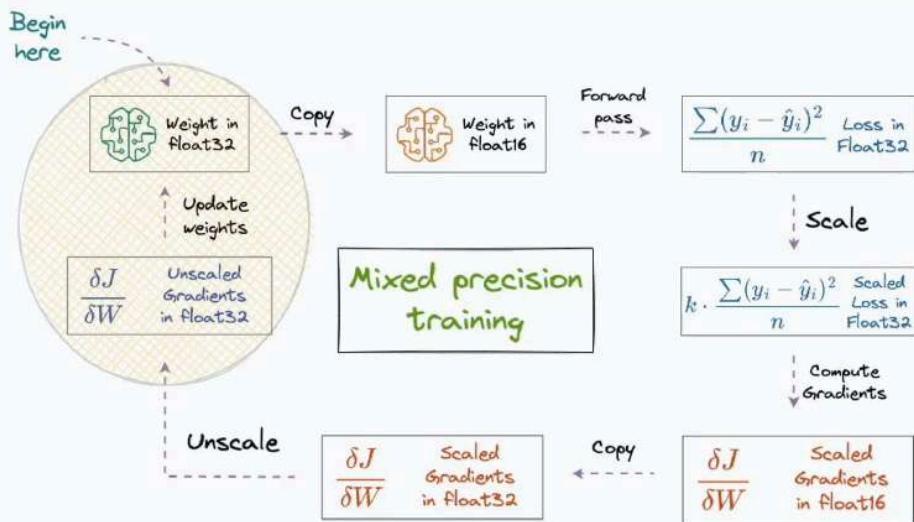
Both the forward and backward propagation are performed using the 16-bit representations of weights and gradients.

Thus, if the model has Φ parameters, then:

- Weights will consume $2^*\Phi$ bytes.
- Gradients will consume $2^*\Phi$ bytes.

Here, the figure “2” represents a memory consumption of 2 bytes/parameter (16-bit).

Moreover, the updates at the end of the backward propagation are still performed under 32-bit for effective computation. I am talking about the circled step in the image below:



Adam is one of the most popular optimizers for model training.

While many practitioners use it just because it is popular, they don't realize that during training, Adam stores two optimizer states to compute the updates — momentum and variance of the gradients:

$$\theta_t = \theta_{t-1} - \alpha \cdot \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$$

Momentum

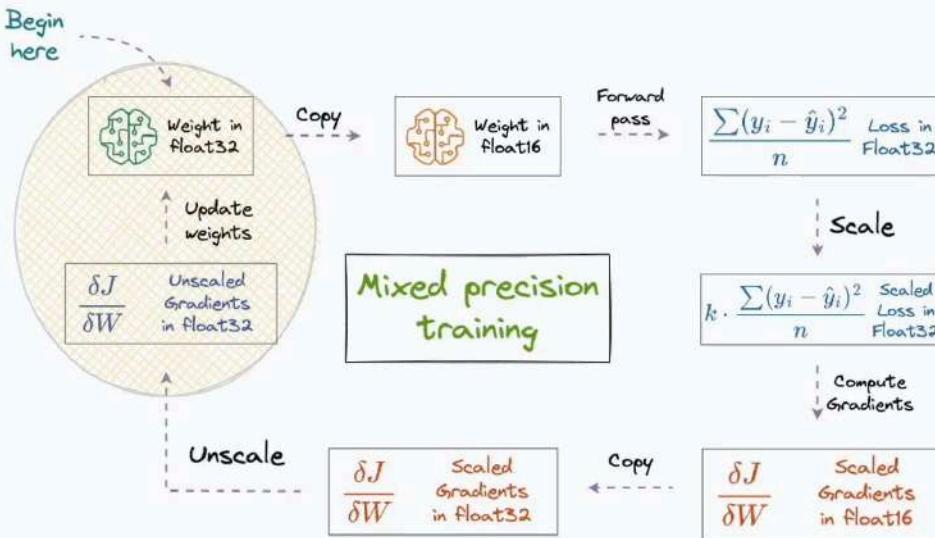
Variance

Thus, if the model has Φ parameters, then these two optimizer states will consume:

- $4^*\Phi$ bytes for momentum.

- Another $4^*\Phi$ bytes for variance.

Here, the figure “4” represents a memory consumption of 4 bytes/parameter (32-bit).



Lastly, as shown in the figure above, the final updates are always adjusted in the 32-bit representation of the model weights. This leads to:

- Another $4^*\Phi$ bytes for model parameters.

Let's sum them up:

$$16 * \phi$$

$$2 * \phi + 2 * \phi + 4 * \phi + 4 * \phi + 4 * \phi$$

Parameters in 16-bit	Gradients in 16-bit	Momentum in 32-bit	Variance in 32-bit	Parameters in 32-bit
-------------------------	------------------------	-----------------------	-----------------------	-------------------------

That's $16^*\Phi$, or 24GB of memory, which is ridiculously higher than the 3GB memory utilized by 16-bit parameters.

And we haven't considered everything yet.

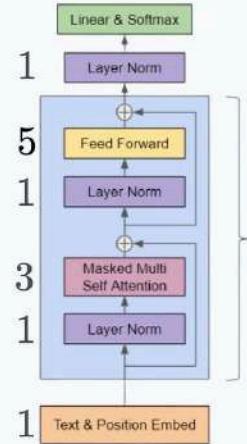
#2) Activations

For big deep learning models, like LLMs, Activations take up significant memory during training.

More formally, the total number of activations computed in one transform block of GPT-2 are:

$$12 * H * S * B$$

Hidden size Sequence length Batch size



Thus, across all transformer blocks, this comes out to be:

$$12 * H * S * B * N$$

Hidden size Sequence length Batch size Transformer layers

This is the configuration for GPT2-XL:

$$12 * 1600 * 1000 * 32 * 48$$

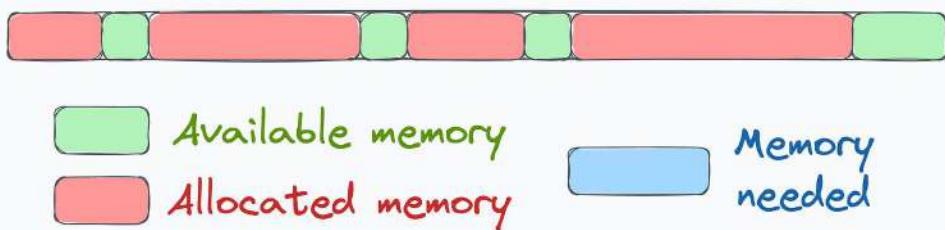
Hidden size Sequence length Batch size Transformer layers

This comes out to be ~30B activations. As each activation is represented in 16-bit, all activations collectively consume 60GB of memory.

With techniques like gradient checkpointing (discussed in the previous chapter), this could be brought down to about 8-9GB at the expense of 25-30% more run-time.

This technique takes complete memory consumption to about 32-35GB range, which I mentioned earlier, for a meager 3GB model, and that too with a pretty small batch size of just 32. On top of this, there are also some more overheads involved, like memory fragmentation.

It occurs when there are small, unused gaps between allocated memory blocks, leading to inefficient use of the available memory.



Memory allocation requests fail because of the unavailability of contiguous memory blocks.

Conclusion

In the above discussion, we considered a relatively small model — GPT-2 (XL) with 1.5 Billion parameters, which is tiny compared to the scale of models being trained these days.

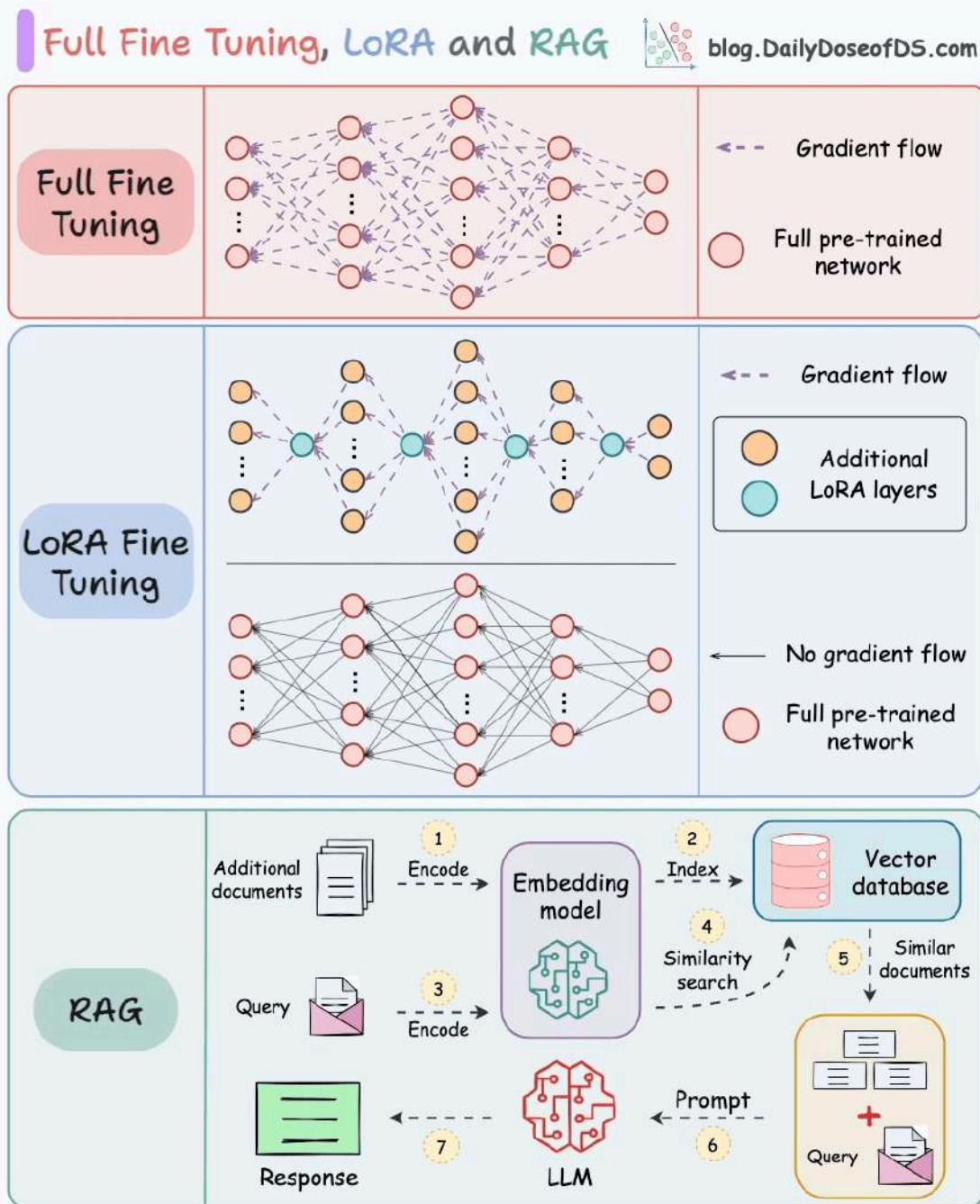
However, the discussion may have helped you reflect on the inherent challenges of building LLMs. Many people often say that GPTs are only about stacking more and more layers in the model and making the network bigger.

If it was that easy, everybody would have been doing it. From this discussion, you may have understood that it's not as simple as appending more layers.

Even one additional layer can lead to multiple GBs of additional memory requirement. Multi-GPU training is at the forefront of these models, which we covered in an earlier chapter in this book.

Full-model Fine-tuning vs. LoRA vs. RAG

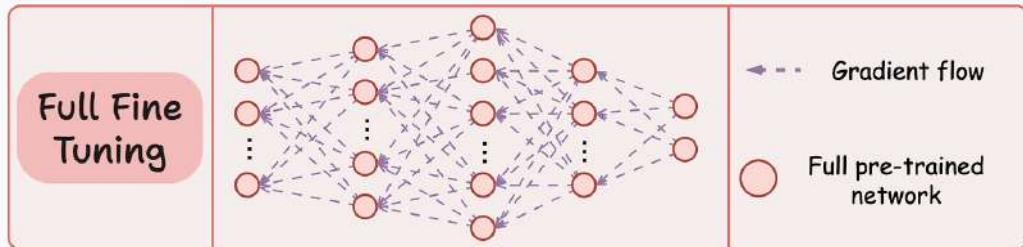
Here's a visual which illustrates "full-model fine-tuning," "fine-tuning with LoRA," and "retrieval augmented generation (RAG)."



All three techniques are used to augment the knowledge of an existing model with additional data.

#1) Full fine-tuning

Fine-tuning means adjusting the weights of a pre-trained model on a new dataset for better performance.

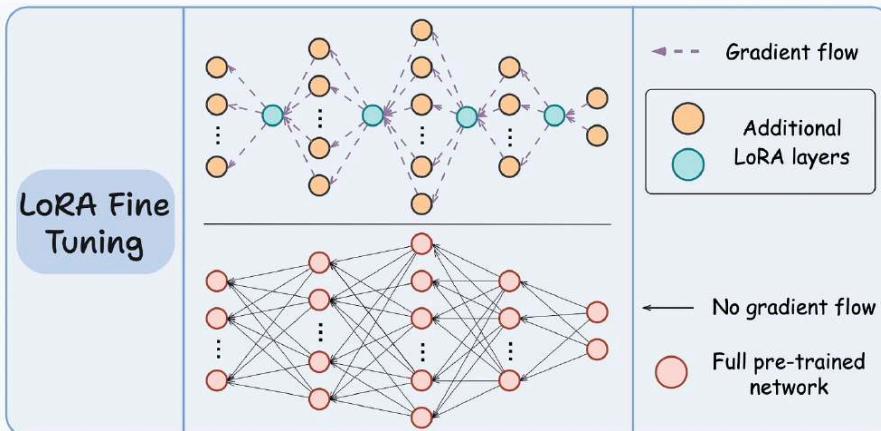


While this fine-tuning technique has been successfully used for a long time, problems arise when we use it on much larger models — LLMs, for instance, primarily because of:

- Their size.
- The cost involved in fine-tuning all weights.
- The cost involved in maintaining all large fine-tuned models.

#2) LoRA fine-tuning

LoRA fine-tuning addresses the limitations of traditional fine-tuning. The core idea is to decompose the weight matrices (some or all) of the original model into low-rank matrices and train them instead. For instance, in the graphic below, the bottom network represents the large pre-trained model, and the top network represents the model with LoRA layers.



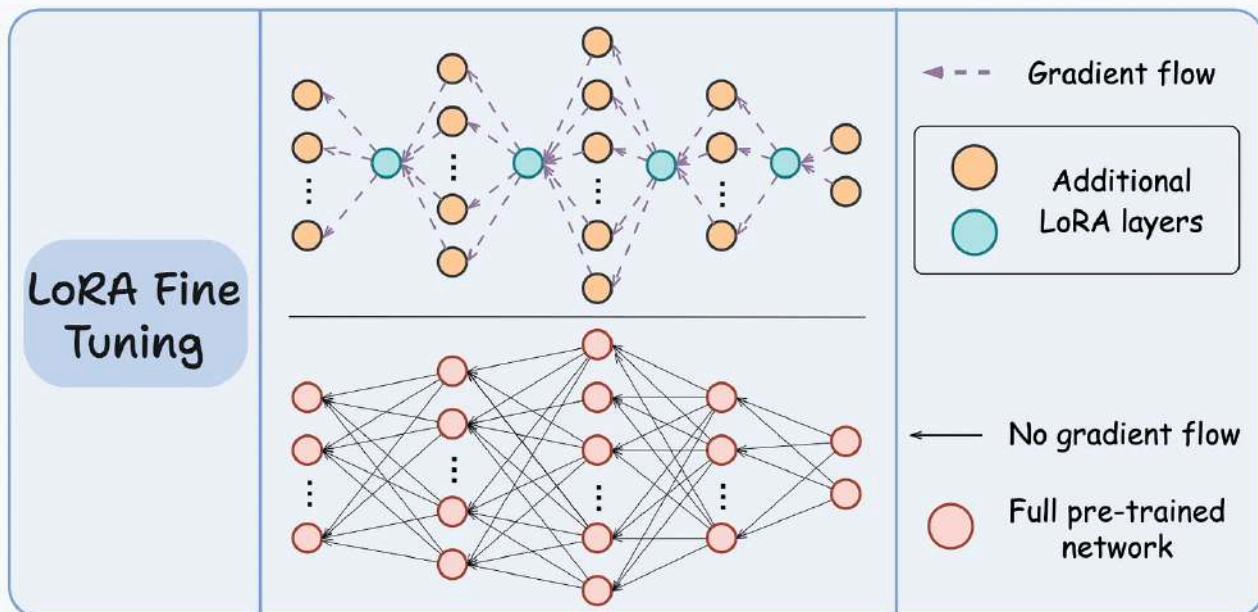
The idea is to train only the LoRA network and freeze the large model.

Looking at the above visual, you might think:

But the LoRA model has more neurons than the original model. How does that help? To understand this, you must make it clear that neurons don't have anything to do with the memory of the network.

They are just used to illustrate the dimensionality transformation from one layer to another.

It is the weight matrices (or the connections between two layers) that take up memory. Thus, we must be comparing these connections instead:

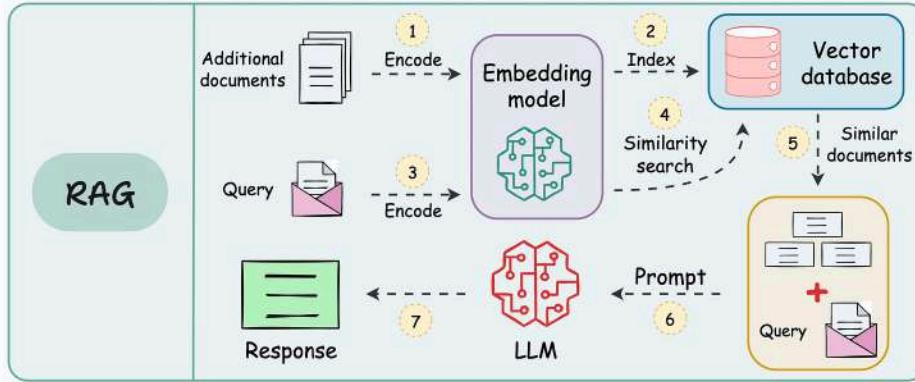


Looking at the above visual, it is pretty clear that the LoRA network has relatively very few connections.

#3) RAG

Retrieval augmented generation (RAG) is another pretty cool way to augment neural networks with additional information, without having to fine-tune the model.

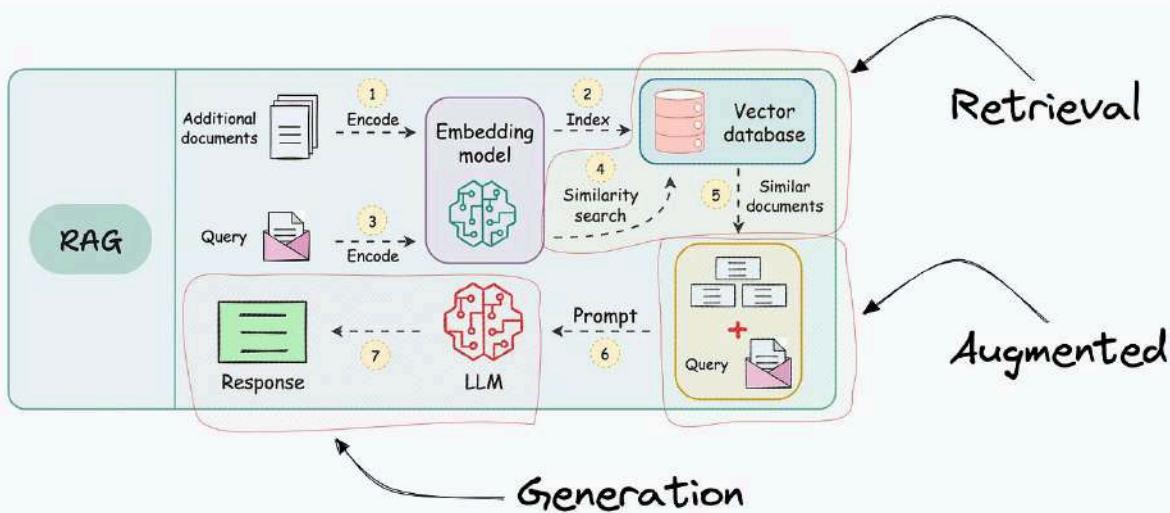
This is illustrated below:



There are 7 steps, which are also marked in the above visual:

- Step 1-2: Take additional data, and dump it in a vector database after embedding. (This is only done once. If the data is evolving, just keep dumping the embeddings into the vector database. There's no need to repeat this again for the entire data)
- Step 3: Use the same embedding model to embed the user query.
- Step 4-5: Find the nearest neighbors in the vector database to the embedded query.
- Step 6-7: Provide the original query and the retrieved documents (for more context) to the LLM to get a response.

In fact, even its name entirely justifies what we do with this technique:



- Retrieval: Accessing and retrieving information from a knowledge source, such as a database or memory.
- Augmented: Enhancing or enriching something, in this case, the text generation process, with additional information or context.
- Generation: The process of creating or producing something, in this context, generating text or language.

Of course, there are many problems with RAG too, such as:

- RAGs involve similarity matching between the query vector and the vectors of the additional documents. However, questions are structurally very different from answers.
- Typical RAG systems are well-suited only for lookup-based question-answering systems. For instance, we cannot build a RAG pipeline to summarize the additional data. The LLM never gets info about all the documents in its prompt because the similarity matching step only retrieves top matches.

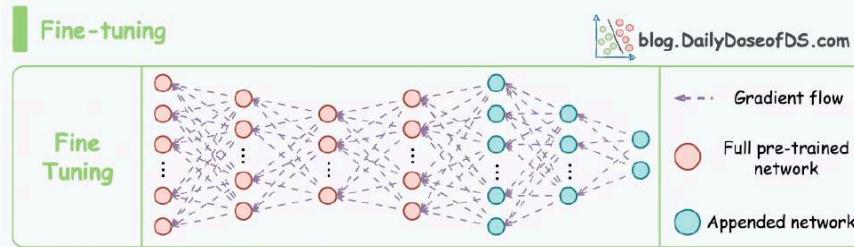
So, it's pretty clear that RAG has both pros and cons.

- We never have to fine-tune the model, which saves a lot of computing power.
- But this also limits the applicability to specific types of systems.

Let's continue the discussion on LLM fine-tuning in the next chapter.

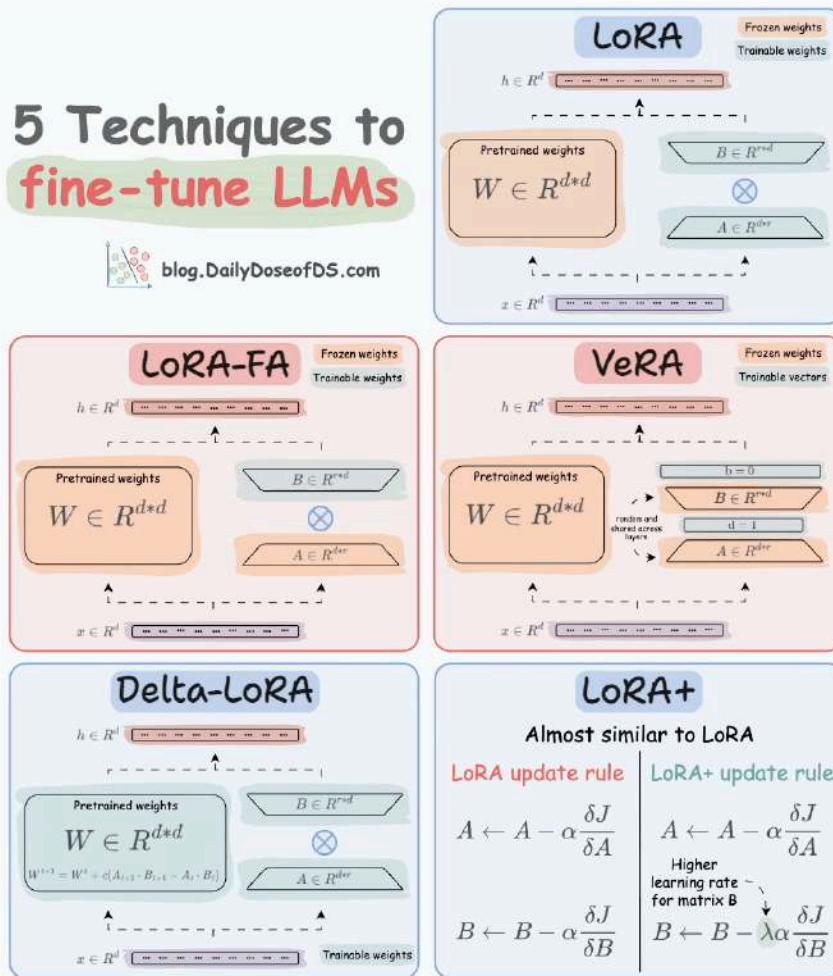
5 LLM Fine-tuning Techniques

Traditional fine-tuning (depicted below) is infeasible with LLMs because these models have billions of parameters and are hundreds of GBs in size, and not everyone has access to such computing infrastructure.

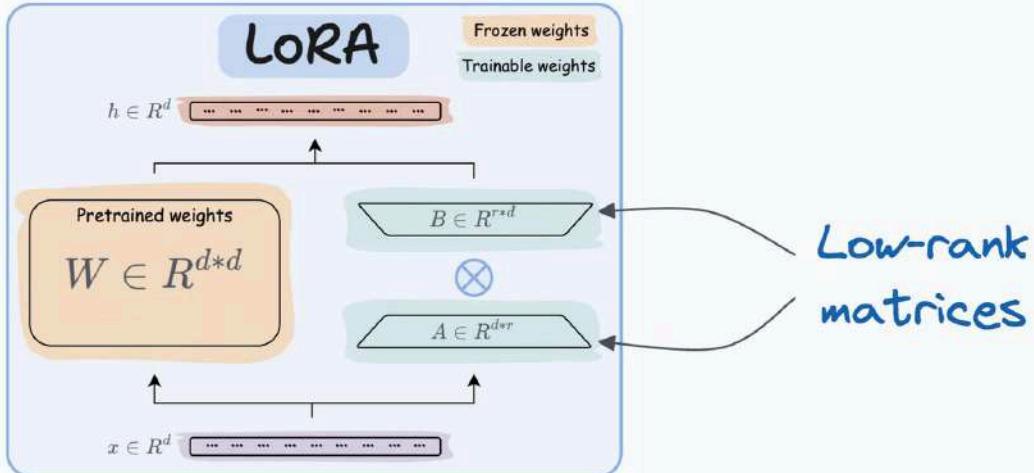


But today, we have many optimal ways to fine-tune LLMs, and five popular techniques are depicted below:

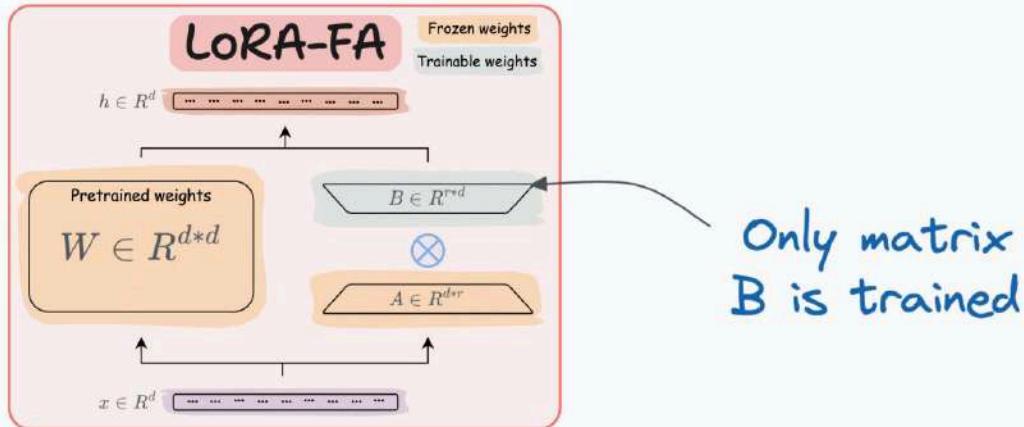
5 Techniques to fine-tune LLMs



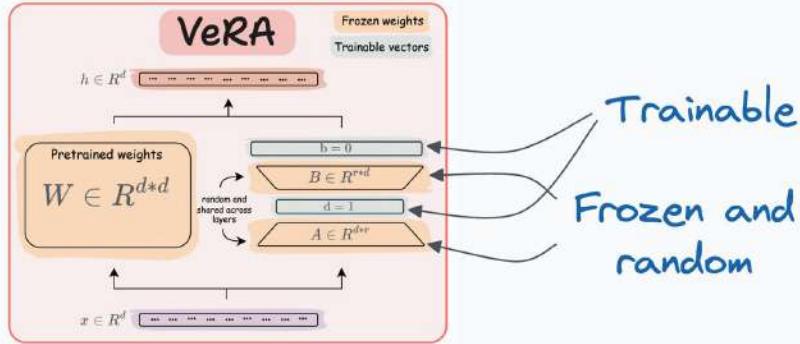
- LoRA: Add two low-rank matrices A and B alongside weight matrices, which contain the trainable parameters. Instead of fine-tuning W, adjust the updates in these low-rank matrices.



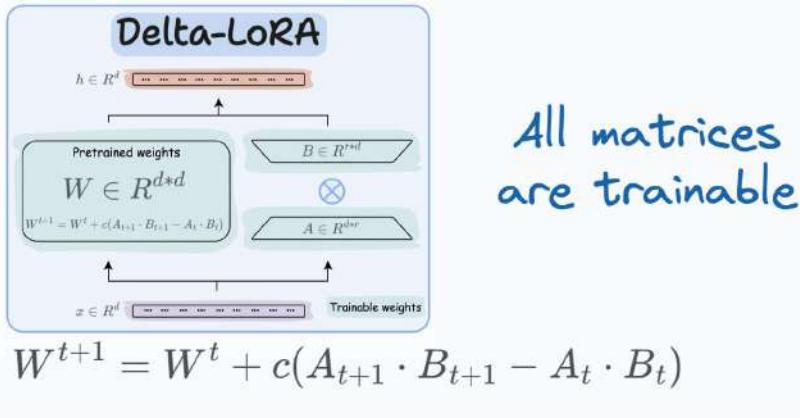
- LoRA-FA: While LoRA considerably decreases the total trainable parameters, it still requires substantial activation memory to update the low-rank weights. LoRA-FA (FA stands for Frozen-A) freezes the matrix A and only updates matrix B.



- VeRA: In LoRA, every layer has a different pair of low-rank matrices A and B, and both matrices are trained. In VeRA, however, matrices A and B are frozen, random, and shared across all model layers. VeRA focuses on learning small, layer-specific scaling vectors, denoted as b and d, which are the only trainable parameters in this setup.



- Delta-LoRA: Here, in addition to training low-rank matrices, the matrix W is also adjusted but not in the traditional way. Instead, the difference (or delta) between the product of the low-rank matrices A and B in two consecutive training steps is added to W :



$$W^{t+1} = W^t + c(A_{t+1} \cdot B_{t+1} - A_t \cdot B_t)$$

- LoRA+: In LoRA, both matrices A and B are updated with the same learning rate. Authors found that setting a higher learning rate for matrix B results in more optimal convergence.



LoRA+ assigns
higher learning
rate for updating
matrix B

Classical ML

ML Fundamentals

Training and Inference Time Complexity of 10 ML Algorithms

Here's the run-time complexity of the 10 most popular ML algorithms.

Time Complexity of 10 Most Popular ML Algorithms		Training	Inference
	Linear Regression (OLS)	$O(nm^2 + m^3)$	$O(m)$
	Linear Regression (SGD)	$O(n_{epoch}nm)$	$O(m)$
	Logistic Regression (Binary)	$O(n_{epoch}nm)$	$O(m)$
	Logistic Regression (Multiclass OvR)	$O(n_{epoch}nmc)$	$O(mc)$
	Decision Tree	$O(n \cdot \log(n) \cdot m)$ $O(n^2 \cdot m)^*$ Worst case	$O(d_{tree})$
	Random Forest Classifier	$O(n_{trees} \cdot n \cdot \log(n) \cdot m)$	$O(n_{trees} \cdot d_{tree})$
	Support Vector Machines (SVMs)	$O(n^2m + n^3)$	$O(m \cdot n_{SV})$
	k-Nearest Neighbors	—	$O(nm)$
$P(B A) = \frac{P(B \cap A)}{P(A)}$	Naive Bayes	$O(nm)$	$O(mc)$
	Principal Component Analysis (PCA)	$O(nm^2 + m^3)$	—
	t-SNE	$O(n^2m)$	—
	KMeans Clustering	$O(iknm)$??
n: samples m: dimensions n_{epoch}: epochs c: classes d_{tree}: depth n_{SV}: Support vectors k: clusters i: iterations			

But why even care about run time? There are multiple reasons why I always care about run time and why you should too.

To begin, we know that everyone is a big fan of sklearn implementations. It literally takes just two (max three) lines of code to run any ML algorithm with sklearn. Yet, in my experience, due to this simplicity, most users often overlook:

- The core understanding of an algorithm.
- The data-specific conditions that allow us to use an algorithm.

For instance, you'll be up for a big surprise if you use SVM or t-SNE on a dataset with plenty of samples.

- SVM's run-time grows cubically with the total number of samples.
- t-SNE's run-time grows quadratically with the total number of samples.

Another advantage of figuring out the run-time is that it helps us understand how an algorithm works end-to-end. Of course, in the above table, I have made some assumptions here and there. For instance:

- In a random forest, all decision trees may have different depths. But here, I have assumed that they are equal.
- During inference in kNN, we first find the distance to all data points. This gives a list of distances of size n (total samples).
 - Then, we find the k-smallest distances from this list.
 - The run-time to determine the k-smallest values may depend on the implementation.
 - Sorting and selecting the k-smallest values will be $O(n \log n)$.
 - But if we use a priority queue, it will take $O(n \log(k))$.
- In t-SNE, there's a learning step. Since the major run-time comes from computing the pairwise similarities in the high-dimensional space, we have ignored that step.

Nonetheless, the table still accurately reflects the general run-time of each of these algorithms.

As an exercise, I would encourage you to derive these run-time complexities yourself. This activity will provide you so much confidence in algorithmic understanding.

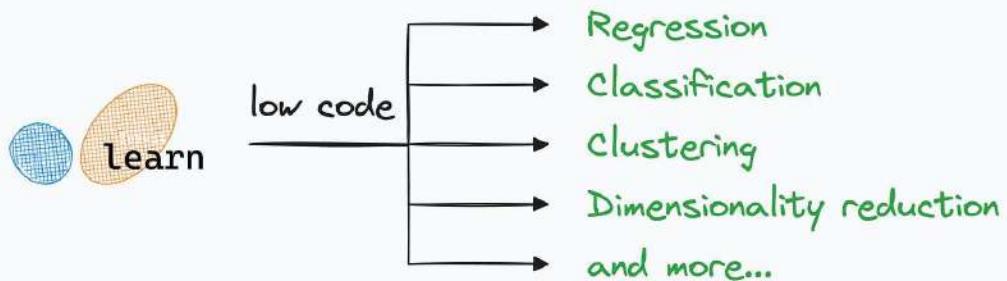
25 Most Important Mathematical Definitions in Data Science

"Is mathematical knowledge important in data science and machine learning?"

This is a question that so many people have, especially those who are just getting started.

Short answer: Yes, it's important, and here's why I say so.

See...these days, one can do "ML" without understanding any mathematical details of an algorithm. For instance (and thanks to sklearn, by the way):



- One can use any clustering algorithm in 2-3 lines of code.
- One can train classification models in 2-3 lines of code.
- And more.

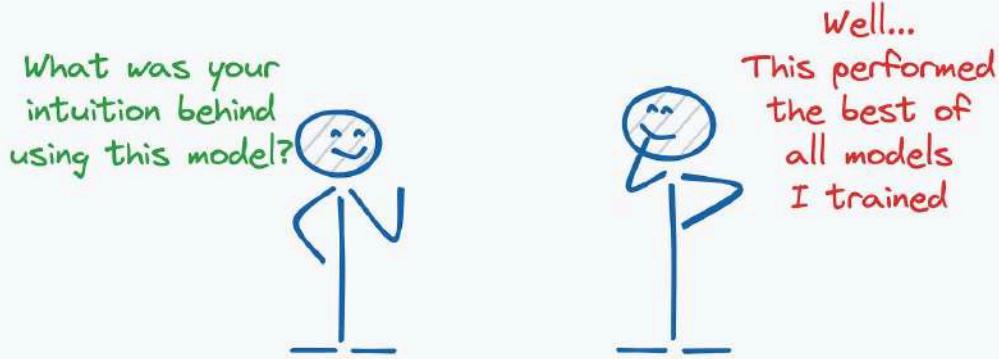
This is both good and bad:

- It's good because it saves us time.
- It's bad because this tempts us to ignore the underlying details.

In fact, I know many data scientists (mainly on the applied side) who do not entirely understand the mathematical details but can still build and deploy models.

Nothing wrong.

However, when I talk to them, I also see some disconnect between “What they are using” and “Why they are using it.”



Due to a lack of understanding of the underlying details:

- They find it quite difficult to optimize their models.
- They struggle to identify potential areas of improvement.
- They take a longer time to debug when things don't work well.
- They do not fully understand the role of specific hyperparameters.
- They use any algorithm without estimating their time complexity first.

If it feels like you are one of them, it's okay. This problem can be solved.

That said, if you genuinely aspire to excel in this field, building a curiosity for the underlying mathematical details holds exponential returns.

- Algorithmic awareness will give you confidence.
- It will decrease your time to build and iterate.

Gradually, you will go from a hit-and-trial approach to “I know what should work.”

To help you take that first step, I prepared the following visual, which lists some of the most important mathematical formulations used in Data Science and Statistics (in no specific order).

Before reading ahead, look at them one by one and calculate how many of them do you already know:

25 Most important Mathematical Definitions in Data Science

blog.DailyDoseofDS.com

1) Gradient Descent $\theta_{j+1} = \theta_j - \alpha \nabla J(\theta_j)$	2) Normal distribution $f(x \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$	3) Z-score $z = \frac{x-\mu}{\sigma}$	
4) Sigmoid $\sigma(x) = \frac{1}{1 + e^{-x}}$	5) Correlation $\text{Correlation} = \frac{\text{Cov}(X, Y)}{\text{Std}(X) \cdot \text{Std}(Y)}$	6) Cosine Similarity $\text{similarity} = \frac{A \cdot B}{\ A\ \ B\ }$	
7) Naive Bayes $P(y x_1, \dots, x_n) = \frac{P(y) \prod_{i=1}^n P(x_i y)}{P(x_1, \dots, x_n)}$	8) MLE $\operatorname{argmax}_{\theta} \prod_{i=1}^n P(x_i \theta)$	9) OLS $\hat{\beta} = (X^T X)^{-1} X^T y$	
10) F1 Score $\frac{2 \cdot P \cdot R}{P + R}$	11) ReLU $\max(0, x)$	12) Softmax $P(y=j x) = \frac{e^{x^T w_j}}{\sum_{k=1}^K e^{x^T w_k}}$	13) R2 score $R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$
14) MSE $\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$	15) MSE + L2 Reg $\text{MSE}_{\text{regularized}} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p \beta_j^2$	16) Eigen vectors $Av = \lambda v$	
17) Entropy $\text{Entropy} = - \sum_i p_i \log_2(p_i)$	18) KMeans $\operatorname{argmin}_S \sum_{i=1}^k \sum_{x \in S_i} \ x - \mu_i\ ^2$	19) KL Divergence $D_{\text{KL}}(P Q) = \sum_{x \in \mathcal{X}} P(x) \log\left(\frac{P(x)}{Q(x)}\right)$	
20) Log-loss $-\frac{1}{N} \sum_{i=1}^N (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i))$	21) SVM $\min_{w,b} \frac{1}{2} \ w\ ^2 + C \sum_{i=1}^n \max(0, 1 - y_i(w \cdot x_i - b))$		
22) Linear regression $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$	23) SVD $A = U \Sigma V^T$		
24) Lagrange multiplier $\max f(x); g(x) = 0$ $L(x, \lambda) = f(x) - \lambda * g(x)$	25) What will you add? • • •		

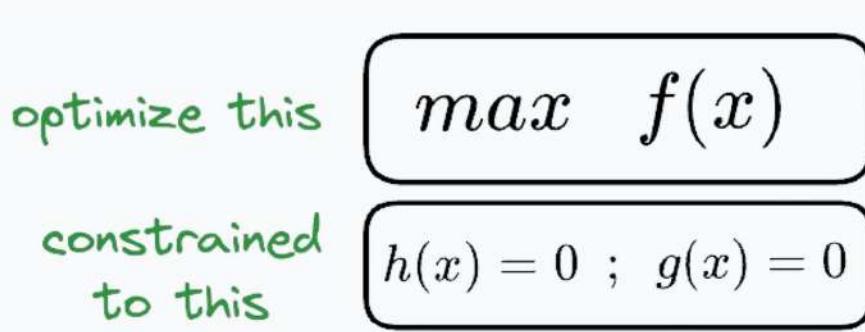
Some of the terms are pretty self-explanatory, so I won't go through each of them, like:

- Gradient Descent, Normal Distribution, Sigmoid, Correlation, Cosine similarity, Naive Bayes, F1 score, ReLU, Softmax, MSE, MSE + L2 regularization, KMeans, Linear regression, SVM, Log loss.

Here are the remaining terms:

- MLE (Maximum Likelihood Estimation): A method for estimating the parameters of a statistical model by maximizing the likelihood of the observed data. We covered it in the previous chapter.

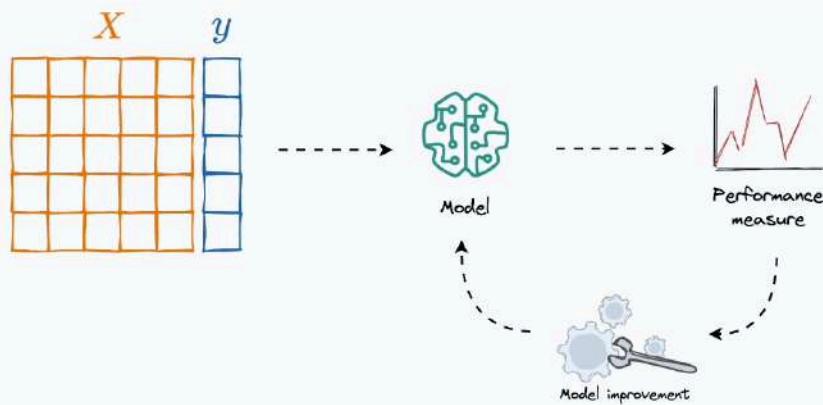
- Z-score: A standardized value that indicates how many standard deviations a data point is from the mean.
- Ordinary Least Squares: A closed-form solution for linear regression obtained using the MLE step mentioned above.
- Entropy: A measure of the uncertainty of a random variable.
- Eigen Vectors: The non-zero vectors that do not change their direction when a linear transformation is applied. It is widely used in dimensionality reduction techniques like PCA.
- R2 (R-squared): A statistical measure that represents the proportion of variance explained by a regression model.
- KL divergence: Assess how much information is lost when one distribution is used to approximate another distribution. It is used as a loss function in the t-SNE algorithm.
- SVD: A factorization technique that decomposes a matrix into three other matrices, often noted as U, Σ , and V. It is fundamental in linear algebra for applications like dimensionality reduction, noise reduction, and data compression.
- Lagrange multipliers: They are commonly used mathematical techniques to solve constrained optimization problems. For instance, consider an optimization problem with an objective function $f(x)$ and assume that the constraints are $g(x)=0$ and $h(x)=0$. Lagrange multipliers solve this.



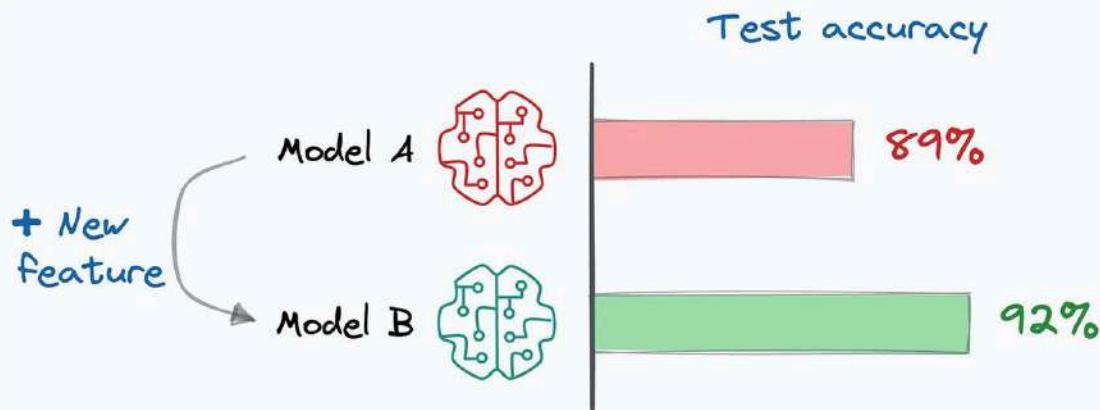
How to Reliably Improve Probabilistic Multiclass-classification Models

ML model building is typically an iterative process. Given some dataset:

- We train a model.
- We evaluate it.
- And we continue to improve it until we are satisfied with the performance.

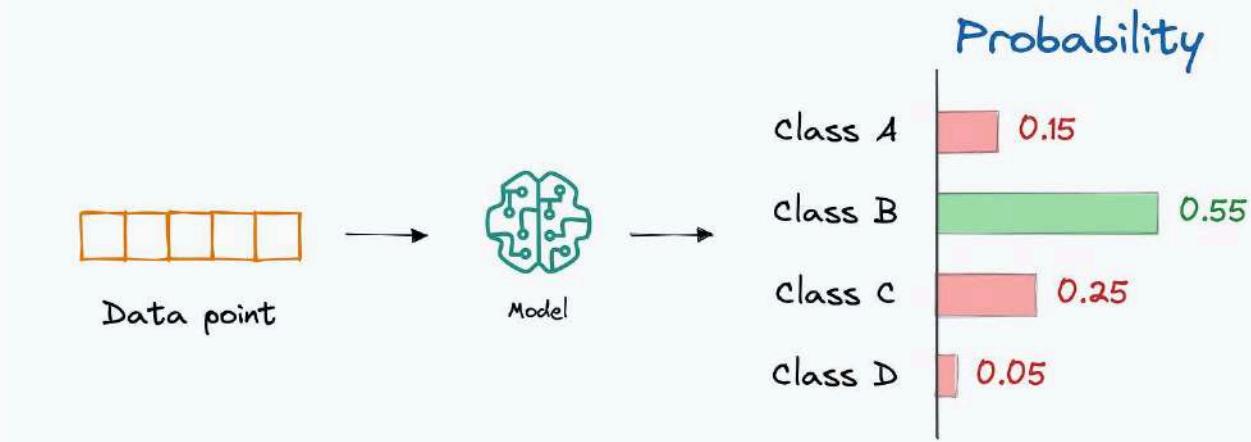


Here, the efficacy of any model improvement strategy (say, introducing a new feature) is determined using some sort of performance metric.



However, I have often observed that when improving probabilistic multiclass-classification models, this technique can be a bit deceptive when the efficacy is determined using “Accuracy.”

Probabilistic multiclass-classification models are those models that output probabilities corresponding to each class, like neural networks.

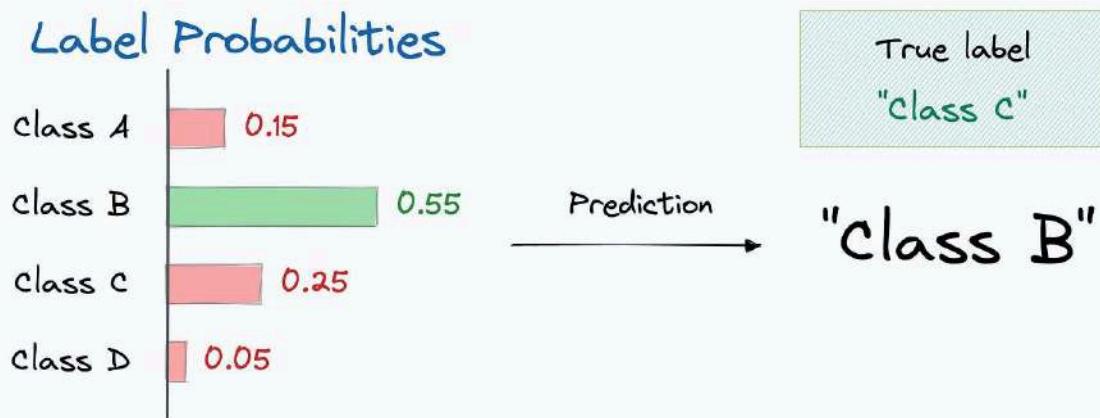


In other words, it is possible that we are actually making good progress in improving the model, but “Accuracy” is not reflecting that (yet).

Let's understand.

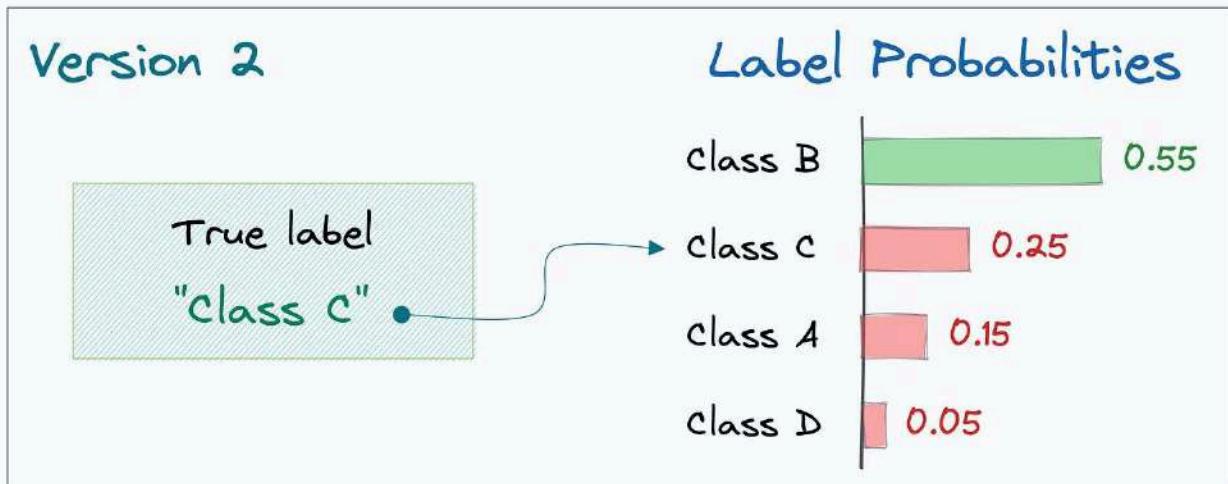
Pitfall of Accuracy

In probabilistic multiclass-classification models, Accuracy is determined using the output label that has the highest probability:

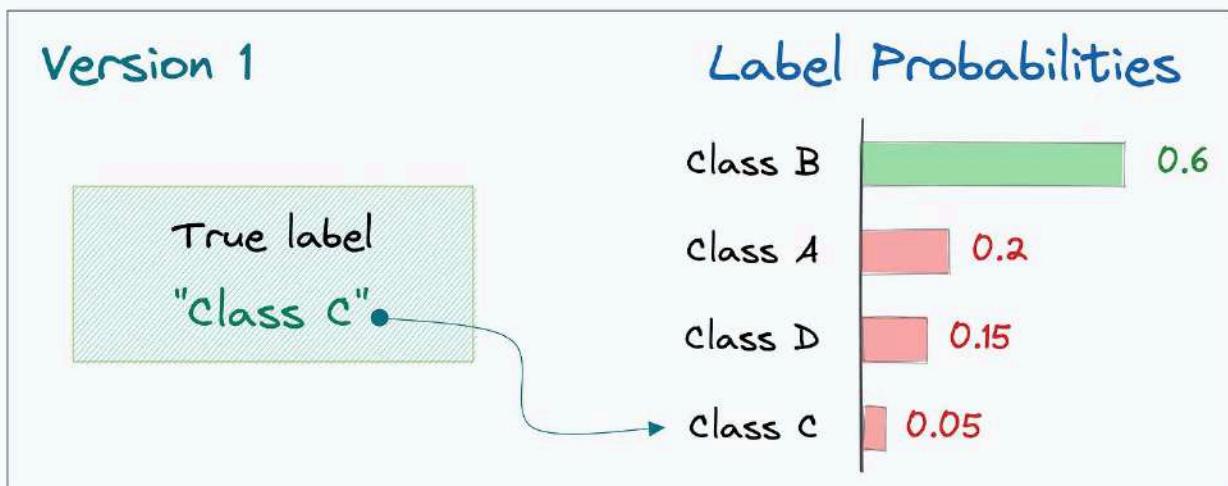


Now, it's possible that the actual label is not predicted with the highest probability by the model, but it's in the top “k” output labels.

For instance, in the image below, the actual label (Class C) is not the highest probability label, but it's at least in the top 2 predicted probabilities (Class B and Class C):



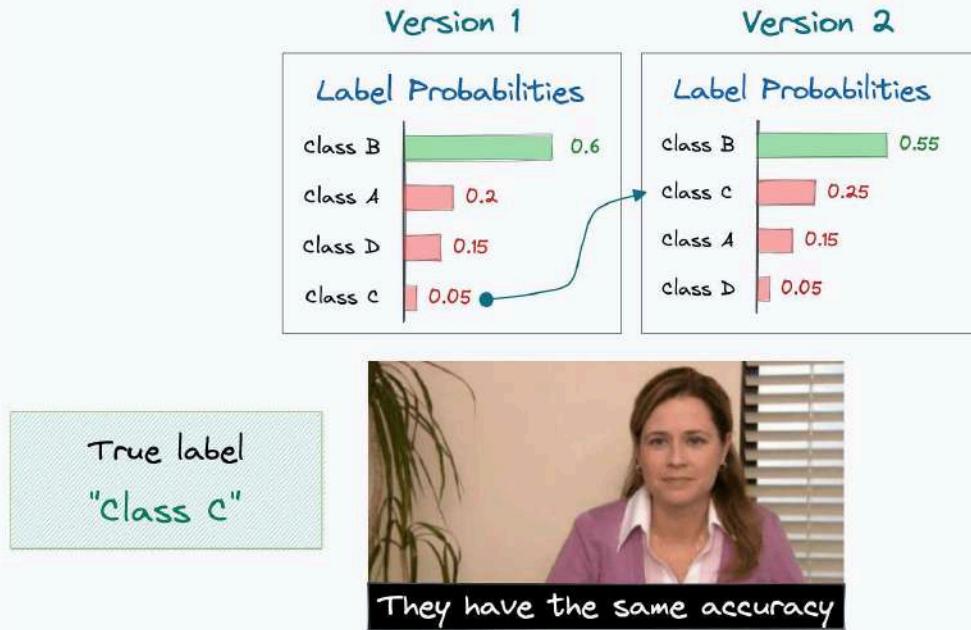
And what if in an earlier version of our model, the output probability of Class C was the lowest, as depicted below:



Now, of course, in both cases, the final prediction is incorrect.

However, while iterating from “Version 1” to “Version 2” using some model improvement techniques, we genuinely made good progress.

Nonetheless, Accuracy entirely discards this as it only cares about the highest probability label.



I hope you understand the problem here.

Solution

Whenever I am building and iteratively improving any probabilistic multiclass classification model, I always use the top-k accuracy score. As the name suggests, it computes whether the correct label is among the top k labels predicted probabilities or not.

As you may have already guessed, top-1 accuracy score is the traditional Accuracy score. This is a much better indicator to assess whether my model improvement efforts are translating into meaningful enhancements in predictive performance or not.

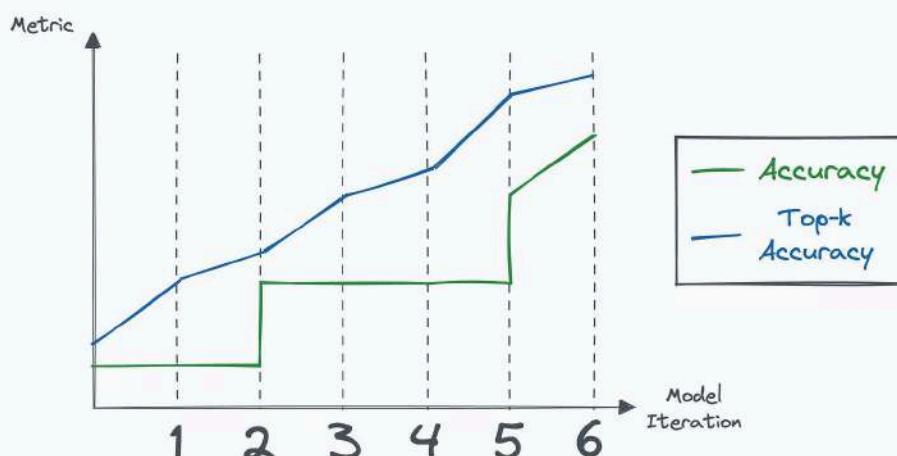
For instance, if the top-3 accuracy score goes from 75% to 90%, this totally suggests that whatever we did to improve the model was effective:

- Earlier, the correct prediction was in the top 3 labels only 75% of the time.
- But now, the correct prediction is in the top 3 labels 90% of the time.

As a result, one can effectively redirect their engineering efforts in the right direction. Of course, what I am saying should only be used to assess the model improvement efforts.

This is because true predictive power will inevitably be determined using traditional model accuracy. So make sure you are gradually progressing on the Accuracy front too.

Ideally, it is expected that “Top-k Accuracy” may continue to increase during model iterations, which reflects improvement in performance. Accuracy, however, may stay the same for a while, as depicted below:



Top-k accuracy score is also available in Sklearn:

sklearn.metrics.top_k_accuracy_score

```
sklearn.metrics.top_k_accuracy_score(y_true, y_score, *, k=2, normalize=True, sample_weight=None, labels=None)
\[source\]
```

Top-k Accuracy classification score.

This metric computes the number of times where the correct label is among the top `k` labels predicted (ranked by predicted scores). Note that the multilabel case isn't covered here.

Read more in the [User Guide](#)

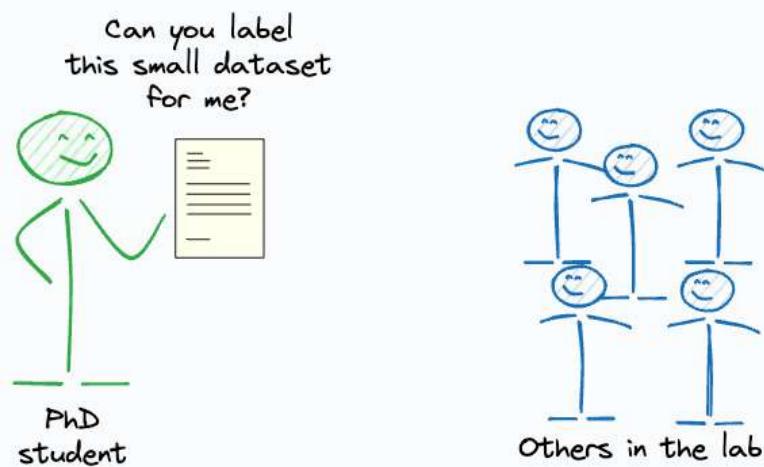
Parameters: `y_true` : *array-like of shape (n_samples,)*
True labels.

`y_score` : *array-like of shape (n_samples,) or (n_samples, n_classes)*
Target scores. These can be either probability estimates or non-thresholded decision values (as returned

Your Entire Model Improvement Efforts Might Be Going in Vain

Back in 2019, I was working with an ML research group in Germany.

One day, a Ph.D. student came up to me (and others in the lab), handed over a small sample of the dataset he was working with, and requested us to label it, despite having true labels.

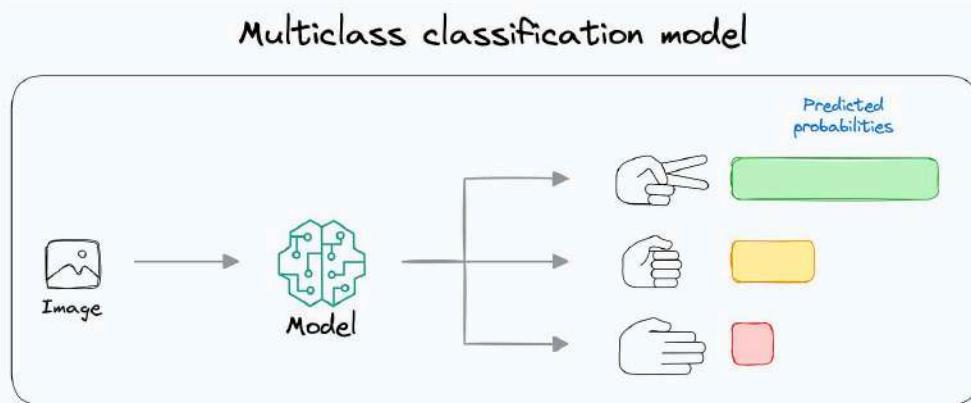


This made me curious about why gathering human labels was necessary for him when he already had ground truth labels available. So I asked.

What I learned that day changed my approach to incremental model improvement, and I am sure you will find this idea fascinating too.

Let me explain what I learned.

Consider we are building a multiclass classification model. Say it's a model that classifies an input image as a rock, paper, or scissors:



For simplicity, let's assume there's no class imbalance. Calculating the class-wise validation accuracies gives us the following results:

Class-wise accuracies

Class	Accuracy
Paper	97%
Rock	82%
Scissor	75%

Question: Which class would you most intuitively proceed to inspect further and improve the model on?

After looking at these results, most people believe that “Scissor” is the worst-performing class and should be inspected further.

Scissor class is giving bad results. I should inspect it.



Class	Accuracy
Paper	97%
Rock	82%
Scissor	75%

But this might not be true. And this is precisely what that Ph.D. student wanted to verify by collecting human labels. Let's say that the human labels give us the following results:

Class-wise human accuracies

Class	Human Accuracy	Model Accuracy
Paper	98%	97% ← -1%
Rock	95%	82% ← -13%
Scissor	77%	75% ← -2%

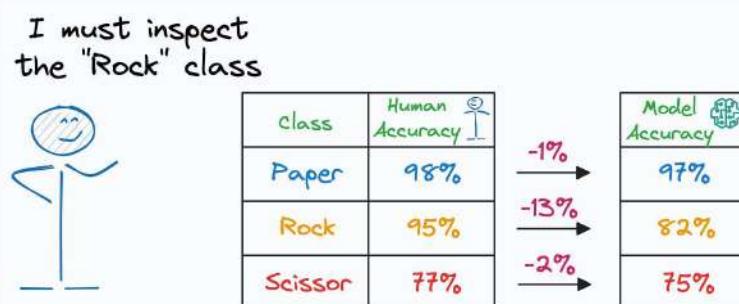
Based on this, do you still think the model performs the worst on the “Scissor” class?

No, right?

I mean, of course, the model has the least accuracy on the “Scissor” class, and I am not denying it. However, with more context, we notice that the model is doing a pretty good job classifying the “Scissor” class. This is because an average human is achieving just 2% higher accuracy in comparison to what our model is able to achieve.

However, the above results astonishingly reveal that it is the “Rock” class instead that demands more attention. The accuracy difference between an average human and the model is way too high (13%). Had we not known this, we would have continued to improve the “Scissor” class, when in reality, “Rock” requires more improvement.

Ever since I learned this technique, I have found it super helpful to determine my next steps for model improvement, if possible. I say “if possible” because I understand that many datasets are hard for humans to interpret and label. Nonetheless, if it is feasible to set up such a “human baseline,” one can get so much clarity into how the model is performing.



As a result, one can effectively redirect their engineering efforts in the right direction.

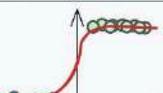
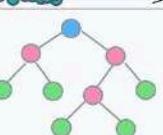
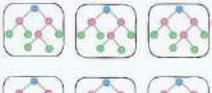
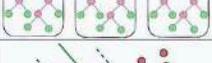
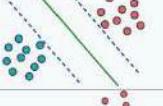
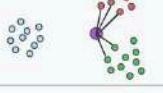
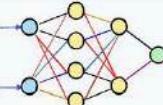
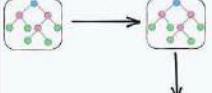
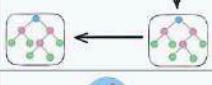
Of course, I am not claiming that this will be universally useful in all use cases.

For instance, if the model is already performing better than the baseline, the model improvements from there on will have to be guided based on past results.

Yet, in such cases, surpassing a human baseline at least helps us validate that the model is doing better than what a human can do.

Loss Function of 16 ML Algos

The below visual depicts the most commonly used loss functions by various ML algorithms.

Loss functions of 16 Most Popular ML Algorithms		 blog.DailyDoseofDS.com
	Linear Regression	Mean Squared Error
	Logistic Regression	Cross-Entropy Loss
	Decision Tree Classifier	Information Gain or Gini impurity
	Decision Tree Regressor	Mean Squared Error
	Random Forest Classifier	Information Gain or Gini impurity
	Random Forest Regressor	Mean Squared Error
	Support Vector Machines (SVMs)	Hinge Loss
	k-Nearest Neighbors	No loss function
$P(B A) = \frac{P(B \cap A)}{P(A)}$	Naive Bayes	No loss function
	Neural Networks	Regression: Mean Squared Error Classification: Cross-Entropy Loss
	AdaBoost	Exponential loss
	Gradient Boosting LightGBM CatBoost XGBoost	Regression: Mean Squared Error Classification: Cross-Entropy Loss
	KMeans Clustering	??

*These are typically used loss functions. It does not mean they are used always in these algorithms.

10 Most Common Loss Function

The below visual depicts some commonly used loss functions in regression and classification tasks.

10 Most Common Loss Functions in Machine Learning



blog.DailyDoseofDS.com

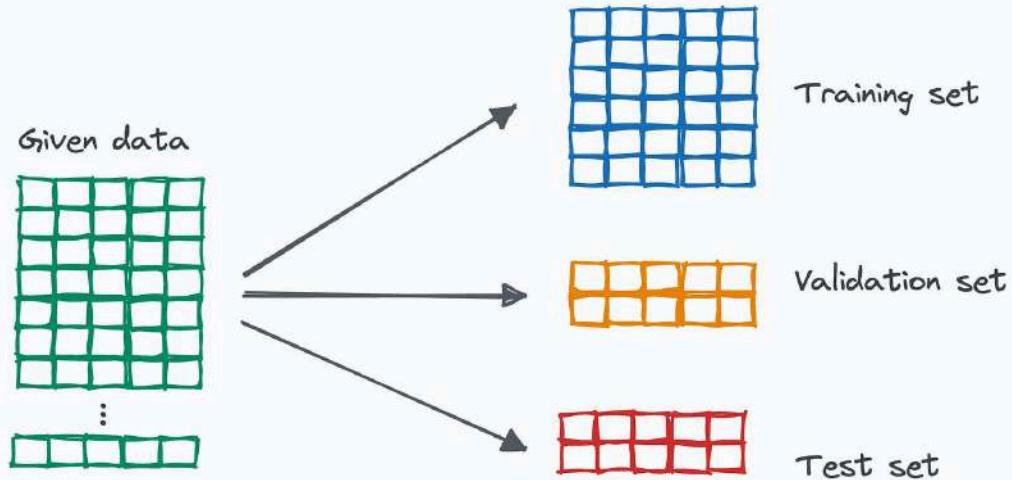
Loss Function Name	Description	Function
Regression Loss Functions		
Mean Bias Error	Captures average bias in prediction. But is rarely used for training.	$\mathcal{L}_{MBE} = \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))$
Mean Absolute Error	Measures absolute average bias in prediction. Also called L1 Loss.	$\mathcal{L}_{MAE} = \frac{1}{N} \sum_{i=1}^N y_i - f(x_i) $
Mean Squared Error	Average squared distance between actual and predicted. Also called L2 Loss.	$\mathcal{L}_{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))^2$
Root Mean Squared Error	Square root of MSE. Loss and dependent variable have same units.	$\mathcal{L}_{RMSE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))^2}$
Huber Loss	A combination of MSE and MAE. It is parametric loss function.	$\mathcal{L}_{\text{Huberloss}} = \begin{cases} \frac{1}{2}(y_i - f(x_i))^2 & : y_i - f(x_i) \leq \delta \\ \delta(y_i - f(x_i) - \frac{1}{2}\delta) & : \text{otherwise} \end{cases}$
Log Cosh Loss	Similar to Huber Loss + non-parametric. But computationally expensive.	$\mathcal{L}_{\text{LogCosh}} = \frac{1}{N} \sum_{i=1}^N \log(\cosh(f(x_i) - y_i))$

Classification Loss Functions (Binary + Multi-class)

Binary Cross Entropy (BCE)	Loss function for binary classification tasks.	$\mathcal{L}_{BCE} = \frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(x_i)) + (1 - y_i) \cdot \log(1 - p(x_i))$
Hinge Loss	Penalizes wrong and right (but less confident) predictions. Commonly used in SVMs.	$\mathcal{L}_{\text{Hinge}} = \max(0, 1 - (f(x) \cdot y))$
Cross Entropy Loss	Extension of BCE loss to multi-class classification.	$\mathcal{L}_{CE} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} \log(f(x_{ij}))$ N : samples; M : classes
KL Divergence	Minimizes the divergence between predicted and true probability distribution	$\mathcal{L}_{KL} = \sum_{i=1}^N y_i \cdot \log\left(\frac{y_i}{f(x_i)}\right)$

How to Actually Use Train, Validation and Test Set

It is pretty conventional to split the given data into train, test, and validation sets.



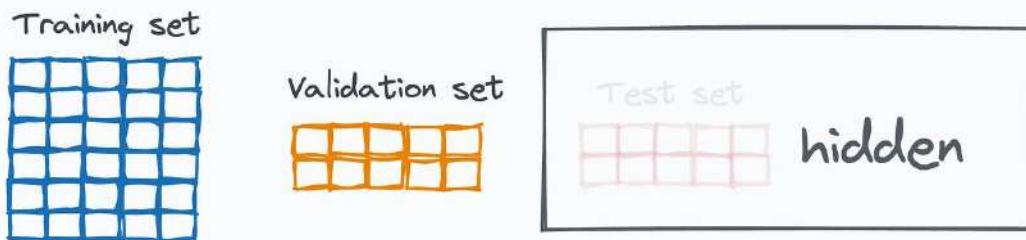
However, there are quite a few misconceptions about how they are meant to be used, especially the validation and test sets.

In this chapter, let's clear them up and see how to truly use train, validation, and test sets.

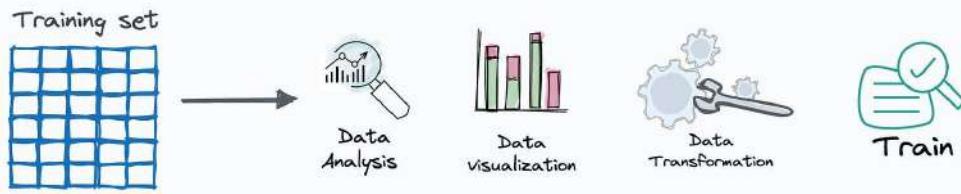
We begin by splitting the data into:

- Train
- Validation
- Test

At this point, just assume that the test data does not even exist. Forget about it instantly.



Begin with the train set. This is your whole world now.



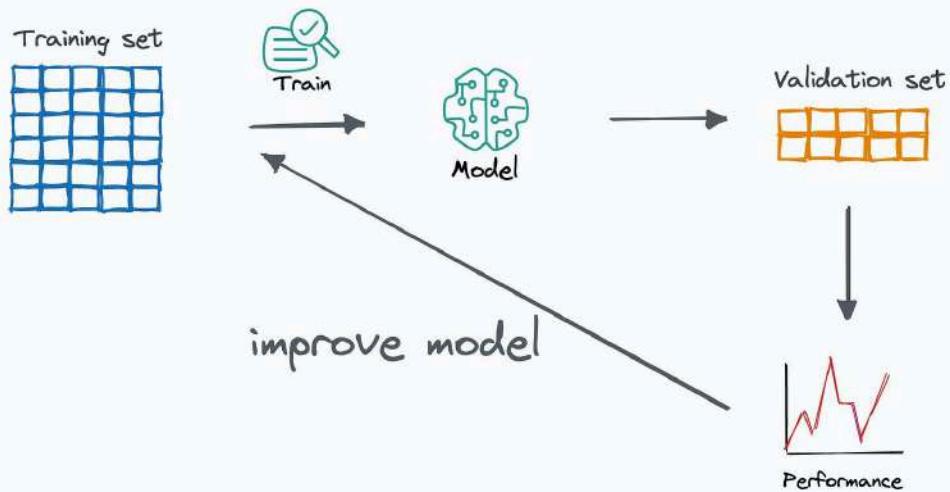
- You analyze it
- You transform it
- You use it to determine features
- You fit a model on it

After modeling, you would want to measure the model's performance on unseen data, wouldn't you?

Bring in the validation set now.



Based on validation performance, improve the model. Here's how you iteratively build your model:



- Train using a train set
- Evaluate it using the validation set
- Improve the model
- Evaluate again using the validation set
- Improve the model again
- and so on.

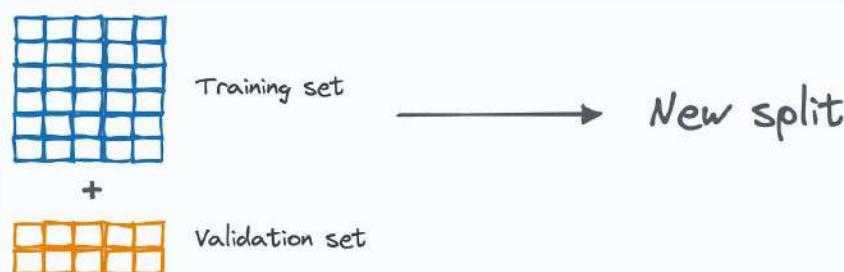
Until...

You reach a point where you start overfitting the validation set.

This indicates that you have exploited (or polluted) the validation set.

No worries.

Merge it with the train set and generate a new split of train and validation.



Note: Rely on cross-validation if needed, especially when you don't have much data. You may still use cross-validation if you have enough data. But it can be computationally intensive.

Now, if you are happy with the model's performance, evaluate it on test data.

✓ What you use a test set for:

- Get a final and unbiased review of the model.

✗ What you DON'T use a test set for:

- Analysis, decision-making, etc.

If the model is underperforming on the test set, no problem. Go back to the modeling stage and improve it.

BUT (and here's what most people do wrong)!

They use the same test set again. This is not allowed!

Think of it this way.

Your professor taught you in the classroom. All in-class lessons and examples are the train set.

The professor gave you take-home assignments, which acted like validation sets.

You got some wrong and some right. Based on this, you adjusted your topic fundamentals, i.e., improved the model.

Now, if you keep solving the same take-home assignment repeatedly, you will eventually overfit it, won't you?

That is why we bring in a new validation set after some iterations.

The final exam day paper is your test set. If you do well, awesome!

But if you fail, the professor cannot give you the exact exam paper next time, can they? This is because you know what's inside.

Of course, by evaluating a model on the test set, the model never gets to "know" the precise examples inside that set. But the issue is that the test set has been exposed now.

Your previous evaluation will inevitably influence any further evaluations on that specific test set. That is why you must always use a specific test set only ONCE.

Once you do, merge it with the train and validation set and generate an entirely new split.

Repeat.

And that is how you use train, validation, and test sets in machine learning.

5 Cross Validation Techniques

Tuning and validating machine learning models on a single validation set can be misleading and sometimes yield overly optimistic results.

This can occur due to a lucky random split of data, which results in a model that performs exceptionally well on the validation set but poorly on new, unseen data.

That is why we often use cross validation instead of simple single-set validation.

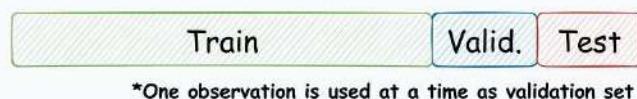
Cross validation involves repeatedly partitioning the available data into subsets, training the model on a few subsets, and validating on the remaining subsets.

The main advantage of cross validation is that it provides a more robust and unbiased estimate of model performance compared to the traditional validation method.

Below are five of the most commonly used and must-know cross validation techniques.

Leave-One-Out Cross Validation

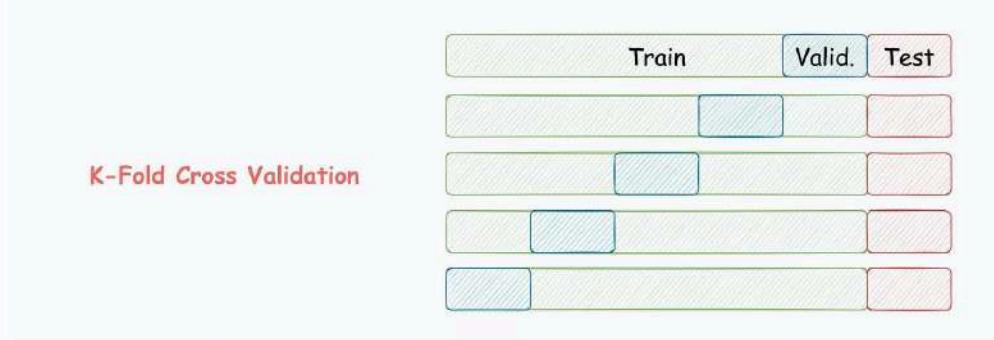
Leave One Out Cross Validation



*One observation is used at a time as validation set

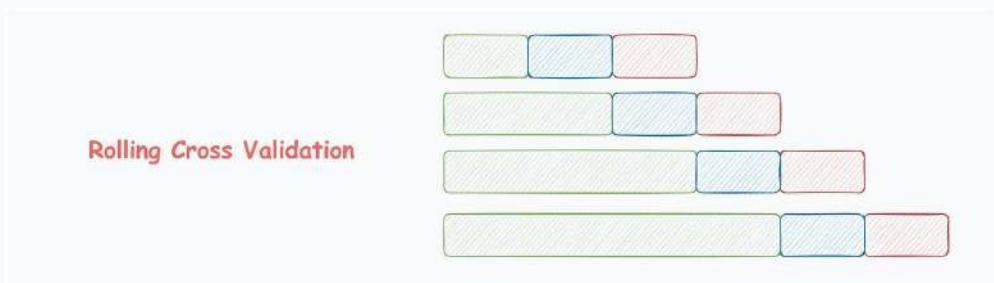
- Leave one data point for validation.
- Train the model on the remaining data points.
- Repeat for all points.
- Of course, as you may have guessed, this is practically infeasible when you have many data points. This is because number of models is equal to number of data points.
- We can extend this to Leave-p-Out Cross Validation, where, in each iteration, p observations are reserved for validation, and the rest are used for training.

K-Fold Cross Validation



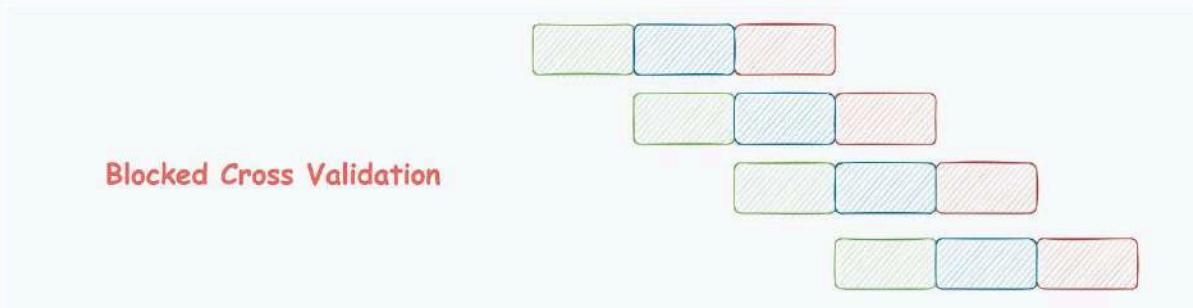
- Split data into k equally-sized subsets.
- Select one subset for validation.
- Train the model on the remaining subsets.
- Repeat for all subsets.

Rolling Cross Validation



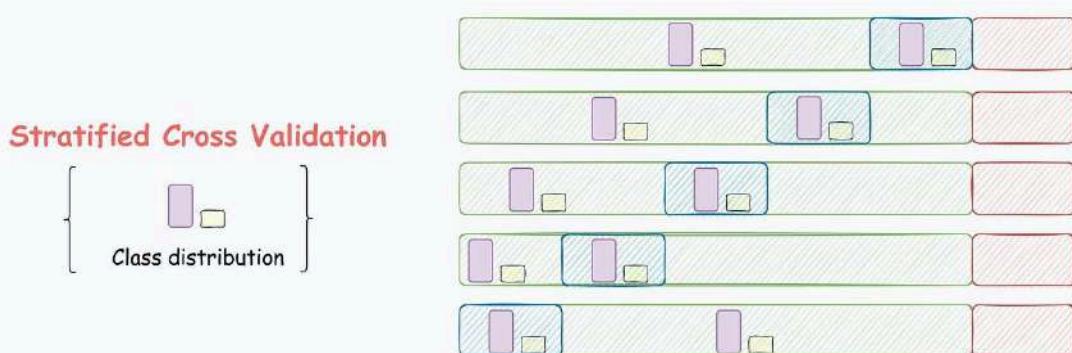
- Mostly used for data with temporal structure.
- Data splitting respects the temporal order, using a fixed-size training window.
- The model is evaluated on the subsequent window.

Blocked Cross Validation



- Another common technique for time-series data.
- In contrast to rolling cross validation, the slice of data is intentionally kept short if the variance does not change appreciably from one window to the next.
- This also saves computation over rolling cross validation.

Stratified Cross Validation



- The above-discussed techniques may not work for imbalanced datasets. Stratified cross validation is mainly used for preserving the class distribution.
- The partitioning ensures that the class distribution is preserved.

Let's continue our discussion on cross validation in the next chapter.

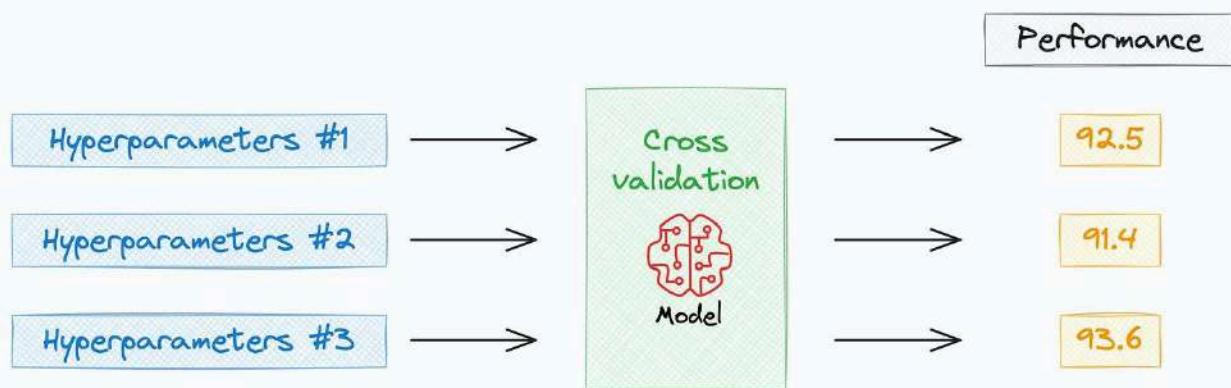
What To Do After Cross Validation?

Let me ask you a question.

But before I do that, I need to borrow your imagination for just a moment.

Imagine you are building some supervised machine learning model. You are using cross-validation to determine an optimal set of hyperparameters.

Essentially, every hyperparameter configuration corresponds to a cross-validation performance:



After obtaining the best hyperparameters, we need to finalize a model (say for production); otherwise, what is the point of all this hassle? Now, here's the question:

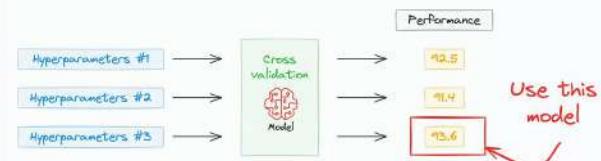
After obtaining the optimal hyperparameters, what would you be more inclined to do:

- 1) Retrain the model again on the entire data (train + validation + test) with the optimal hyperparameters?
If we do this, remember that we can't reliably validate this new model as there is no unseen data left.



2) Just proceed with the best-performing model based on cross-validation performance itself. If we do this, remember that we are leaving out important data, which we could have trained our model with.

What would you do?



Recommended path

My strong preference has almost always been “retraining a new model with entire data.”



There are, of course, some considerations to keep in mind, which I have learned through the models I have built and deployed. That said, in most cases, retraining is the ideal way to proceed.

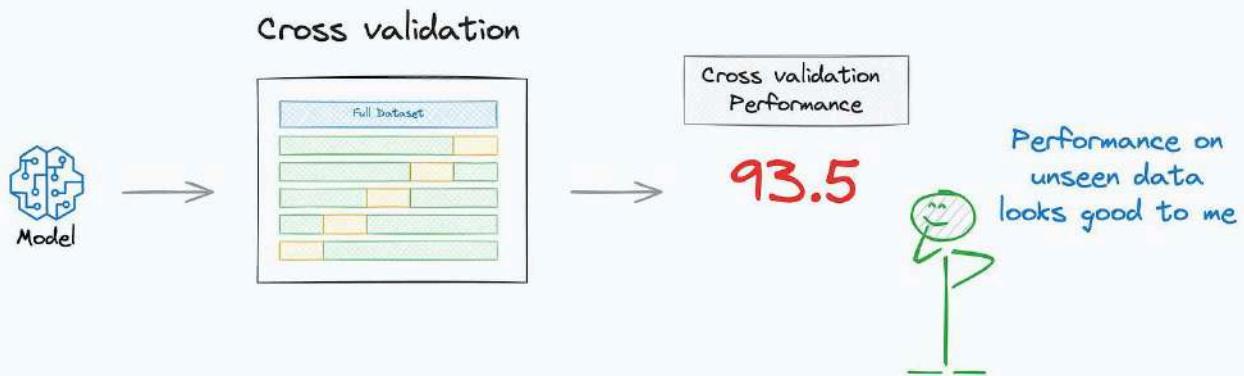
Let me explain.

Why retrain the model?

We would want to retrain a new model because, in a way, we are already satisfied with the cross-validation performance, which, by its very nature, is an out-of-fold metric.

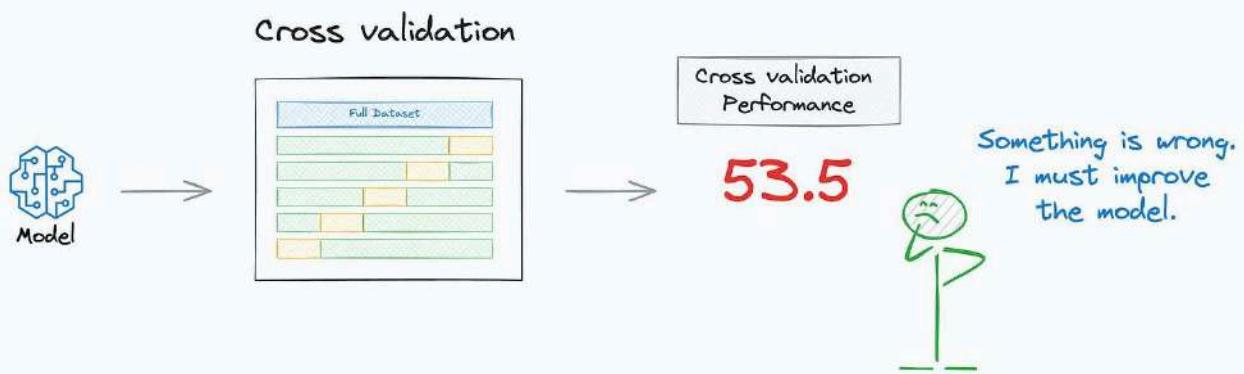
An out-of-fold data is data that has not been seen by the model during the training. An out-of-fold metric is the performance on that data.

In other words, we already believe that the model aligns with how we expect it to perform on unseen data.



Thus, incorporating this unseen validation set in the training data and retraining the model will MOST LIKELY have NO effect on its performance on unseen data after deployment (assuming a sudden covariate shift hasn't kicked in, which is a different issue altogether).

If, however, we were not satisfied with the cross-validation performance itself, we wouldn't even be thinking about finalizing a model in the first place.

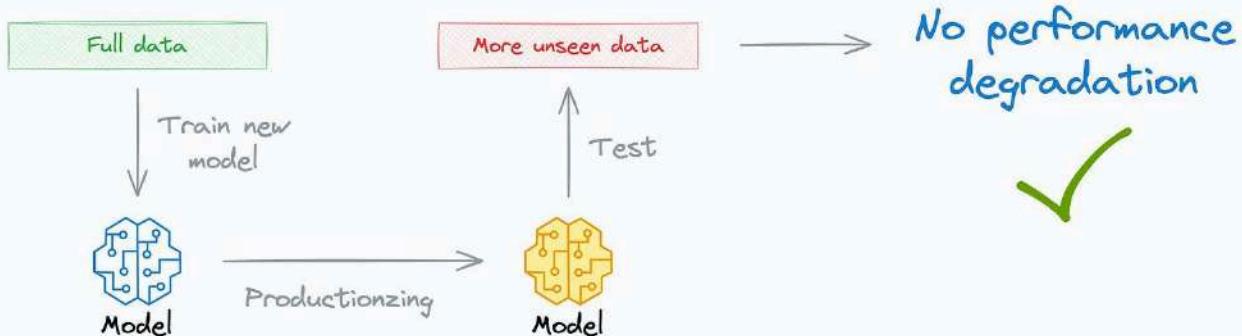


Instead, we would be thinking about ways to improve the model by working on feature engineering, trying new hyperparameters, experimenting with different models, and more.

The reasoning makes intuitive sense as well.

It's hard for me to recall any instance where retraining did something disastrously bad to the overall model.

In fact, I vividly remember one instance wherein, while I was productionizing the model (it took me a couple of days after retraining), the team had gathered some more labeled data.

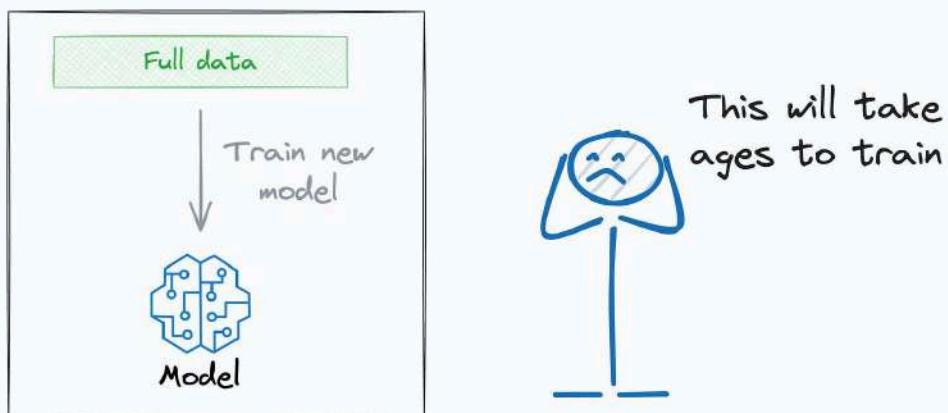


The model didn't show any performance degradation when I evaluated it (just to double-check). As an added benefit, this also helped ensure that I had made no errors while productionizing my model.

Some considerations

Here, please note that it's not a rule that you must always retrain a new model.

The field itself and the tasks one can solve are pretty diverse, so one must be open-minded while solving the problem at hand. One of the reasons I wouldn't want to retrain a new model is that it takes days or weeks to train the model.



In fact, even if we retrain a new model, there are MANY business situations in which stakes are just too high.

Thus, one can never afford to be negligent about deploying a model without re-evaluating it — transactional fraud, for instance.

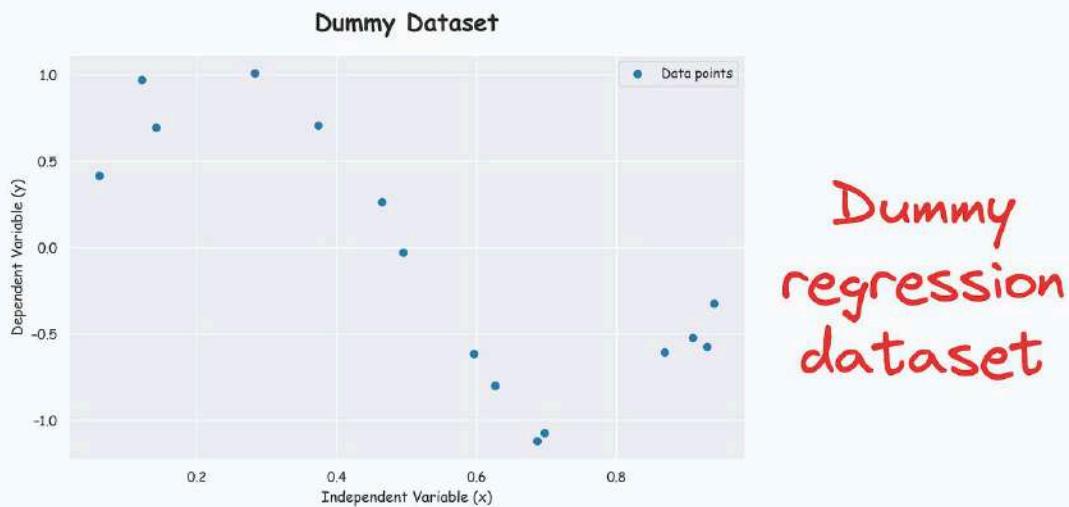
In such cases, I have seen that while a team works on productionizing the model, data engineers gather some more data in the meantime.

Before deploying, the team would do some final checks on that dataset.

The newly gathered data is then considered in the subsequent iterations of model improvements.

Double Descent vs. Bias-Variance Trade-off

It is well-known that as the number of model parameters increases, we typically overfit the data more and more. For instance, consider fitting a polynomial regression model trained on this dummy dataset below:



In case you don't know, this is called a polynomial regression model:

Polynomial
regression model

$$y = \theta_0 + \theta_1 * x^1 + \theta_2 * x^2 + \dots + \theta_m * x^m$$

It is expected that as we'll increase the degree (m) and train the polynomial regression model:

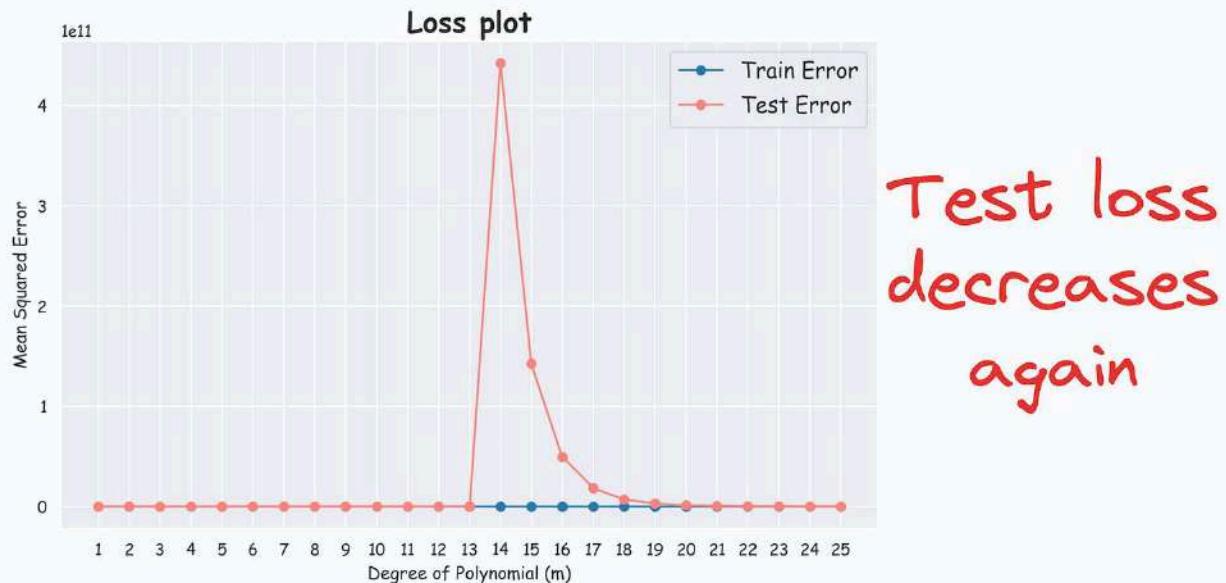
- The training loss will get closer and closer to zero.
- The test (or validation) loss will first reduce and then get bigger and bigger.

This is because, with a higher degree, the model will find it easier to contort its regression fit through each training data point, which makes sense.

In fact, this is also evident from the following loss plot:



But notice what happens when we continue to increase the degree (m):



That's strange, right?

Why does the test loss increase to a certain point but then decrease?

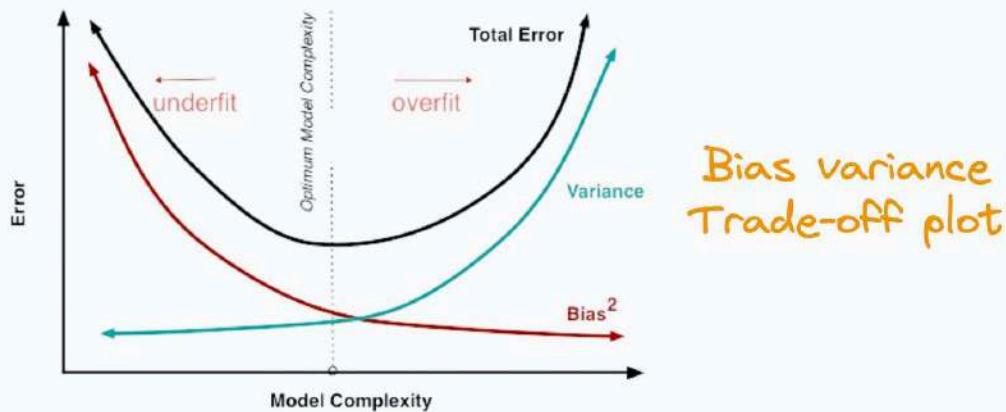
This was not expected, was it?

Well...what you are seeing is called the “double descent phenomenon,” which is quite commonly observed in many ML models, especially deep learning models.

It shows that, counterintuitively, increasing the model complexity beyond the point of interpolation can improve generalization performance.

In fact, this whole idea is deeply rooted to why LLMs, although massively big (billions or even trillions of parameters), can still generalize pretty well.

And it’s hard to accept it because this phenomenon directly challenges the traditional bias-variance trade-off we learn in any introductory ML class:



Putting it another way, training very large models, even with more parameters than training data points, can still generalize well.

To the best of my knowledge, this is still an open question, and it isn’t entirely clear why neural networks exhibit this behavior.

There are some theories around regularization, however, such as this one:

It could be that the model applies some sort of implicit regularization, with which, it can precisely focus on an apt number of parameters for generalization.

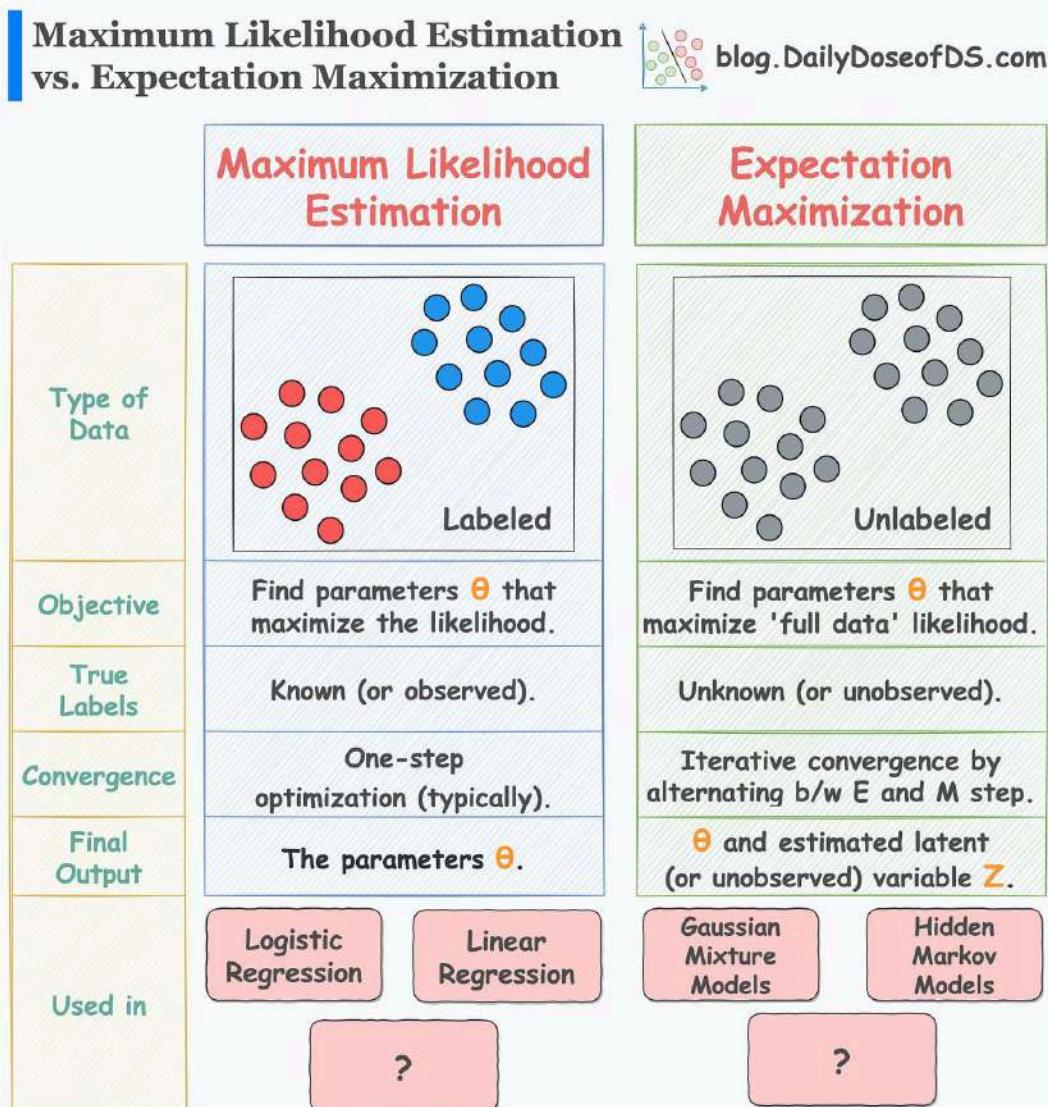
But to be honest, nothing is clear yet.

Statistical Foundations

MLE vs. EM – What's the Difference?

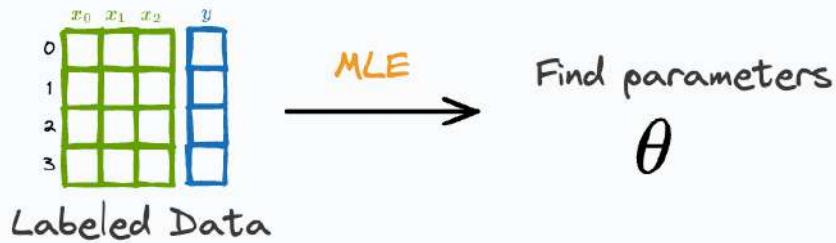
Maximum likelihood estimation (MLE) and expectation maximization (EM) are two popular techniques to determine the parameters of statistical models.

Due to its applicability in MANY statistical models, I have seen it being asked in plenty of data science interviews as well, especially the distinction between the two. The following visual summarizes how they work:



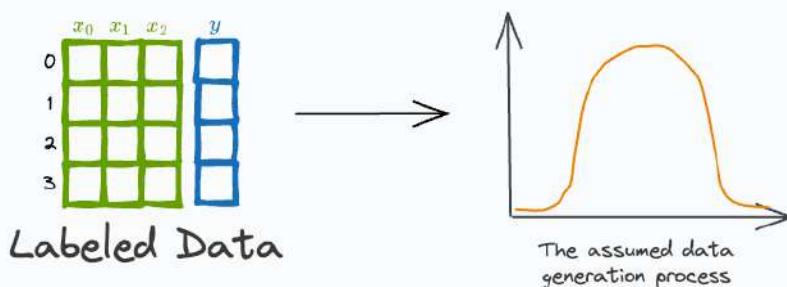
Maximum likelihood estimation (MLE)

MLE starts with a labeled dataset and aims to determine the parameters of the statistical model we are trying to fit.



The process is pretty simple and straightforward. In MLE, we:

- Start by assuming a data generation process. Simply put, this data generation process reflects our belief about the distribution of the output label (y), given the input (X).



- Next, we define the likelihood of observing the data. As each observation is independent, the likelihood of observing the entire data is the same as the product of observing individual observations:

From the assumed data generation process

Labeled Data

Likelihood function

$$\xrightarrow{\quad} L(\theta) = \prod P(y|X; \theta)$$

- The likelihood function above depends on parameter values (θ). Our objective is to determine those specific parameter values that maximize the likelihood function. We do this as follows:

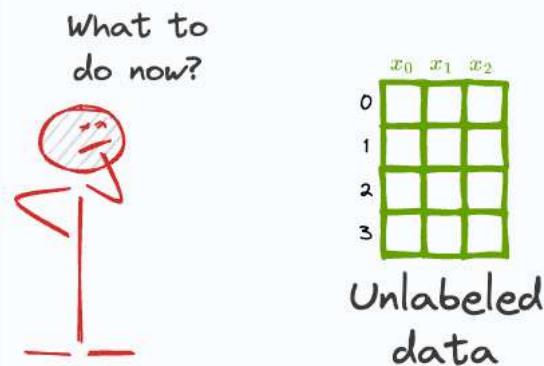
$$\frac{\delta L(\theta)}{\delta \theta} = 0 \longrightarrow \text{Solve for } \theta$$

This gives our parameter estimates that would have most likely generated the given data.

That was pretty simple, wasn't it?

But what do we do if we don't have true labels?

We still want to estimate the parameters, don't we?

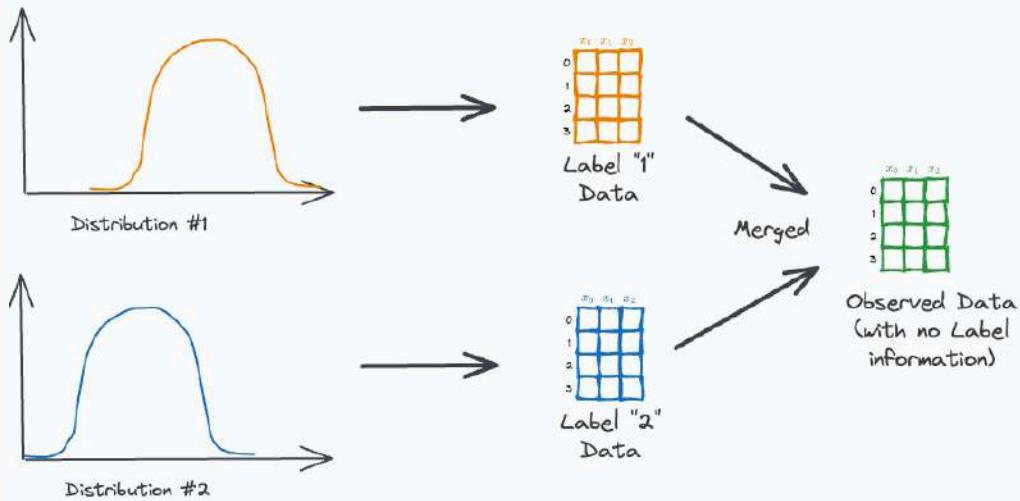


MLE, as you may have guessed, will not be applicable. The true label (y), being unobserved, makes it impossible to define a likelihood function like we did earlier.

In such cases, advanced techniques like expectation maximization are pretty helpful.

Expectation maximization (EM)

EM is an iterative optimization technique to estimate the parameters of statistical models. It is particularly useful when we have an unobserved (or hidden) label. One example situation could be as follows:



As depicted above, we assume that the data was generated from multiple distributions (a mixture). However, the observed/complete data does not contain that information. In other words, the observed dataset does not have information about whether a specific row was generated from distribution 1 or distribution 2.

Had it contained the label (y) information, we would have already used MLE.

EM helps us with parameter estimates of such datasets. The core idea behind EM is as follows:

- Make a guess about the initial parameters (θ).
- Expectation (E) step: Compute the posterior probabilities of the unobserved label (let's call it 'z') using the above parameters.

$$P(z = 1 | X; \theta)$$

$$P(z = 0 | X; \theta)$$

- Here, ‘z’ is also called a latent variable, which means hidden or unobserved.
- Relating it to our case, we know that the true label exists in nature. But we don’t know what it is.
- Thus, we replace it with a latent variable ‘z’ and estimate its posterior probabilities using the guessed parameters.
- Given that we now have a proxy (not precise, though) for the true label, we can define an “expected likelihood” function. Thus, we use the above posterior probabilities to do so:

$$\begin{array}{ccc}
 & \text{dependent} \\
 L(\theta) & \xrightarrow{\text{on}} & P(z = 1|X; \theta) \\
 & \text{Expected} \\
 & \text{Likelihood} \\
 & \text{function} & P(z = 0|X; \theta)
 \end{array}$$

- Maximization (M) step: So now we have a likelihood function to work with. Maximizing it with respect to the parameters will give us a new estimate for the parameters (θ').
- Next, we use the updated parameters (θ') to recompute the posterior probabilities we defined in the expectation step.

$$\begin{aligned}
 & P(z = 1|X; \theta') \\
 & P(z = 0|X; \theta')
 \end{aligned}$$

- We will update the likelihood function (L) using the new posterior probabilities.
- Again, maximizing it will give us a new estimate for the parameters (θ).
- And this process goes on and on until convergence.

The point is that in expectation maximization, we repeatedly iterate between the E and the M steps until the parameters converge. A good thing about EM is that it always converges. Yet, at times, it might converge to a local extrema.

MLE vs. EM is a popular question asked in many data science interviews.

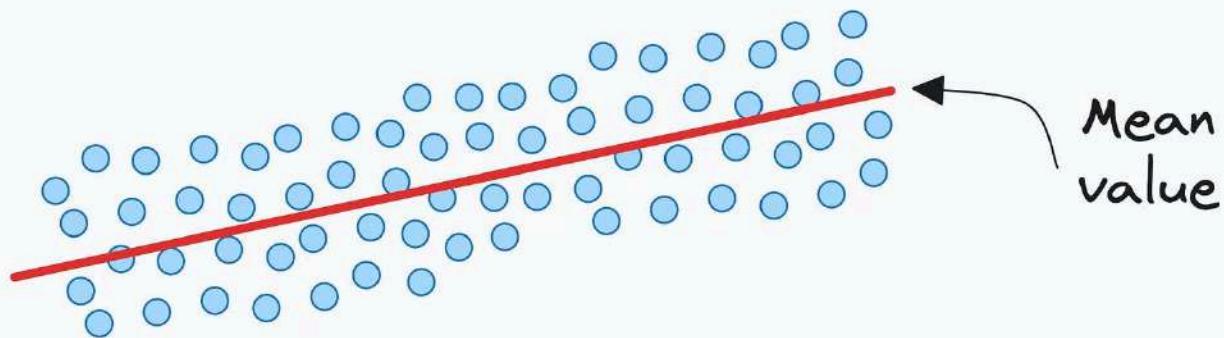
Confidence Interval and Prediction Interval

Statistical estimates always have some uncertainty.

For instance, a linear regression model never predicts an actual value.

Consider a simple example of modeling house prices just based on its area. A prediction wouldn't tell the true value of a house based on its area. This is because different houses of the same size can have different prices.

Instead, what it predicts is the **mean value** related to the outcome at a particular input.



The point is...

There's always some uncertainty involved in statistical estimates, and it is important to communicate it.

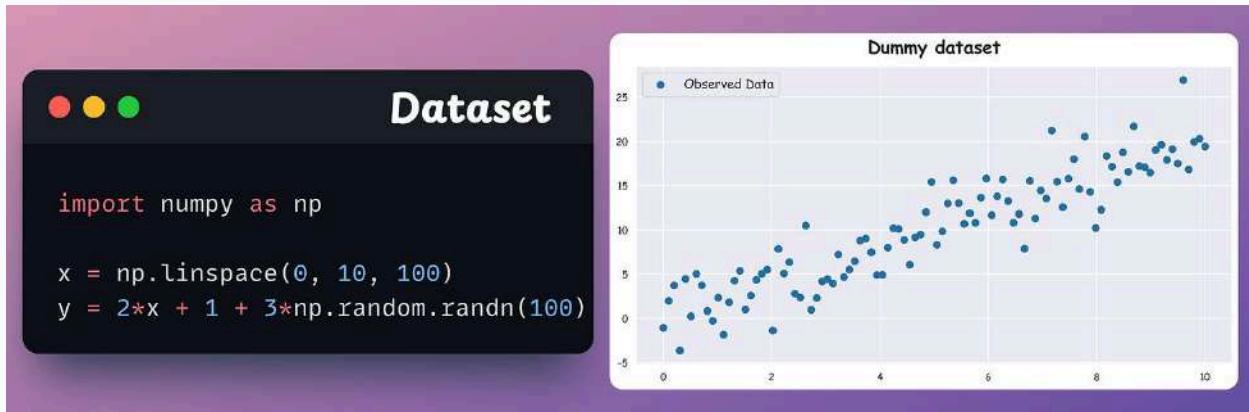
In this specific case, there are two types of uncertainties:

- The uncertainty in estimating the true mean value.
- The uncertainty in estimating the true value.

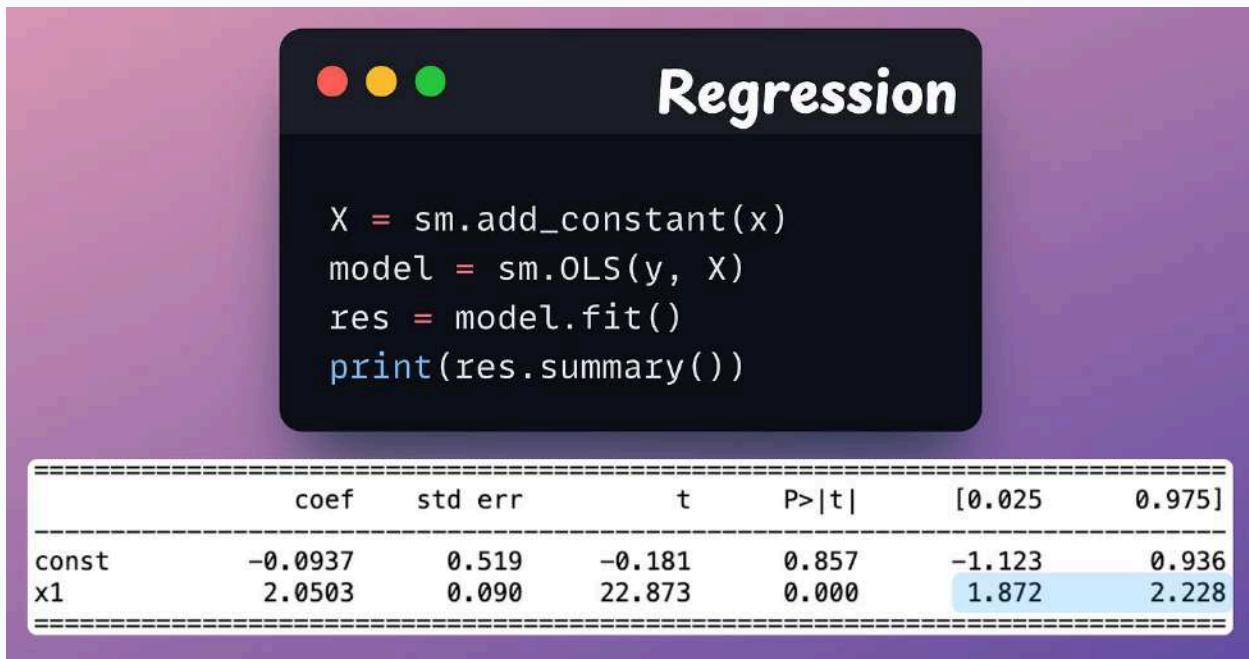
Confidence interval and prediction interval help us capture these uncertainties.

Let's understand.

Consider the following dummy dataset:



Let's fit a linear regression model using statsmodel and print a part of the regression summary:

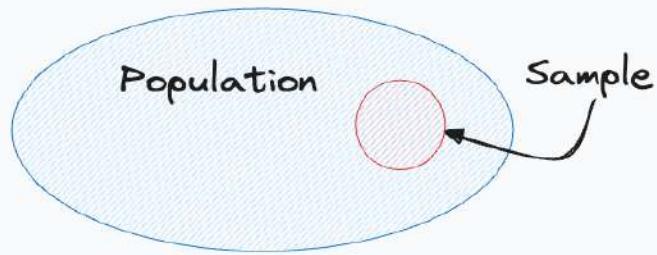


Notice that the coefficient of the predictor “x1” is **2.0503** with a **95%** interval of **[1.872, 2.228]**.

It is a 95% interval because $0.975 - 0.025 = 0.95$.

This is known as the confidence interval, which comes from **sampling uncertainty**.

More specifically, this uncertainty arises because the data we used above for modeling is just a **sample** of the population.



So, if we gathered more such samples and fit an OLS to each sample, the true coefficient (which we can only know if we had the data for the entire population) would lie 95% of the time in this confidence interval.

Next, we use this model to make a prediction as follows:

	mean	mean_se	mean_ci_lower	mean_ci_upper	obs_ci_lower	obs_ci_upper
0	7.62	0.391	6.832	8.419	1.06	14.19
Predicted value			95% Confidence interval		95% Prediction interval	

- The predicted value is **7.62 (mean)**.
- The **95% confidence interval** is **[6.832, 8.419]**.
- The **95% prediction interval** is **[1.06, 14.19]**.

The confidence interval we saw above was for the coefficient, so what does the confidence interval represent in this case?

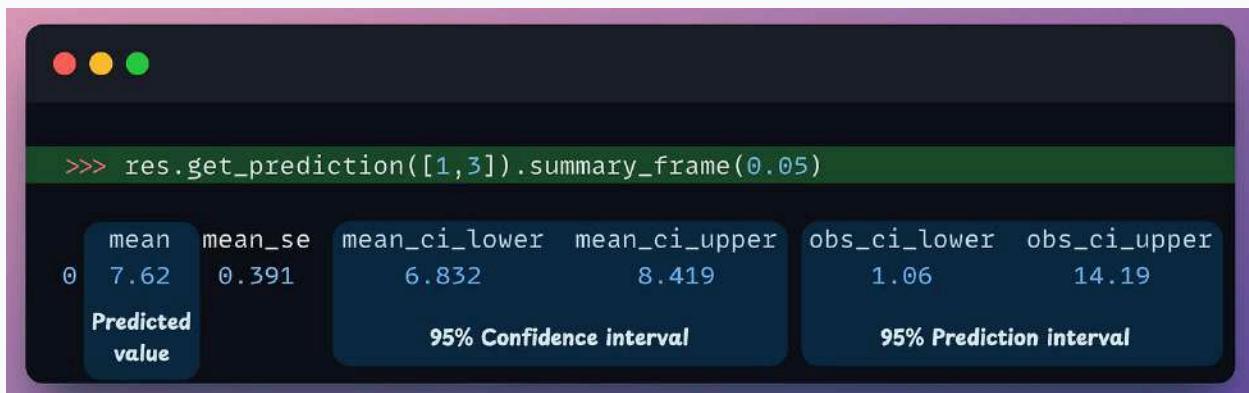
Similar to what we discussed above, the data is just a **sample** of the population.

The regression fit obtained by this sample produced a prediction (some mean value) for the input $x=3$.

However, if we gathered more such samples and fit an OLS to each dataset, the **true mean value** (which we can only know if we had the data for the entire

population) for this specific input ($x=3$) would lie 95% of the time in this confidence interval.

Coming to the prediction interval...



```
>>> res.get_prediction([1,3]).summary_frame(0.05)
```

	mean	mean_se	mean_ci_lower	mean_ci_upper	obs_ci_lower	obs_ci_upper
0	7.62	0.391	6.832	8.419	1.06	14.19

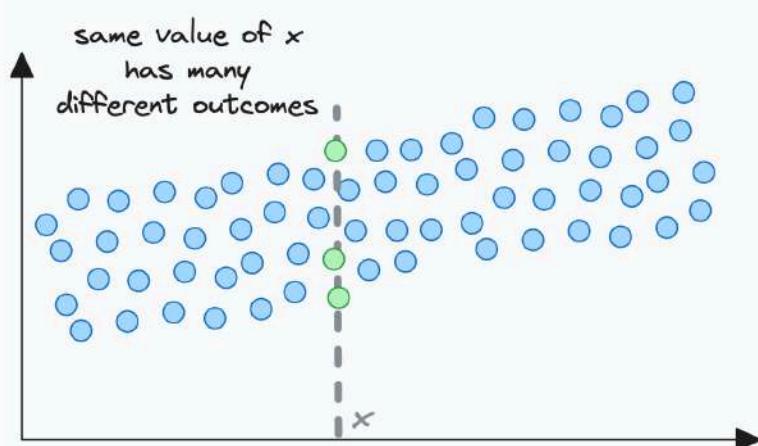
Predicted value **95% Confidence interval** **95% Prediction interval**

...we notice that it is wider than the confidence interval. Why is it, and what does this interval tell?

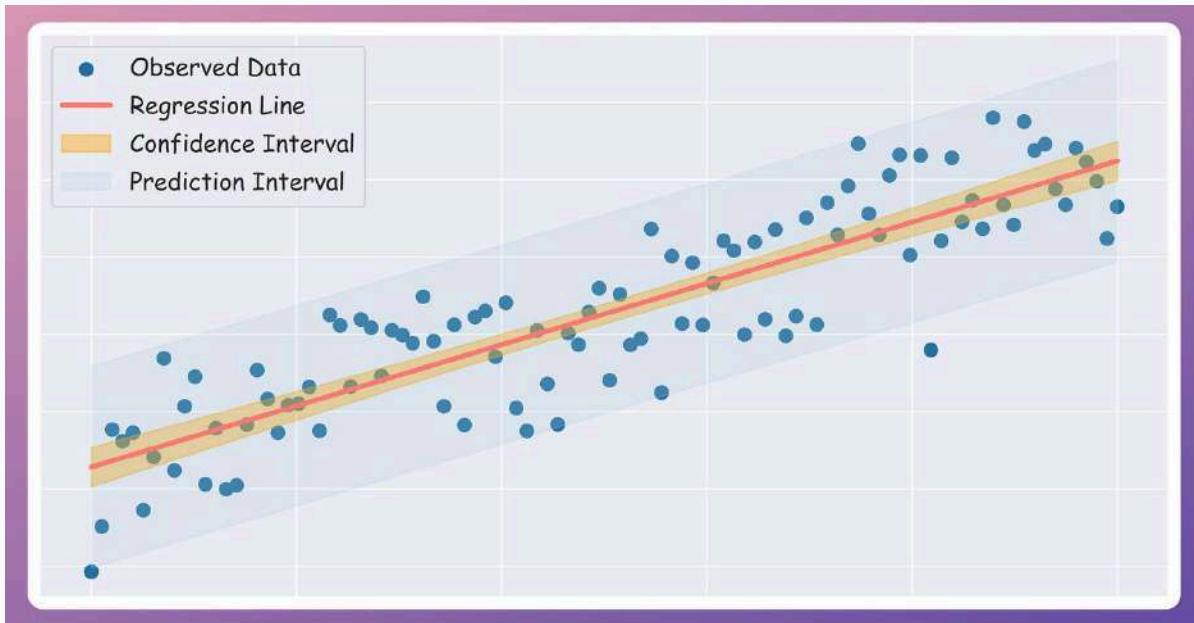
What we saw above with confidence interval was about estimating the true population mean at a specific input.

What we are talking about now is obtaining an interval where the **true value** for an input can lie.

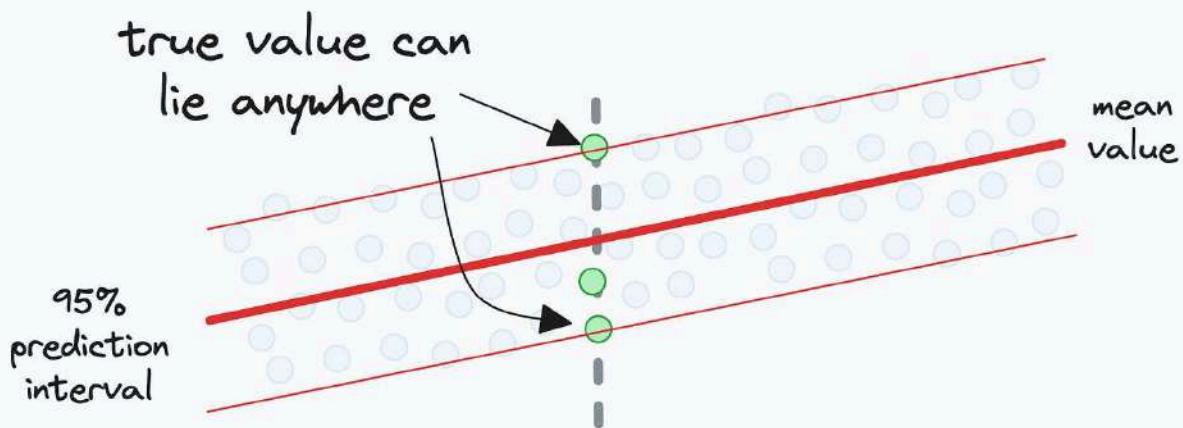
Thus, this additional uncertainty appears because in our dataset, for the same value of input x , there can be multiple different values of the outcome. This is depicted below:



Thus, it is wider than the confidence interval. Plotting it across the entire input range, we get the following plot:



Given that the model is predicting a mean value (as depicted below), we have to represent the prediction uncertainty that the actual value can lie anywhere in the prediction interval:



A 95% prediction interval tells us that we can be 95% sure that the actual value of this observation will fall within this interval.

So to summarize:

- A confidence interval captures the sampling uncertainty. More data means less sampling uncertainty, which in turn leads to a smaller interval.
- In addition to the sampling uncertainty, the prediction interval also represents the uncertainty in estimating the true value of a particular data point. Thus, it is wider than the confidence interval.

Communicating these uncertainties is quite crucial in decision-making because it provides a clearer understanding of the reliability and precision of predictions.

This transparency allows stakeholders to make more informed decisions by considering the range of possible outcomes and the associated risks.

Why is OLS Called an Unbiased Estimator?

The OLS estimator for linear regression (shown below) is known as an unbiased estimator.

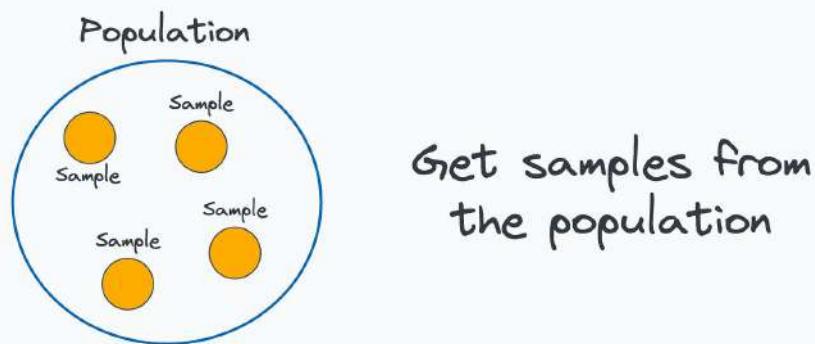
$$\theta = (X^T X)^{-1} X^T y$$

- What do we mean by that?
- Why is OLS called such?

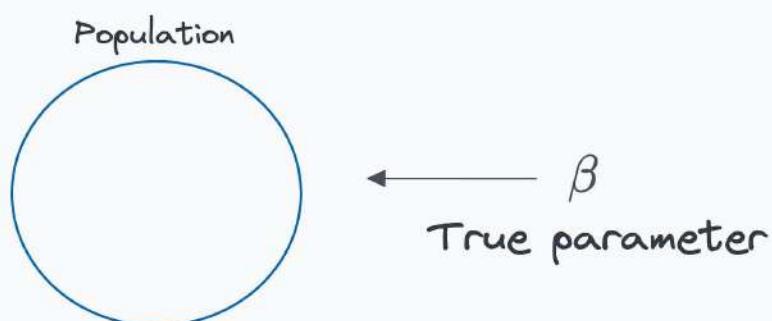
Background

The goal of statistical modeling is to make conclusions about the whole population.

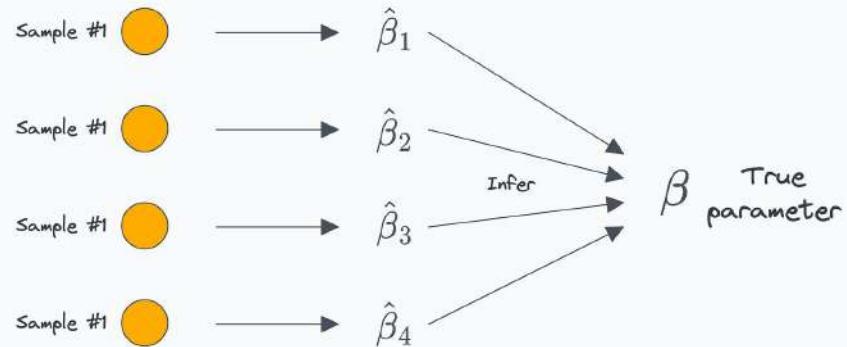
However, it is pretty obvious that observing the entire population is impractical.



In other words, given that we cannot observe (or collect data of) the entire population, we cannot obtain the true parameter (β) for the population:



Thus, we must obtain parameter estimates ($\hat{\beta}$) on samples and infer the true parameter (β) for the population from those estimates:



And, of course, we want these sample estimates ($\hat{\beta}$) to be reliable to determine the actual parameter (β).

The OLS estimator ensures that.

Let's understand how!

True population model

When using a linear regression model, we assume that the response variable (Y) and features (X) for the entire population are related as follows:

$$Y = \beta * X + \epsilon$$

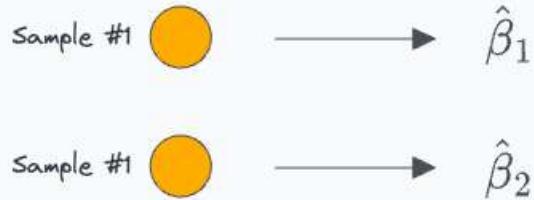
- β is the true parameter that we are not aware of.
- ϵ is the error term.

Expected Value of OLS Estimates

The closed-form solution of OLS is given by:

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

What's more, as discussed above, using OLS on different samples will result in different parameter estimates:



Let's find the expected value of OLS estimates $E[\hat{\beta}]$.

Simply put, the expected value is the average value of the parameters if we run OLS on many samples.

This is given by:

$$E[\hat{\beta}] = E[(X^T X)^{-1} X^T Y]$$

Substitute $\hat{\beta}$ as the OLS solution

Here, substitute $Y = \beta X + \varepsilon$:

$$\begin{aligned} E[\hat{\beta}] &= E[(X^T X)^{-1} X^T Y] \\ &= E[(X^T X)^{-1} X^T (\beta X + \epsilon)] \end{aligned}$$

If you are wondering how we can substitute $Y = \beta X + \varepsilon$ when we don't know what β is, then here's the explanation:

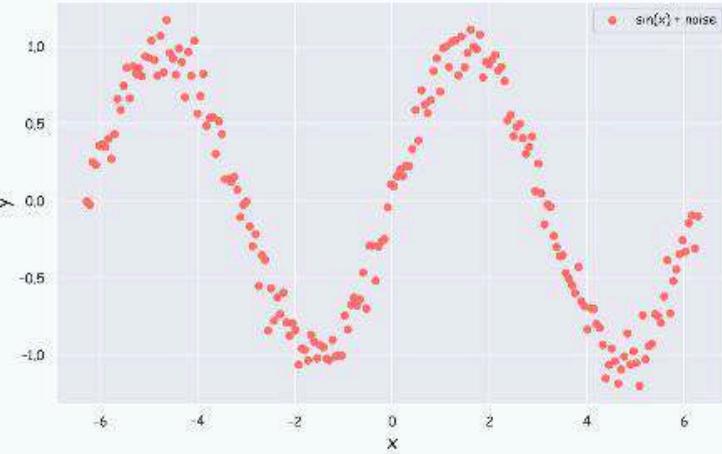
See, we can do that substitution because even if we don't know the parameter β for the whole population, we know that the sample was drawn from the population.

Thus, the equation in terms of the true parameters ($Y = \beta X + \varepsilon$) still holds for the sample.

Let me give you an example.

Say the population data was defined by $y = \sin(x) + \varepsilon$. Of course, we wouldn't know this, but just keep that aside for a second.

$$y = \sin(x) + \epsilon$$



Now, even if we were to draw samples from this population data, the true equation $y = \sin(x) + \varepsilon$ would still be valid on the sampled data points, wouldn't it?

The same idea has been extended for expected value.

Coming back to the following:

$$\begin{aligned} E[\hat{\beta}] &= E[(X^T X)^{-1} X^T Y] \\ &= E[(X^T X)^{-1} X^T (\beta X + \epsilon)] \end{aligned}$$

Let's open the inner parenthesis:

$$E[\hat{\beta}] = E[(X^T X)^{-1} X^T X \beta] + E[(X^T X)^{-1} X^T \epsilon]$$

↓ Some constant
 I

Simplifying, we get:

$$E[\hat{\beta}] = E[\beta] + (X^T X)^{-1} X^T E[\epsilon]$$

zero

β is constant so its expected value is β

And finally, what do we get?

$$E[\hat{\beta}] = \beta$$

The expected value of parameter estimates on the samples equals the true parameter value β .

And this is precisely what the definition of an unbiased estimator is.

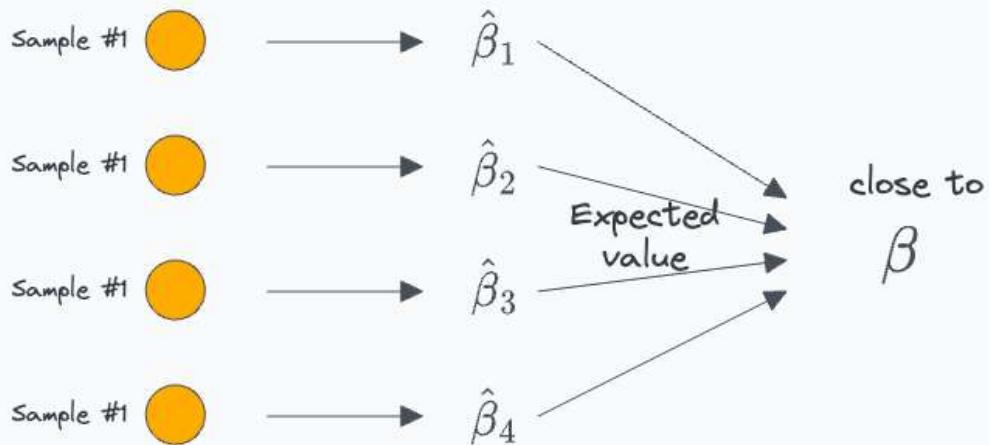
More formally, an estimator is called unbiased if the expected value of the parameters is equal to the actual parameter value.

And that is why we call OLS an unbiased estimator.

An important takeaway

Many people misinterpret unbiasedness with the idea that the parameter estimates from a single run of OLS on a sample are equal to the true parameter values.

Don't make that mistake.



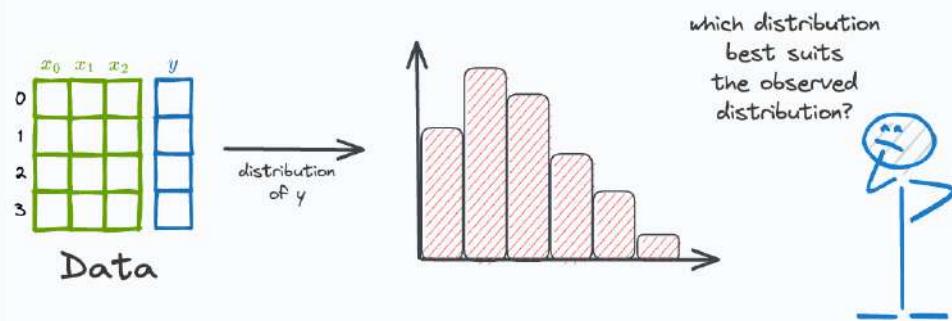
Instead, unbiasedness implies that if we were to generate OLS estimates on many different samples (drawn from the same population), then the expected value of obtained estimates will be equal to the true population parameter.

And, of course, all this is based on the assumption that we have good representative samples and that the assumptions of linear regression are not violated.

Bhattacharyya Distance

Assessing the similarity between two probability distributions is quite helpful at times. For instance, imagine we have a labeled dataset (X, y) .

By analyzing the label (y) distribution, we may hypothesize its distribution before building a statistical model.



We also looked at this in an earlier chapter on [generalized linear models \(GLMs\)](#).

While visual inspection is often helpful, this approach is quite subjective and may lead to misleading conclusions.

Is this more like Poisson or Binomial?



Thus, it is essential to be aware of quantitative measures as well. Bhattacharyya distance is one such reliable measure.

It quantifies the similarity between two probability distributions.

The core idea is to approximate the overlap between two distributions, which measures the “closeness” between the two distributions under consideration.

Bhattacharyya distance is measured as follows:

For two discrete probability distributions

$$D_B(P, Q) = -\ln \left(\sum_{x \in X} \sqrt{P(x) \cdot Q(x)} \right)$$

For two continuous probability distributions (replace summation with an integral):

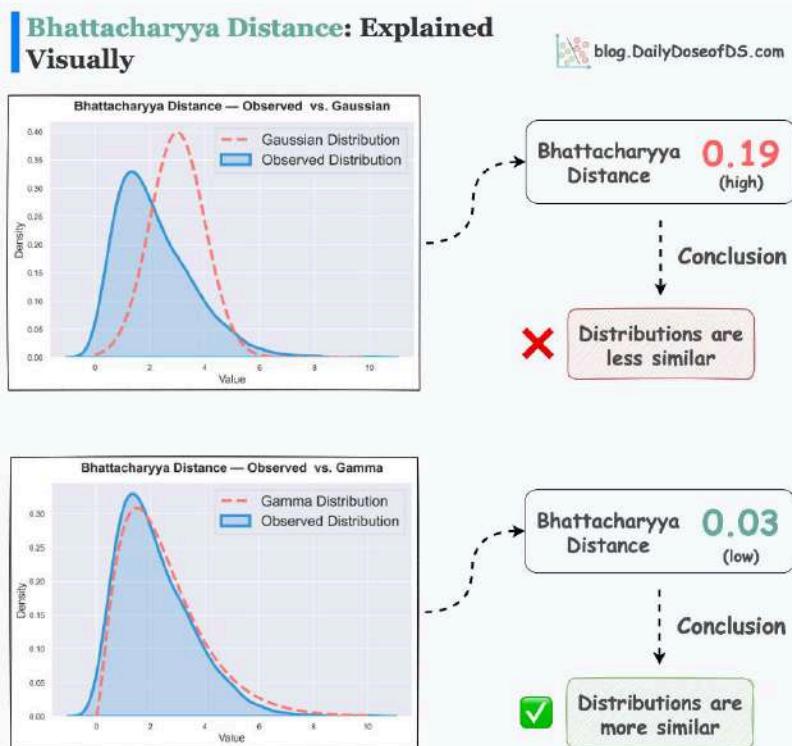
$$D_B(P, Q) = -\ln \left(\int \sqrt{P(x) \cdot Q(x)} dx \right)$$

Its effectiveness is evident from the image below.

Here, we have an observed distribution (Blue). Next, we measure its distance from:

- Gaussian $\rightarrow 0.19$.
- Gamma $\rightarrow 0.03$.

A high Bhattacharyya distance indicates less overlap or more dissimilarity. This lets us conclude that the observed distribution resembles a Gamma distribution.



The results also resonate with visual inspection.

KL Divergence vs. Bhattacharyya distance

Now, many often get confused between KL Divergence and Bhattacharyya distance. Effectively, both are quantitative measures to determine the “similarity” between two distributions.

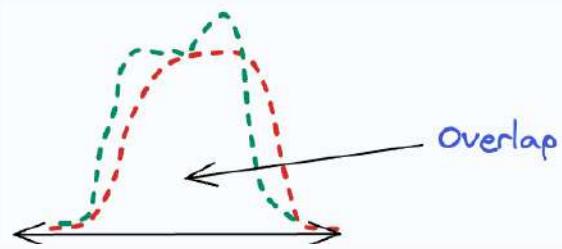
However, their notion of “similarity” is entirely different.

The core idea behind KL Divergence is to assess how much information is lost when one distribution is used to approximate another distribution.



The more information is lost, the more the KL Divergence and, consequently, the less the “similarity”. Also, approximating a distribution Q using P may not be the same as doing the reverse — P using Q. This makes KL Divergence asymmetric in nature.

Moving on,
Bhattacharyya distance
measures the overlap
between two
distributions.



This “overlap” is often interpreted as a measure of closeness (or distance) between the two distributions under consideration. Thus, Bhattacharyya distance primarily serves a distance metric, like Euclidean, for instance.

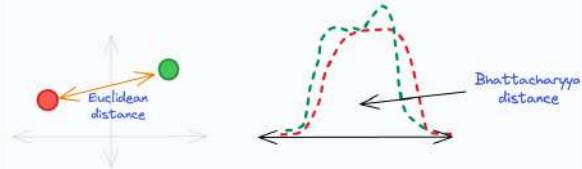
Being a distance metric, it is symmetric in nature. We can also verify this from the distance formula:

$$D_B(P, Q) = -\ln \left(\int \sqrt{P(x) \cdot Q(x)} dx \right)$$

↓

$$D_B(P, Q) = D_B(Q, P) \quad \text{Symmetric}$$

Just like we use Euclidean distance to find the distance between two points, we can use Bhattacharyya distance to find the distance between two distributions.



But if we intend to measure the amount of information lost when we approximate one distribution using another, KL divergence is more apt. In fact, KL divergence also serves as a loss function in machine learning algorithms at times (in t-SNE).

- Say we have an observed distribution (P) and want to approximate it with another simpler distribution Q .
- So, we can define a simpler parametric distribution Q .
- Next, we can measure the information lost by approximating P using Q with KL divergence.
- As we want to minimize the information lost, we can use KL divergence as our objective function.
- Finally, we can use gradient descent to determine the parameters of Q such that we minimize the KL divergence.

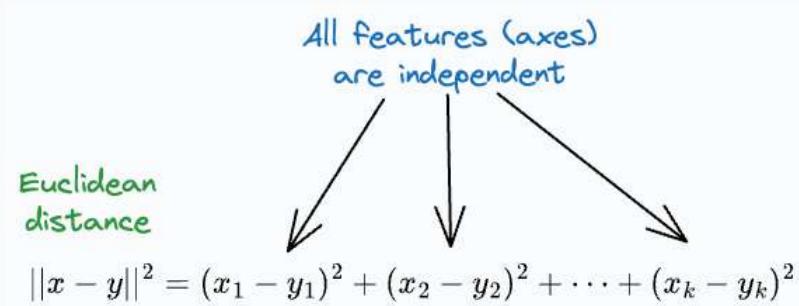
Bhattacharyya distance has many applications, not just in machine learning but in many other domains. For instance:

- Using this distance, we can simplify complex distributions to simple ones if the distance is low.
- In image processing, Bhattacharyya distance is often used for image matching. By comparing the color or texture distributions of images, it helps identify similar objects or scenes, etc.

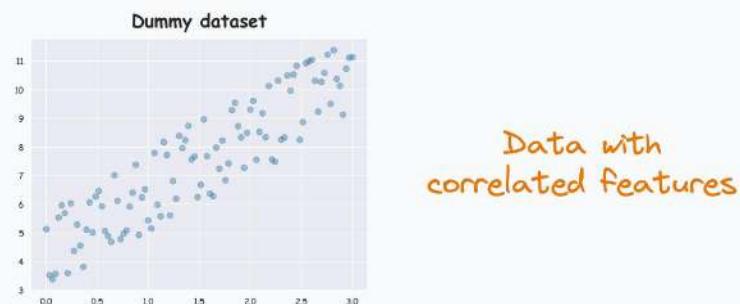
The only small caveat is that Bhattacharyya distance does not satisfy the triangle inequality, so that's something to keep in mind.

Why Prefer Mahalanobis Distance Over Euclidean distance?

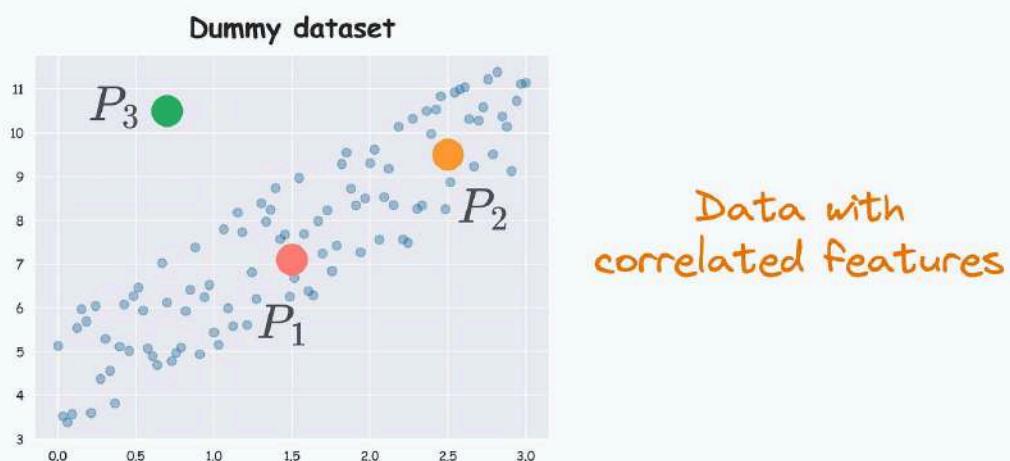
During distance calculation, Euclidean distance assumes independent axes.



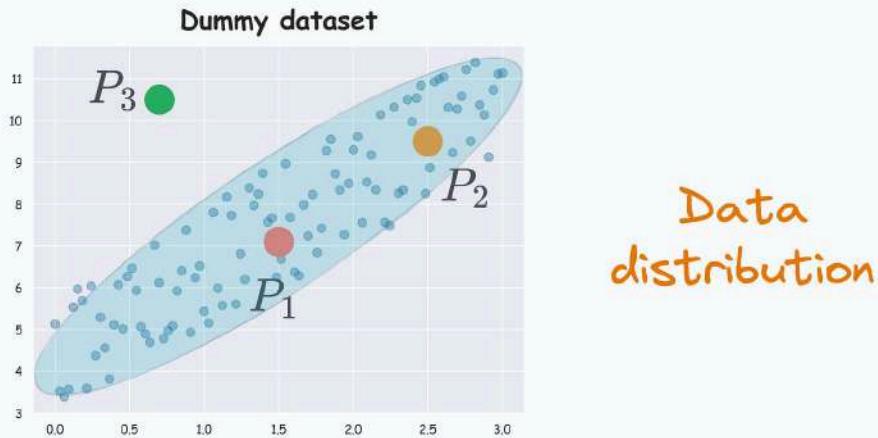
Thus, Euclidean distance will produce misleading results if your features are correlated. For instance, consider this dummy dataset below:



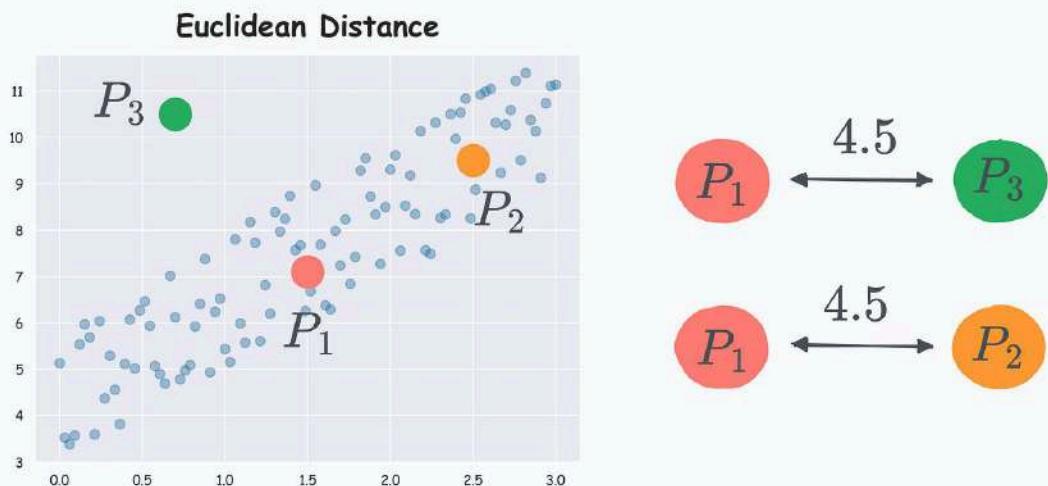
Clearly, the features are correlated. Here, consider three points marked P1, P2, and P3 in this dataset.



Considering the data distribution, something tells us that P2 is closer to P1 than P3. This is because P2 lies more within the data distribution than P3.



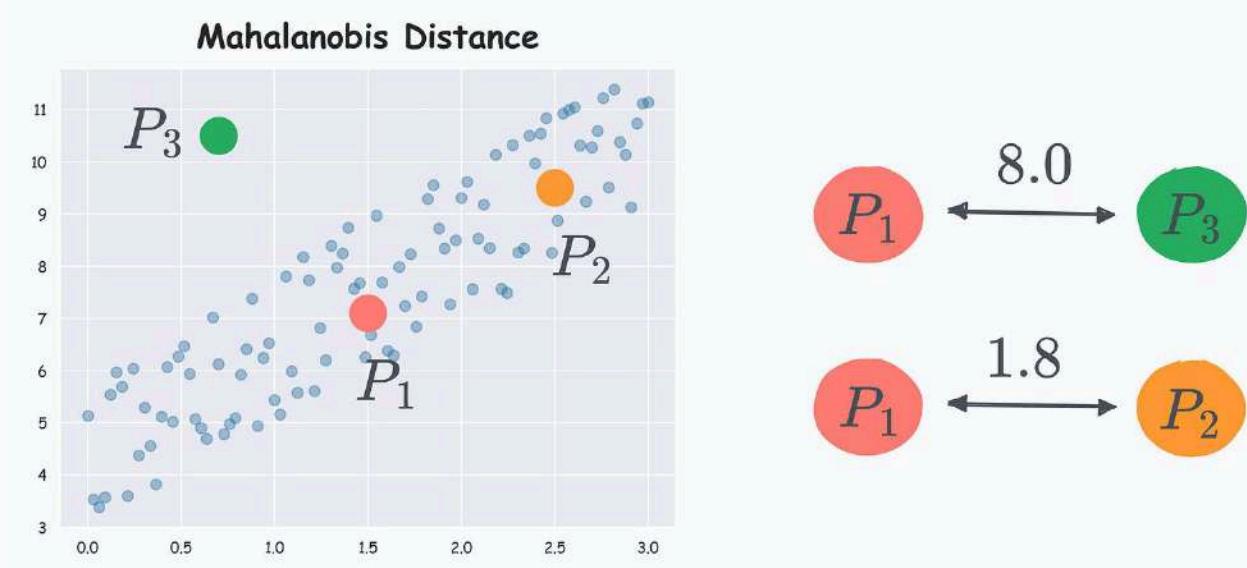
Yet, Euclidean distance ignores this, and P2 and P3 come out to be equidistant to P1, as depicted below:



Mahalanobis distance addresses this limitation. It is a distance metric that takes into account the data distribution.

As a result, it can measure how far away a data point is from the distribution, which Euclidean can not.

Referring to the earlier dataset again, with Mahalanobis distance, P2 comes out to be closer to P1 than P3.



How does it work?

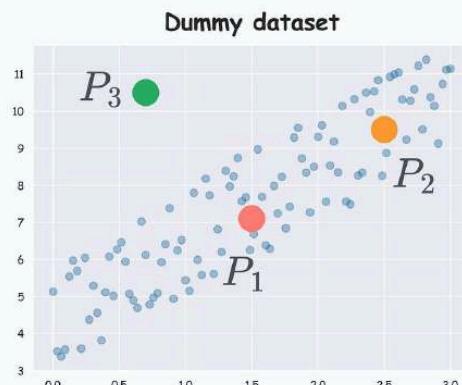
In a gist, the objective is to construct a new coordinate system with independent and orthogonal axes. The steps are:

- Step 1: Transform the columns into uncorrelated variables.
- Step 2: Scale the new variables to make their variance equal to 1.
- Step 3: Find the Euclidean distance in this new coordinate system.

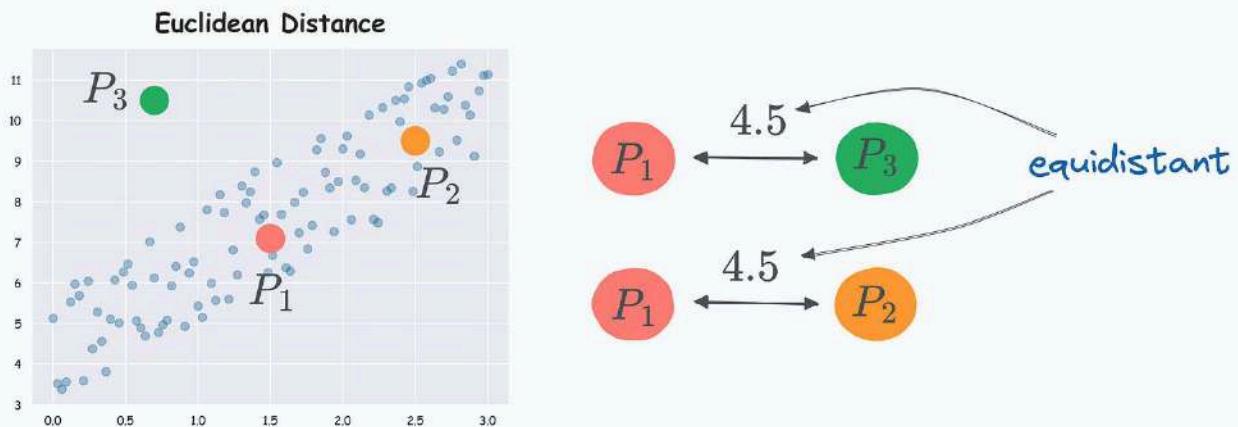
So, eventually, we do use Euclidean distance. However, we first transform the data to ensure that it obeys the assumptions of Euclidean distance.

Uses

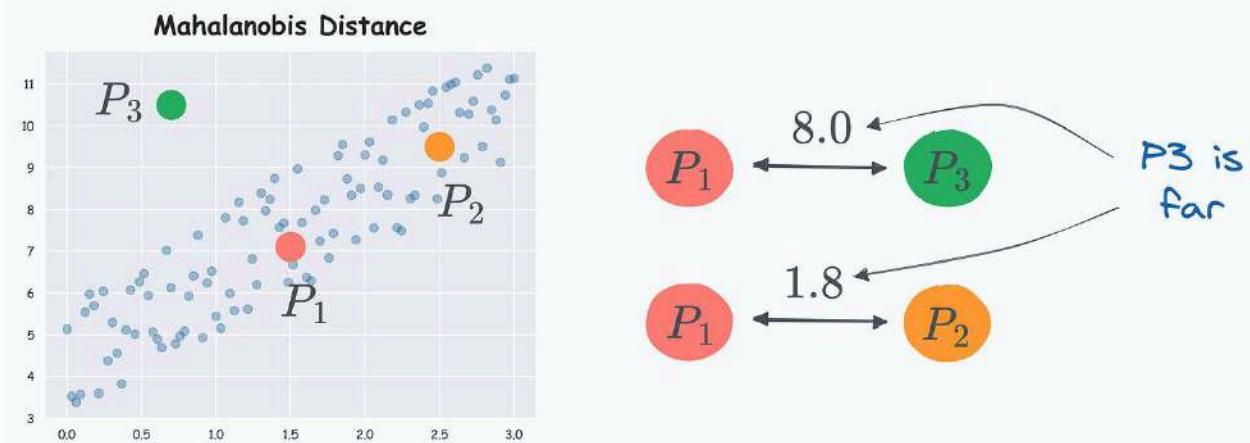
One of the most common use cases of Mahalanobis distance is outlier detection. Reconsidering the dataset we discussed earlier where P3 is clearly an outlier.



If we consider P1 as the distribution's centroid and use Euclidean distance, we will infer that P3 is not an outlier as both P2 and P3 are equidistant to P1.



Using Mahalanobis distance, however, provides a clearer picture:



This becomes more useful in a high-dimensional setting where visualization is infeasible.

Another use case we typically do not hear of often, but that exists is a variant of kNN that is implemented with Mahalanobis distance instead.

Scipy implements the Mahalanobis distance, which you can check here:
<https://bit.ly/3LjAymm>.

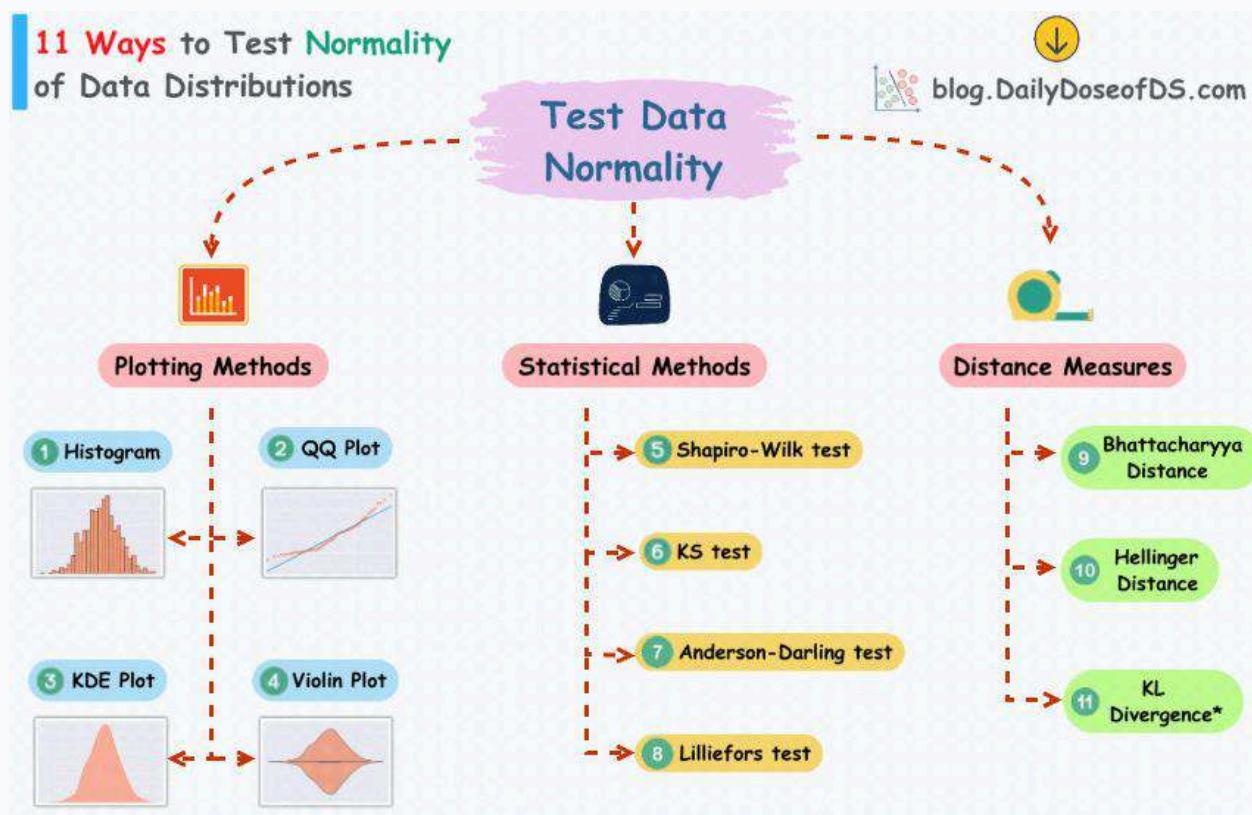
11 Ways to Determine Data Normality

Many ML models assume (or work better) under the presence of normal distribution.

For instance:

- Linear regression assumes residuals are normally distributed.
- At times, transforming the data to normal distribution can be beneficial.
- Linear discriminant analysis (LDA) is derived under the assumption of normal distribution, etc.

Thus, being aware of the ways to test normality is extremely crucial for data scientists. The visual below depicts the 11 essential ways to test normality.



Let's understand these in this chapter.

#1) Plotting Methods (*self-explanatory*)

- Histogram
- QQ Plot (We shall cover it in the plotting section of this book).
- KDE Plot
- Violin Plot

While plotting is often reliable, it is a subjective approach and prone to errors.

Thus, we must know reliable quantitative measures as well.

#2) Statistical Methods:

- 1) Shapiro-Wilk test:
 - Finds a statistic using the correlation between the observed data and the expected values under a normal distribution.
 - The p-value indicates the likelihood of observing such a correlation if the data were normally distributed.
 - A high p-value indicates a normal distribution.
- 2) KS test:
 - Measures the max difference between the cumulative distribution functions (CDF) of observed and normal distribution.
 - The output statistic is based on the max difference between the two CDFs.
 - A high p-value indicates a normal distribution.
- 3) Anderson-Darling test:
 - Measures the differences between the observed data and the expected values under a normal distribution.
 - Emphasizes the differences in the tail of the distribution.
 - This makes it particularly effective at detecting deviations in the extreme values.

4) Lilliefors test:

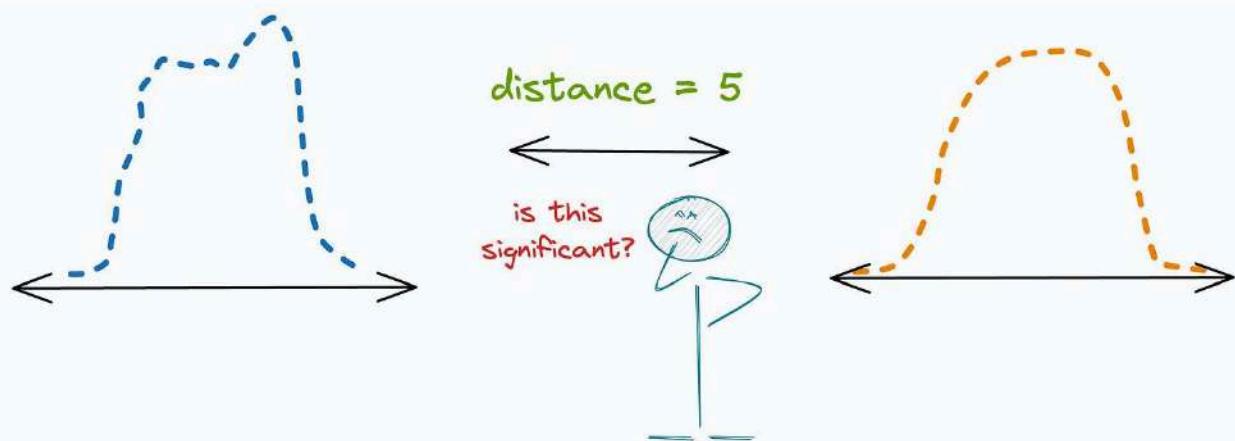
- It is a modification of the KS test.
- The KS test is appropriate in situations where the parameters of the reference distribution are known.
- If the parameters are unknown, Lilliefors is recommended.
- Get started: [Statsmodel Docs](#).

#3) Distance Measures

Distance measures are another reliable and more intuitive way to test normality.

But they can be a bit tricky to use.

See, the problem is that a single distance value needs more context for interpretability.



For instance, if the distance between two distributions is 5, is this large or small?

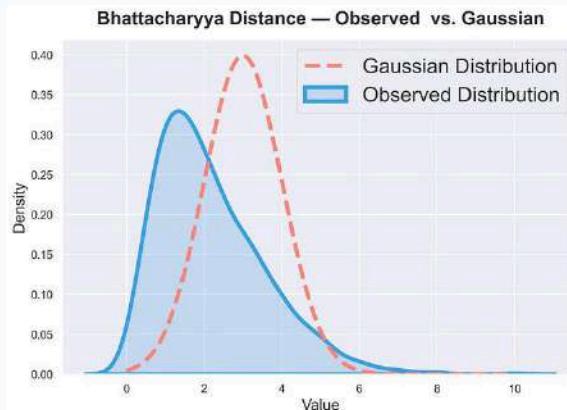
We need more context.

I prefer using these measures as follows:

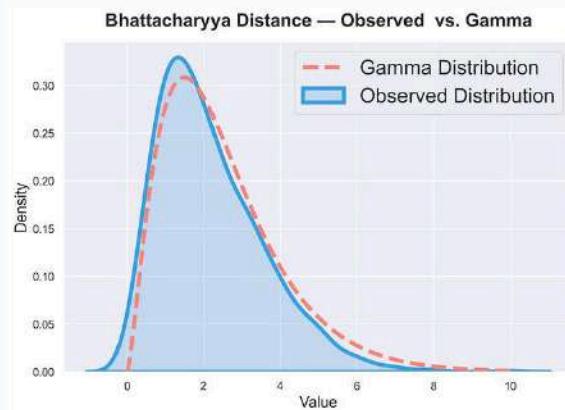
- Find the distance between the observed distribution and multiple reference distributions.
- Select the reference distribution with the minimum distance to the observed distribution.

Here are a few distance common and useful measures:

1) Bhattacharyya distance:



Bhattacharyya
distance 0.19



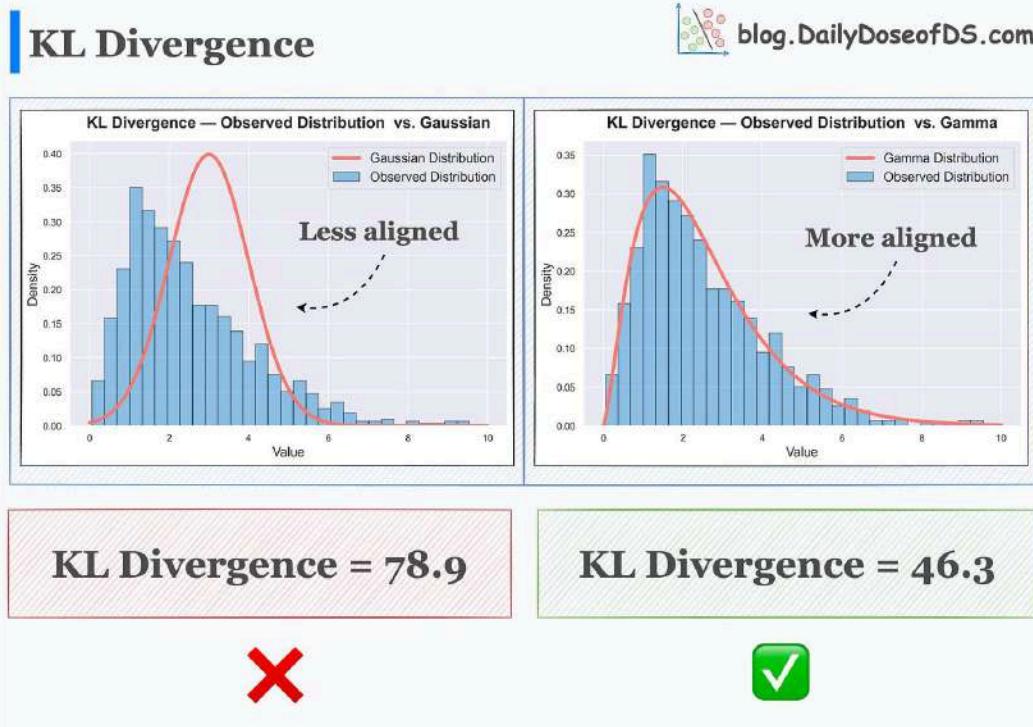
Bhattacharyya
distance 0.03

- Measure the overlap between two distributions.
- This “overlap” is often interpreted as closeness between two distributions.
- Choose the distribution that has the least Bhattacharyya distance to the observed distribution.

2) Hellinger distance:

- It is used quite similar to how we use the Bhattacharyya distance
- The difference is that Bhattacharyya distance does not satisfy triangular inequality.
- But Hellinger distance does.

3) KL Divergence:



- It is not entirely a "distance metric" per se, but can be used in this case.
- Measure information lost when one distribution is approximated using another distribution.
- The more information is lost, the more the KL Divergence.
- Choose the distribution that has the least KL divergence from the observed distribution.

KL divergence is used as a loss function in the t-SNE algorithm.

Probability vs. Likelihood

In data science and statistics, many folks often use “probability” and “likelihood” interchangeably.

However, likelihood and probability DO NOT convey the same meaning.

And the misunderstanding is somewhat understandable, given that they carry similar meanings in our regular language.

While writing this chapter, I searched for their meaning in the Cambridge Dictionary. Here's what it says:

- Probability: the level of possibility of something happening or being true.
- Likelihood: the chance that something will happen.

If you notice closely, “likelihood” is the only synonym of “probability”.

The screenshot shows the Cambridge Dictionary entry for the word "probability". The word is in large blue bold letters at the top. Below it, the part of speech is listed as "noun [C or U]". The pronunciation is given as UK /prə'bɪlətɪ/ US /prə'be'bɪlətɪ/. There is a yellow "Add to word list" button with three horizontal lines. A blue arrow points from the word "likelihood" in the synonym section below to the word "probability" in the main definition. The main definition is in bold: "the level of possibility of something happening or being true:". Below this, a bulleted list provides examples: "What is the probability of winning?", "The probability of getting all the answers correct is about one in ten.", "There's a high/strong probability (that) (= it is very likely that) she'll be here.", and "Until yesterday, the project was just a possibility, but now it has become a real probability (= it is likely to happen.)."

Anyway.

In my opinion, it is crucial to understand that probability and likelihood convey very different meanings in data science and statistics.

Let's understand!

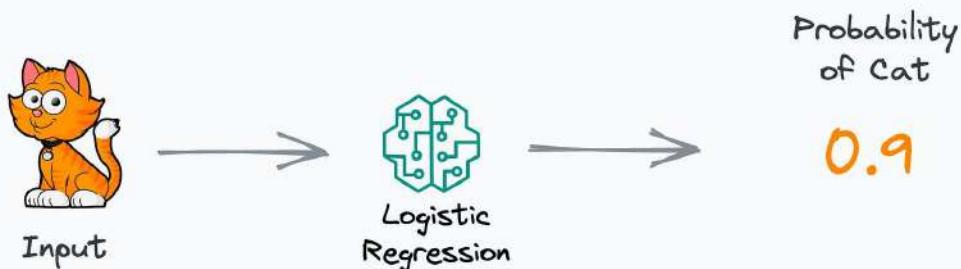
Probability is used in contexts where you wish to know the possibility/odds of an event.

For instance, what is the:

- Probability of obtaining an even number in a die roll?
- Probability of drawing an ace of diamonds from a deck?
- and so on...

When translated to ML, probability can be thought of as:

- What is the probability that a transaction is fraud?
- What is the probability that an image depicts a cat?
- and so on...



Essentially, many classification models, like logistic regression or a classification neural network, etc., assign the probability of a specific label to an input.

When calculating probability, the model's parameters are known. Also, we assume that they are trustworthy.

For instance, to determine the probability of a head in a coin toss, we mostly assume and trust that it is a fair coin.

Likelihood, on the other hand, is about explaining events that have already occurred.

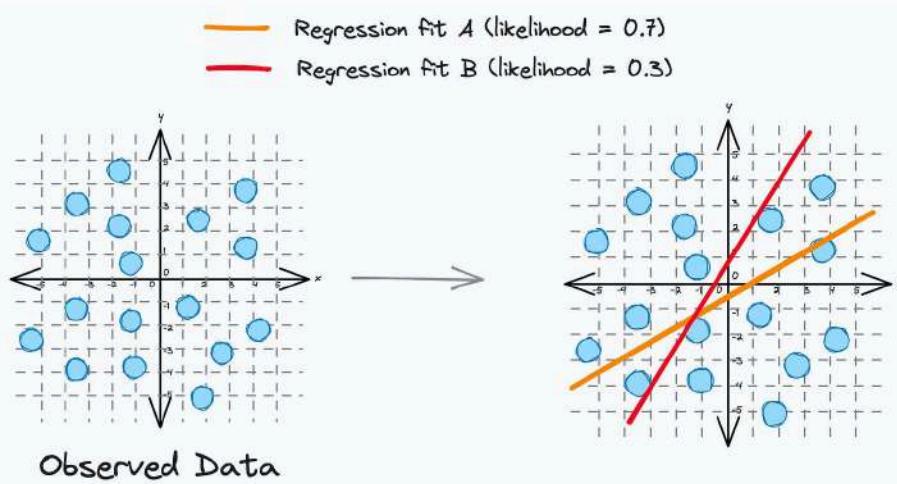
Unlike probability (where parameters are known and assumed to be trustworthy)...

...likelihood helps us determine if we can trust the parameters in a model based on the observed data.

Let me elaborate more on that.

Assume you have collected some 2D data and wish to fit a straight line with two parameters — slope (m) and intercept (c).

Here, likelihood is defined as the support provided by a data point for some particular parameter values in your model.



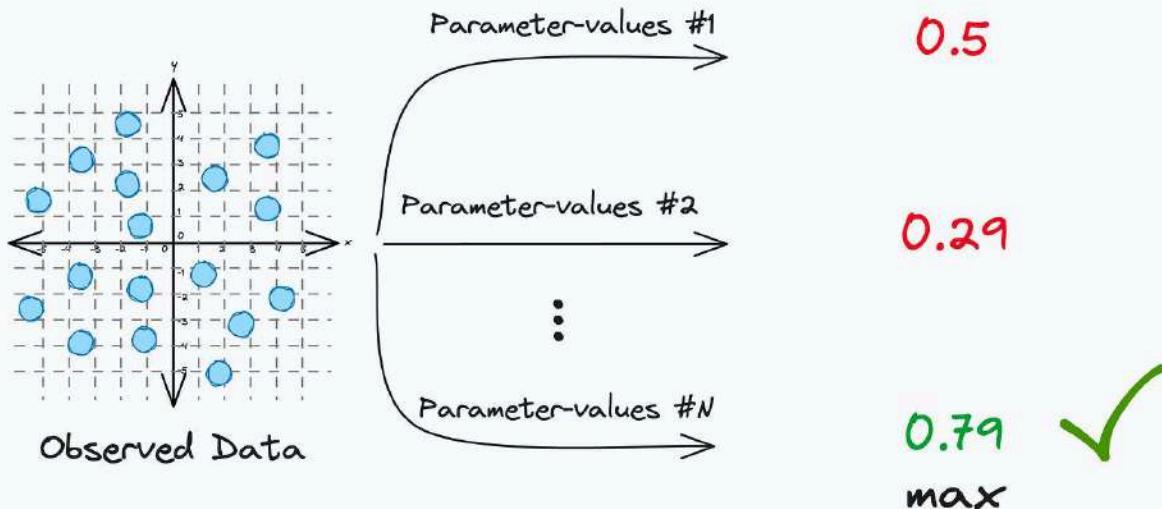
Here, you will ask questions like:

- If I model this data with the parameters:
 - $m=2$ and $c=1$, what is the likelihood of observing the data?
 - $m=3$ and $c=2$, what is the likelihood of observing the data?
 - and so on...

The above formulation popularly translates into the maximum likelihood estimation (MLE), which we discussed here: ([MLE vs. EM — What's the Difference?](#))

In maximum likelihood estimation, you have some observed data and you are trying to determine the specific set of parameters (θ) that maximize the likelihood of observing the data.

Likelihood



Using the term “likelihood” is like:

- I have a possible explanation for my data. (In the above illustration, “explanation” can be thought of as the parameters you are trying to determine)
- How well does my explanation explain what I’ve already observed? This is precisely quantified with likelihood.

For instance:

- Observation: The outcomes of 10 coin tosses are “HHHHHHHHHTHH”.
- Explanation: I think it is a fair coin ($p=0.5$).
- What is the likelihood that my above explanation is true based on the observed data?

To summarize...

It is immensely important to understand that in data science and statistics, likelihood and probability DO NOT convey the same meaning.

As explained above, they are pretty different.

In probability:

- We determine the possibility of an event.
- We know the parameters associated with the event and assume them to be trustworthy.

In likelihood:

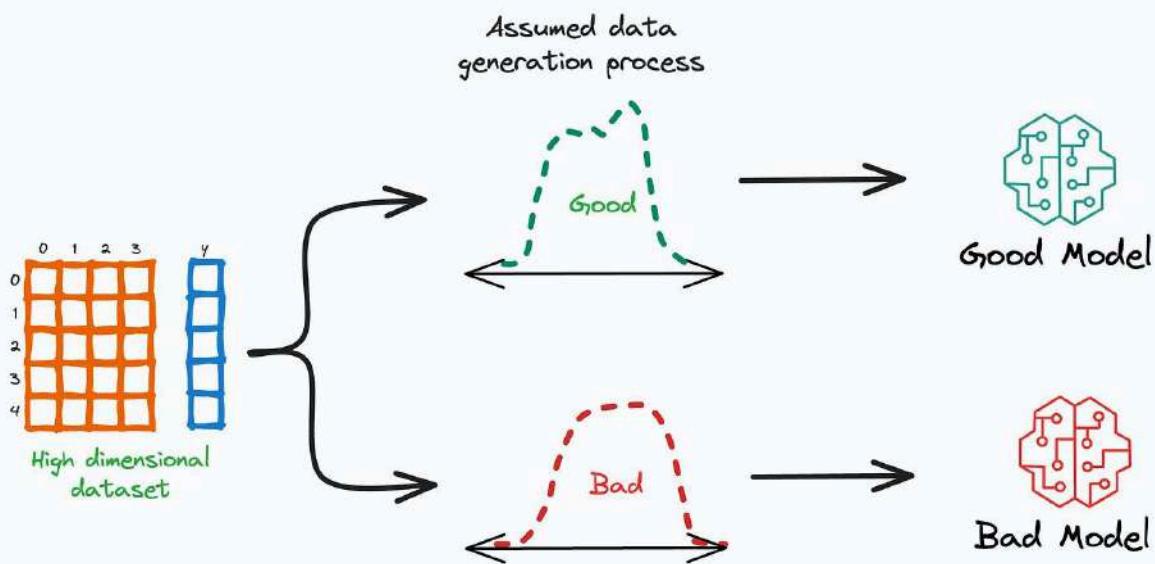
- We have some observations.
- We have an explanation (or parameters).
- Likelihood helps us quantify whether the explanation is trustworthy.

Hope that helped!

11 Key Probability Distributions in Data Science

Statistical models assume an underlying data generation process.

This is exactly what lets us formulate the generation process, using which we define the maximum likelihood estimation (MLE) step.



Thus, when dealing with statistical models, the model performance becomes entirely dependent on:

- Your understanding of the data generation process.
- The distribution you chose to model data with, which, in turn, depends on how well you understand various distributions.

We shall also look at this in the chapter on GLMs ([this chapter](#)).

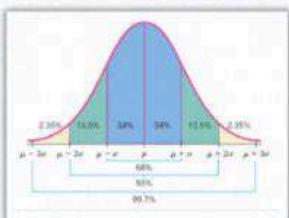
Thus, it is crucial to be aware of some of the most important distributions and the type of data they can model.

The visual below depicts the 11 most important distributions in data science:

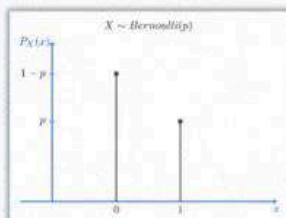
11 Most Important Distributions in Data Science



1 Normal Distribution



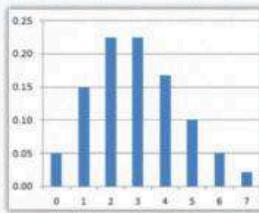
2 Bernoulli Distribution



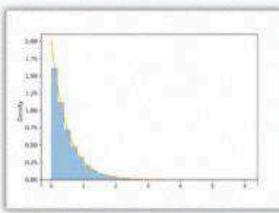
3 Binomial Distribution



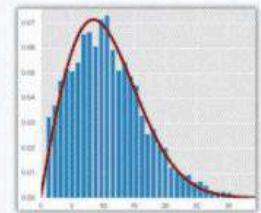
4 Poisson Distribution



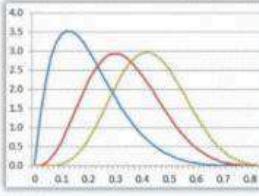
5 Exponential Distribution



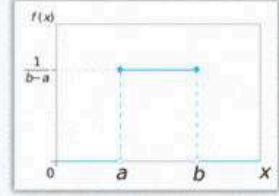
6 Gamma Distribution



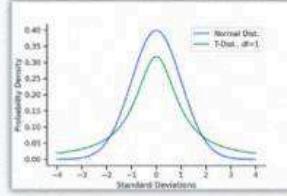
7 Beta Distribution



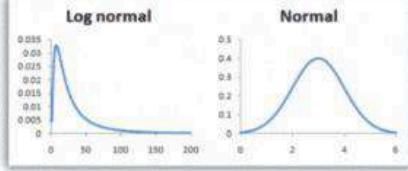
8 Uniform Distribution



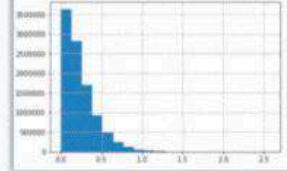
9 Student t-distribution



10 Log Normal Distribution



11 Weibull Distribution



Let's understand them briefly and how they are used.

Normal Distribution

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$$

- The most widely used distribution in data science.
- Characterized by a symmetric bell-shaped curve
- It is parameterized by two parameters—mean and standard deviation.
- Example: Height of individuals.

Bernoulli Distribution

$$P(X=x) = \begin{cases} p & \text{for } x=1 \\ 1-p & \text{for } x=0 \end{cases}$$

- A discrete probability distribution that models the outcome of a binary event.
- It is parameterized by one parameter—the probability of success.
- Example: Modeling the outcome of a single coin flip.

Binomial Distribution

$$P(X=x) = \frac{n!}{x!(n-x)!} p^x (1-p)^{(n-x)}$$

- It is Bernoulli distribution repeated multiple times.
- A discrete probability distribution that represents the number of successes in a fixed number of independent Bernoulli trials.
- It is parameterized by two parameters—the number of trials and the probability of success.

Poisson Distribution

$$P(X) = \frac{\lambda^x e^{-\lambda}}{X!}$$

- A discrete probability distribution that models the number of events occurring in a fixed interval of time or space.
- It is parameterized by one parameter—lambda, the rate of occurrence.
- Example: Analyzing the number of goals a team will score during a specific time period.

Exponential Distribution

$$f(x, \lambda) = \begin{cases} \lambda \cdot e^{-(\lambda x)} & x \geq 0 \\ 0 & x < 0 \end{cases}$$

- A continuous probability distribution that models the time between events occurring in a Poisson process.
- It is parameterized by one parameter—lambda, the average rate of events.
- Example: Analyzing the time between goals scored by a team.

Gamma Distribution

$$f(x; \alpha, \beta) = \begin{cases} \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} e^{-x/\beta} & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

- It is a variation of the exponential distribution.
- A continuous probability distribution that models the waiting time for a specified number of events in a Poisson process.
- It is parameterized by two parameters—alpha (shape) and beta (rate).
- Example: Analysing the time it would take for a team to score three goals.

Beta Distribution

PDF	$\frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)}$
	where $B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}$ and Γ is the Gamma function.

- It is used to model probabilities, thus, it is bounded between [0,1].
- Differs from Binomial in this respect that in Binomial, probability is a parameter.
- But in Beta, the probability is a random variable.

Uniform Distribution

$$f(x) = \begin{cases} \frac{1}{b-a}, & x \in [a, b] \\ 0, & otherwise \end{cases}$$

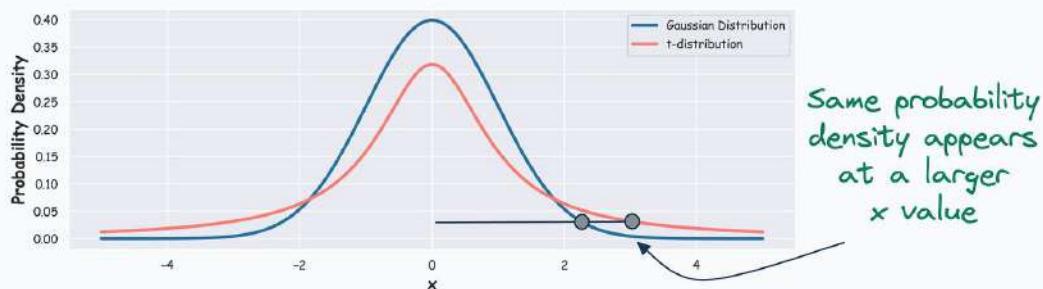
- All outcomes within a given range are equally likely.
- It can be continuous or discrete.
- It is parameterized by two parameters: a (min value) and b (max value).
- Example: Simulating the roll of a fair six-sided die, where each outcome (1, 2, 3, 4, 5, 6) has an equal probability.

Log-Normal Distribution

$$X \sim N(\mu, \sigma^2) \quad Y = e^X \\ Y \sim Lognormal(\mu, \sigma^2)$$

- A continuous probability distribution where the logarithm of the variable follows a normal distribution.
- It is parameterized by two parameters—mean and standard deviation.
- Example: Typically, in stock returns, the natural logarithm follows a normal distribution.

Student t-distribution



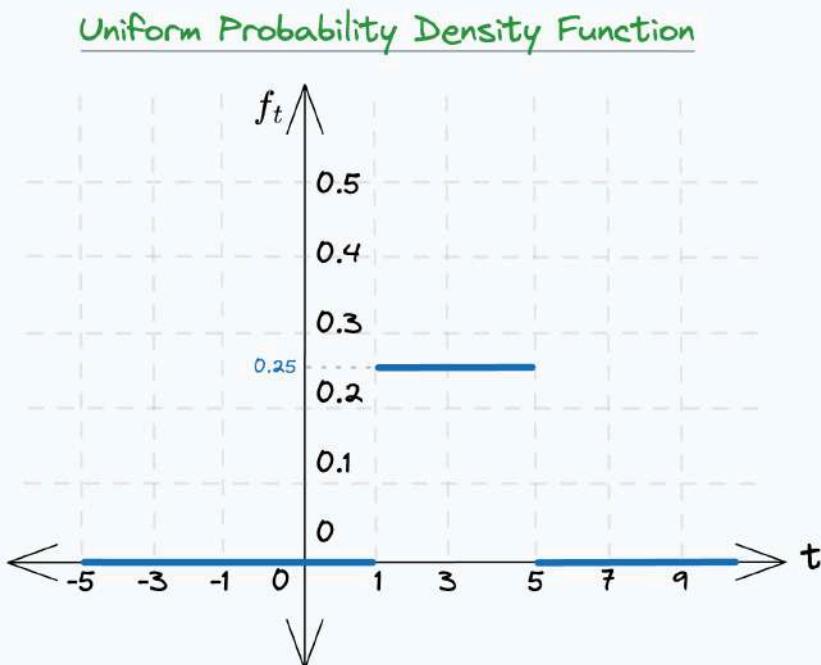
- It is similar to normal distribution but with longer tails (shown above).
- It is used in t-SNE to model low-dimensional pairwise similarities.

Weibull

- Models the waiting time for an event.
- Often employed to analyze time-to-failure data.

A Common Misinterpretation of Continuous Probability Distributions

Consider the following probability density function of a continuous probability distribution. Say it represents the time one may take to travel from point A to B.



- For simplicity, we are assuming a uniform distribution in the interval [1,5].
- Essentially, it says that it will take somewhere between 1 and 5 minutes to go from A to B. Never more, never less.

Thus, the probability density function (PDF) can be written as follows:

$$f_T(t) = \begin{cases} \frac{1}{5-1} & : 1 \leq t \leq 5 \\ 0 & : \text{otherwise} \end{cases}$$

Answer the following question for me:

Q) What is the probability that one will take precisely three minutes $P(T=3)$ to reach point B?

- A) $1/4$ (or 0.25)
- B) Area under the curve from $t=[1,3]$.
- C) Area under the curve from $t=[3,5]$.

Decide on an answer before you read further.

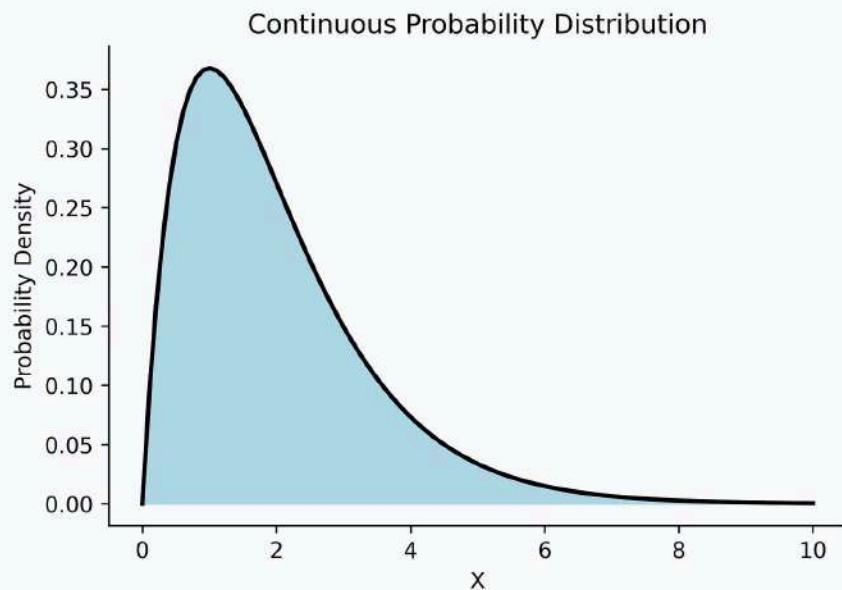
Well, all of the above answers are wrong.

The correct answer is ZERO.

And I intentionally kept only wrong answers here so that you never forget something fundamentally important about continuous probability distributions.

Let's dive in!

The probability density function of a continuous probability distribution may look as follows:



Some conditions for this probability density function are:

- It should be defined for all real numbers (can be zero for some values).

$$PDF \rightarrow f_X(x) ; x \in \mathbb{R}$$

This is in contrast to a discrete probability distribution which is only defined for a list of values.

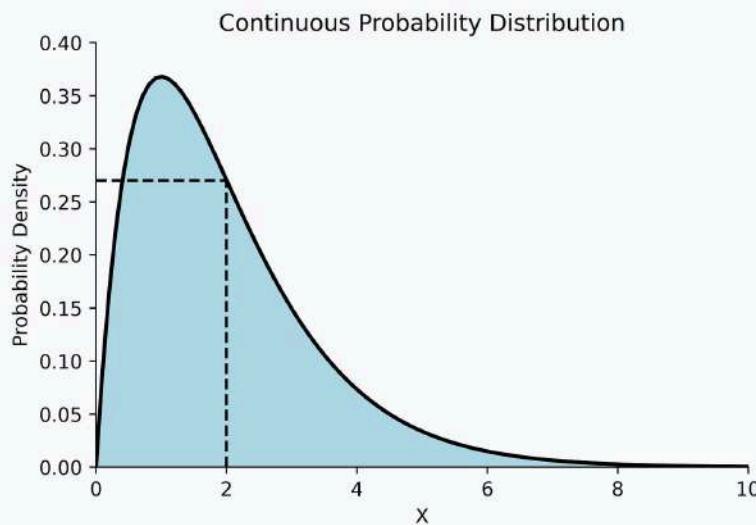
- The area should be 1.

$$\int_{-\infty}^{\infty} f_X(x) dx = 1$$

- The function should be non-negative for all real values.

$$f_X(x) \geq 0 \quad \forall x$$

Here, many folks often misinterpret that the probability density function represents the probability of obtaining a specific value.



For instance, by looking at the above probability density function, many incorrectly conclude that the probability of the random variable X being 2 is close to 0.27.

$$P(X = 2) = 0.27$$

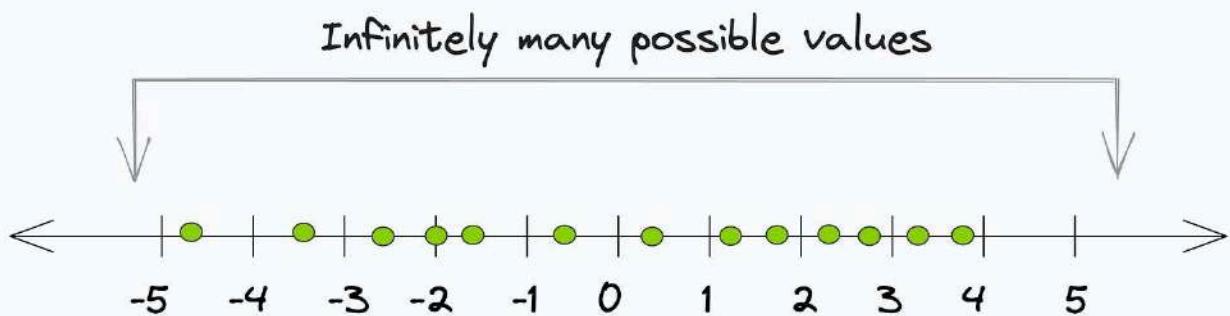
But contrary to this common belief, a probability density function:

- DOES NOT depict the probabilities of a specific value.
- is not meant to depict a discrete random variable.

Instead, a probability density function:

- depicts the rate at which probabilities accumulate around each point.
- is only meant to depict a continuous random variable.

Now, there are infinitely possible values that a continuous random variable may take.



So the probability of obtaining a specific value is always zero (or infinitesimally small).

Thus, answering our original question, the probability that one will take three minutes to reach point B is ZERO.

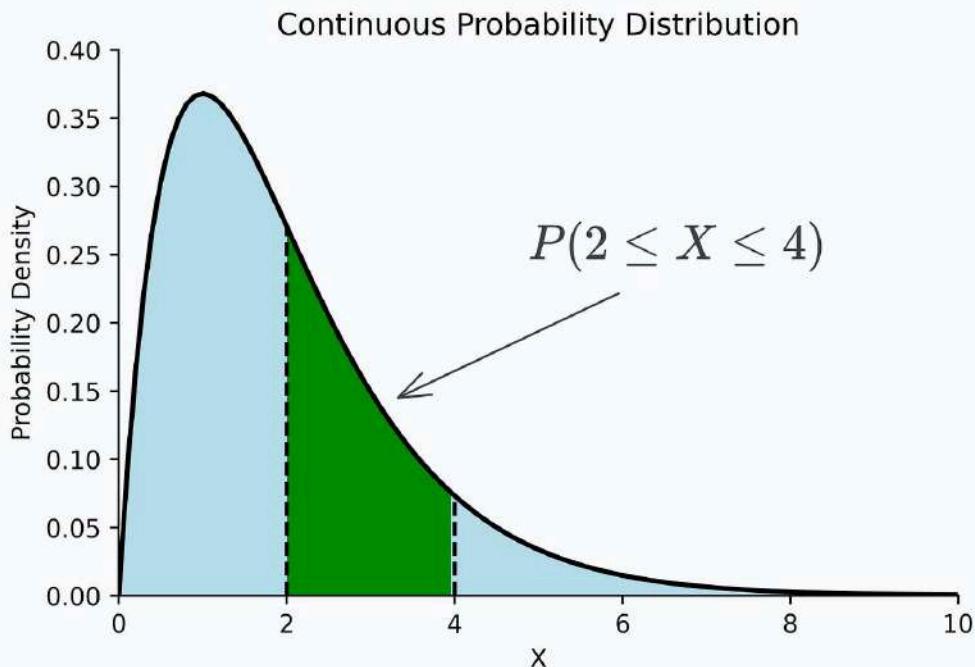
So what is the purpose of using a probability density function?

In statistics, a PDF is used to calculate the probability over an interval of values.

Thus, we can use it to answer questions such as...

- What is the probability that it will take between:
 - 3 to 4 minutes to reach point B from point A, or,
 - 2 to 4 minutes to reach point B from point A, and so on...

And we do this using integrals.



More formally, the probability that a random variable X will take values in the interval $[a,b]$ is:

$$P(a \leq X \leq b) = \int_a^b f_X(x) dx$$

Simply put, it's the area under the curve from $[a,b]$.

From the above probability estimation over an interval, we can also verify that the probability of obtaining a specific value is indeed zero.

By substituting $b=a$, we get:

$$P(a \leq X \leq b) = \int_a^b f_X(x) dx$$

Substitute $b = a$

$$P(a \leq X \leq a) = \int_a^a f_X(x) dx$$

$$P(a \leq X \leq a) = 0$$

$P(a) = 0$

To summarize, always remember that in a continuous probability distribution:

- The probability density function does not depict the exact probability of obtaining a specific value.
- Estimating the probability for a precise value of the random value makes no sense because it is infinitesimally small.

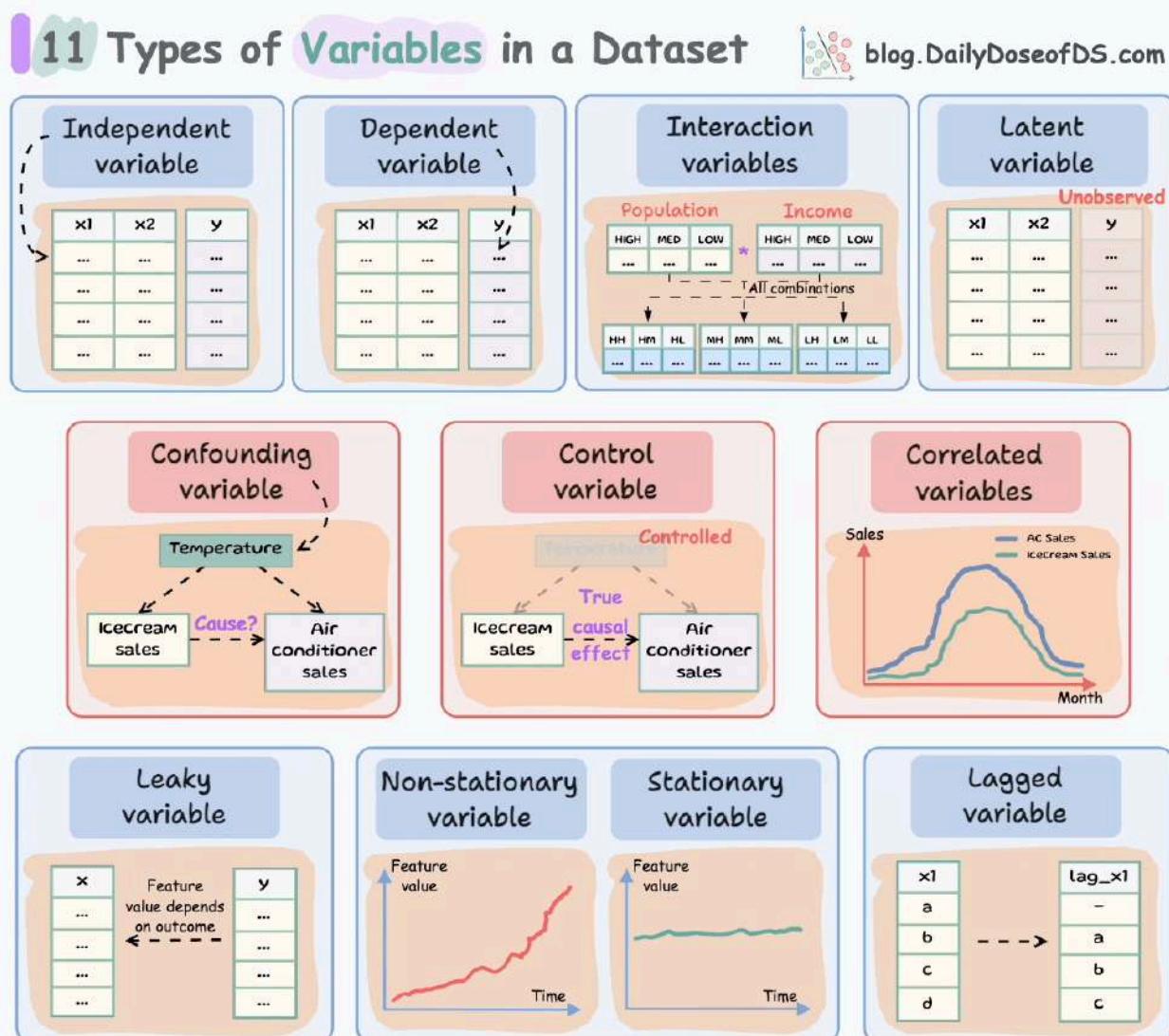
Instead, we use the probability density function to calculate the probability over an interval of values.

Feature Definition, Engineering and Selection

11 Types of Variables in a Dataset

In any tabular dataset, we typically categorize the columns as either a feature or a target.

However, there are so many variables that one may find/define in their dataset, which I want to discuss in this chapter. These are depicted in the image below:



#1-2) Independent and dependent variables

These are the most common and fundamental to ML.

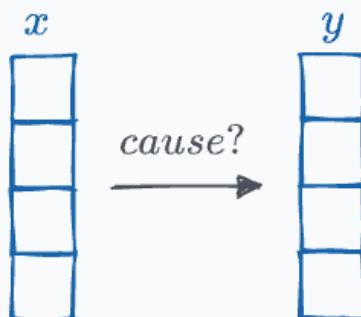
Independent variables are the features that are used as input to predict the outcome. They are also referred to as predictors/features/explanatory variables.



The dependent variable is the outcome that is being predicted. It is also called the target, response, or output variable.

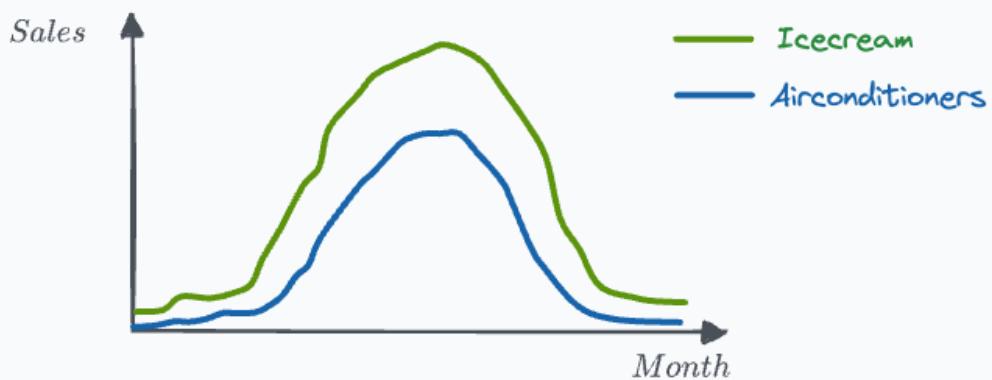
#3-4) Confounding and correlated variables

Confounding variables are typically found in a cause-and-effect study (causal inference).



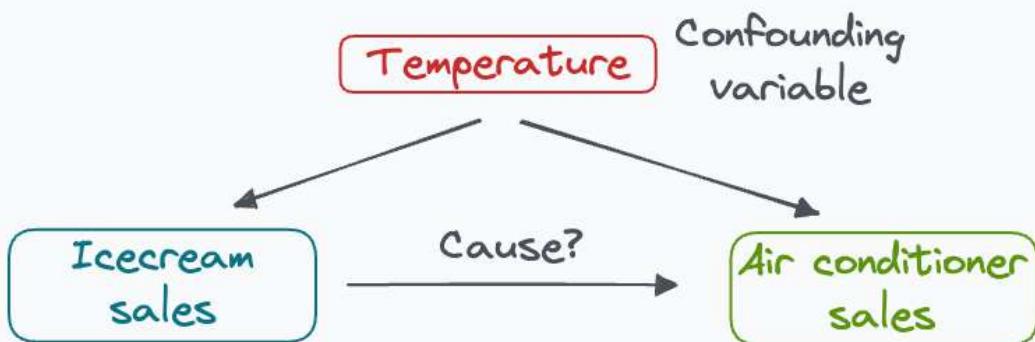
These variables are not of primary interest in the cause-and-effect equation but can potentially lead to spurious associations.

To exemplify, say we want to measure the effect of ice cream sales on the sales of air conditioners.



As you may have guessed, these two measurements are highly correlated.

However, there's a confounding variable — temperature, which influences both ice cream sales and the sales of air conditioners.



To study the true causal impact, it is essential to consider the confounder (temperature). Otherwise, the study will produce misleading results.

In fact, it is due to the confounding variables that we hear the statement: "Correlation does not imply causation."

In the above example:

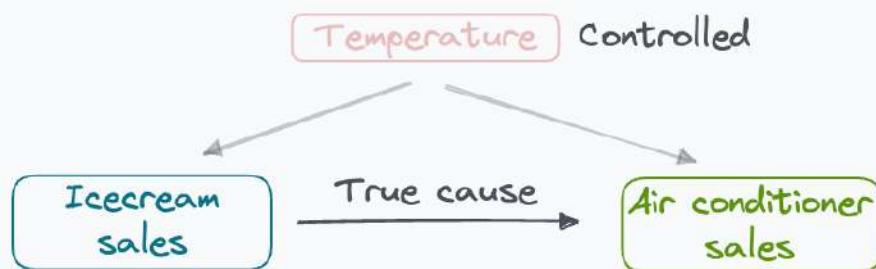
- There is a high correlation between ice cream sales and sales of air conditioners.
- But the sales of air conditioners (effect) are NOT caused by ice cream sales.

Also, in this case, the air conditioner and ice cream sales are correlated variables.

More formally, a change in one variable is associated with a change in another.

#5) Control variables

In the above example, to measure the true effect of ice cream sales on air conditioner sales, we must ensure that the temperature remains unchanged throughout the study.



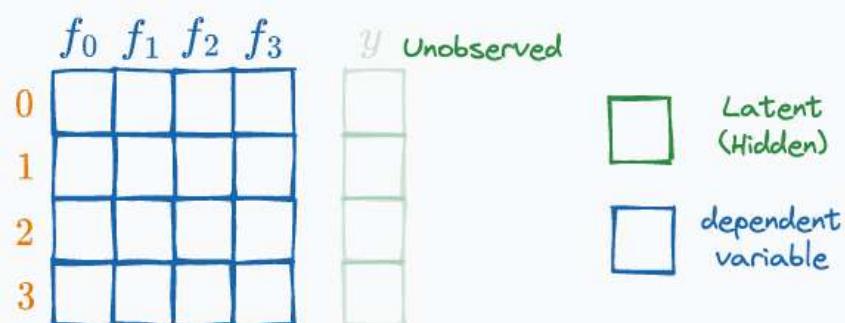
Once controlled, temperature becomes a control variable.

More formally, these are variables that are not the primary focus of the study but are crucial to account for to ensure that the effect we intend to measure is not biased or confounded by other factors.

#6) Latent variables

A variable that is not directly observed but is inferred from other observed variables.

For instance, we use clustering algorithms because the true labels do not exist, and we want to infer them somehow.



The true label is a latent variable in this case.

Another common example of a latent variable is “intelligence.”

Intelligence itself cannot be directly measured; it is a latent variable.

However, we can infer intelligence through various observable indicators such as test scores, problem-solving abilities, and memory retention.

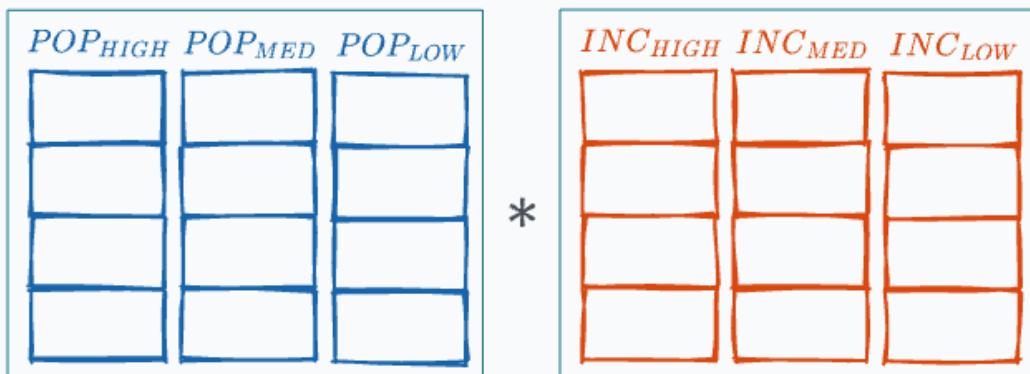
#7) Interaction variables

As the name suggests, these variables represent the interaction effect between two or more variables, and are often used in regression analysis.

Here's an instance I remember using them in.

In a project, I studied the impact of population density and income levels on spending behavior.

- I created three groups for population density — HIGH, MEDIUM, and LOW (one-hot encoded).
- Likewise, I created three groups for income levels — HIGH, MEDIUM, and LOW (one-hot encoded).



To do regression analysis, I created interaction variables by cross-multiplying both one-hot columns.

This produced 9 interaction variables:

- Population-High and Income-High
- Population-High and Income-Med
- Population-High and Income-Low
- Population-Med and Income-High
- and so on...

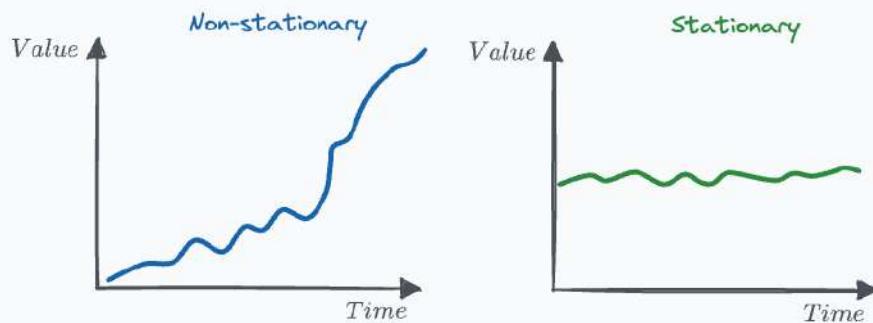
Conducting the regression analysis on interaction variables revealed more useful insights than what I observed without them.

To summarize, the core idea is to study two or more variables together rather than independently.

#8-9) Stationary and Non-Stationary variables:

The concept of stationarity often appears in time-series analysis.

Stationary variables are those whose statistical properties (mean, variance) DO NOT change over time.



On the flip side, if a variable's statistical properties change over time, they are called non-stationary variables.

Preserving stationarity in statistical learning is critical because these models are fundamentally reliant on the assumption that samples are identically distributed.

But if the probability distribution of variables is evolving over time, (non-stationary), the above assumption gets violated.

That is why, typically, using direct values of the non-stationary feature (like the absolute value of the stock price) is not recommended.

Instead, I have always found it better to define features in terms of relative changes:

$$\frac{\delta P}{P} \rightarrow \text{relative change in stock price}$$

#10) Lagged variables

Talking of time series, lagged variables are pretty commonly used in feature engineering and data analytics.

As the name suggests, a lagged variable represents previous time points' values of a given variable, essentially shifting the data series by a specified number of periods/rows.

	f_0	f_1
15-Jan 2024		a
15-Feb 2024		b
15-Mar 2024		c
15-Apr 2024		d
15-May 2024		e

30-day lag on f_1

	f_0	f_1	lag- f_1
15-Jan 2024		a	
15-Feb 2024		b	
15-Mar 2024		c	
15-Apr 2024		d	
15-May 2024		e	

For instance, when predicting next month's sales figures, we might include the sales figures from the previous month as a lagged variable.

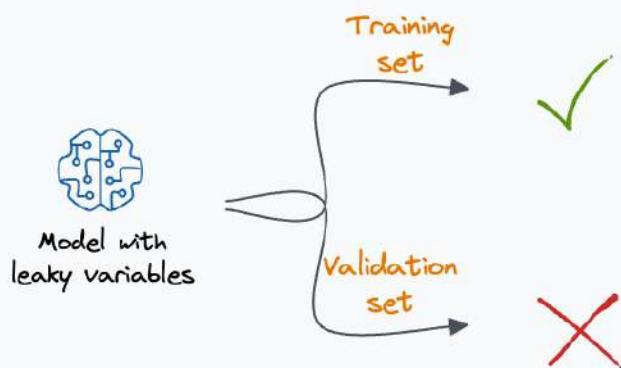
Lagged features may include:

- 7-day lag on website traffic to predict current website traffic.
- 30-day lag on stock prices to predict the next month's closing prices.
- And so on...

#11) Leaky variables

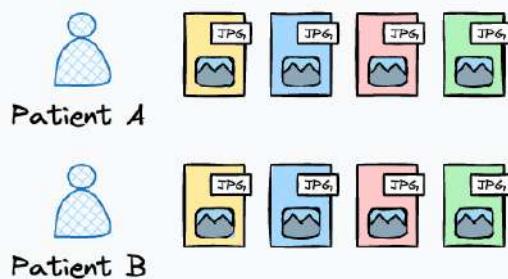
Yet again, as the name suggests, these variables (unintentionally) provide information about the target variable that would not be available at the time of prediction.

This leads to overly optimistic model performance during training but fails to generalize to new data.



Consider a dataset containing medical imaging data.

Each sample consists of multiple images (e.g., different views of the same patient's body part), and the model is intended to detect the severity of a disease.



In this case, randomly splitting the images into train and test sets will result in data leakage.

This is because images of the same patient will end up in both the training and test sets, allowing the model to “see” information from the same patient during training and testing.

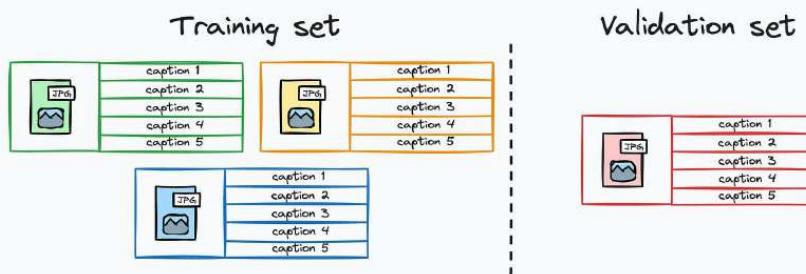
Here's a paper that committed this mistake (and later corrected it):

3.1. Training

We use the ChestX-ray14 dataset released by Wang et al. (2017) which contains 112,120 frontal-view X-ray images of 30,805 unique patients. Wang et al. (2017) annotate each image with up to 14 different thoracic pathology labels using automatic extraction methods on radiology reports. We label images that have pneumonia as one of the annotated pathologies as positive examples and label all other images as negative examples for the pneumonia detection task. We randomly split the entire dataset into 80% training, and 20% validation.

To avoid this, a patient must only belong to the test or train/val set, not both.

This is called group splitting:



Creating forward-lag features is another way leaky variables get created unintentionally at times:

	f_0	f_1		f_0	f_1	$\text{lag-}f_1$
15-Jan 2024		a		15-Jan 2024		b
15-Feb 2024		b		15-Feb 2024		c
15-Mar 2024		c		15-Mar 2024		d
15-Apr 2024		d		15-Apr 2024		e
15-May 2024		e		15-May 2024		

A horizontal arrow points from the original table to the second table, labeled '30-day lag on f_1 '.

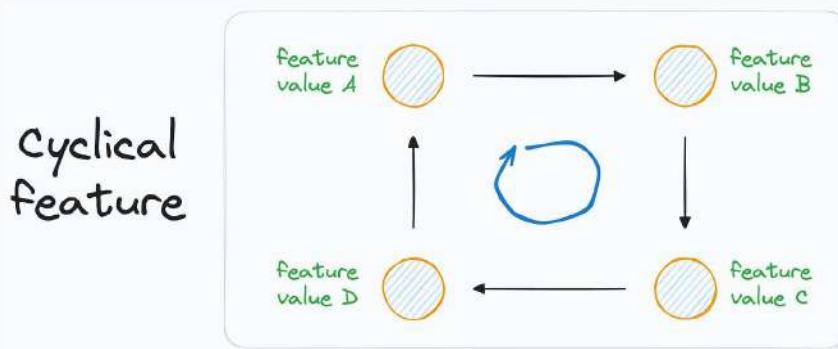
Let's get into more detail about the issue with random splitting below.

Cyclical feature encoding

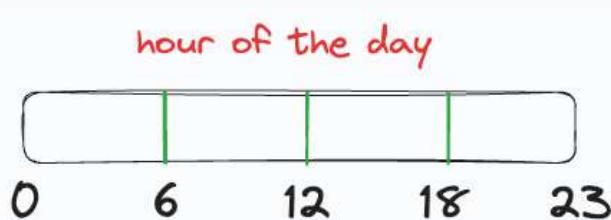
In typical machine learning datasets, we mostly find features that progress from one value to another: For instance:

- Numerical features like age, income, transaction amount, etc.
- Categorical features like t-shirt size, income groups, age groups, etc.

However, there is one more type of feature, which, in most cases, deserves special feature engineering effort but is often overlooked. These are cyclical features, i.e., features with a recurring pattern (or cycle).



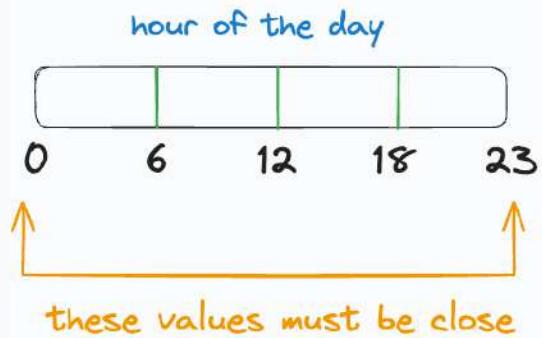
Unlike other features that progress continuously (or have no inherent order), cyclical features exhibit periodic behavior and repeat after a specific interval. For instance, the hour-of-the-day, the day-of-the-week, and the month-of-an-year are all common examples of cyclical features. Talking specifically about, say, the hour-of-the-day, its value can range between 0 to 23:



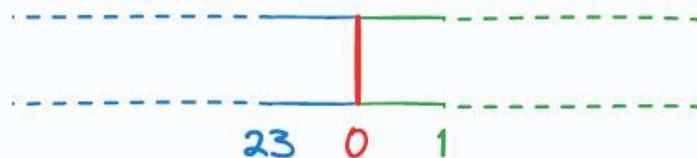
If we DON'T consider this as a cyclical feature and don't utilize appropriate feature engineering techniques, we will lose some really critical information.

To understand better, consider this:

Realistically speaking, the values “23” and “0” must be close to each other in our “ideal” feature representation of the hour-of-the-day.



Moreover, the distance between “0” and “1” must be the same as the distance between “23” and “0”.



$$\text{distance}(23, 0) = \text{distance}(0, 1)$$

However, standard representation does not fulfill these properties. Thus, the value “23” is far from “0”. In fact, the distance property isn’t satisfied either.

Now, think about it for a second. Intuitively speaking, don’t you think this feature deserves special feature engineering, i.e., one that preserves the inherent natural property?

I am sure you do!

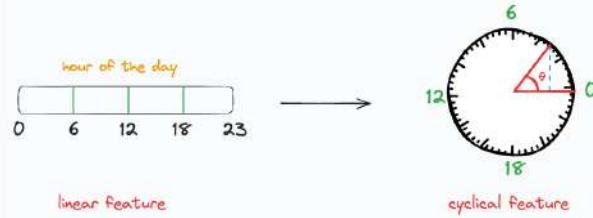
Let’s understand how we typically do it.

Cyclical feature encoding

One of the most common techniques to encode such a feature is using trigonometric functions, specifically, sine and cosine. These are helpful because sine and cosine are periodic, bounded, and defined for all real values.

Of course, even other trigonometric functions are also periodic, but they are also undefined for some values, like, $\tan(\pi/2)$.

For instance, consider representing the linear hour-of-the-day feature as a cyclical feature:



The central angle (2π) represents 24 hours. Thus, the linear feature values can be easily converted into cyclical features as follows:

```
import numpy as np

df['sin_hour'] = np.sin(df['hour'] * 2 * np.pi / 24)
df['cos_hour'] = np.cos(df['hour'] * 2 * np.pi / 24)
```

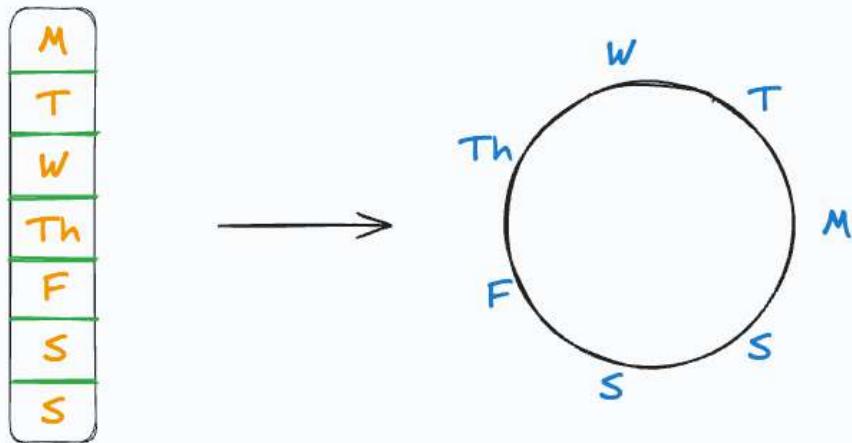
The benefit of doing this is how neatly the engineered feature satisfies the properties we discussed earlier:

$$|\sin((0 - 23) * \frac{2\pi}{24})| = |\sin((0 - 1) * \frac{2\pi}{24})|$$

$$|\cos((0 - 23) * \frac{2\pi}{24})| = |\cos((0 - 1) * \frac{2\pi}{24})|$$

As depicted above, the distance between the cyclical feature representation of “23” and “0” is the same as the distance between “0” and “1”. The standard linear representation of the hour-of-the-day feature, however, violates this property, which results in loss of information...

...or rather, I should say that the standard linear representation of the hour-of-the-day feature results in an underutilization of information, which the model can benefit from. Had it been the day-of-the-week instead, the central angle (2π) must have represented 7 days.



The same idea can be extended to all sorts of cyclical features you may find in your dataset:

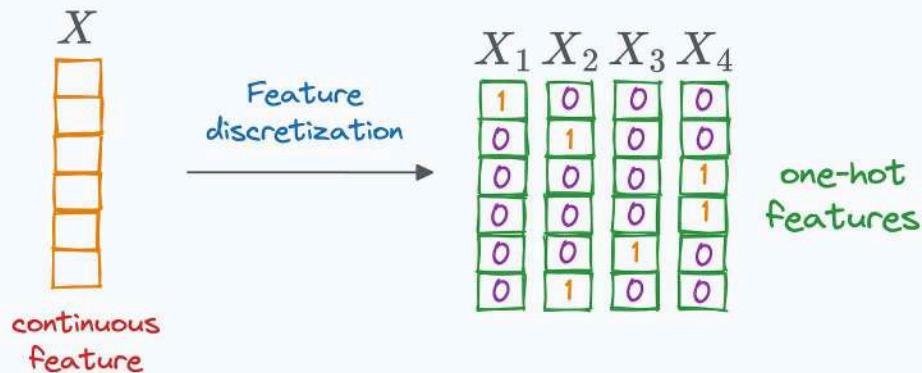
- Wind direction, if represented categorically, will go in this order: N, NE, E, SE, S, SW, W, NW, and then back to N.
- Phases of the moon, like new moon, first quarter, full moon, and last quarter, can be represented as categories with a cyclical order.
- Seasons, such as spring, summer, fall, and winter, are categorical features with a cyclical pattern as they repeat annually.

The point is that as you will inspect the dataset features, you will intuitively know which features are cyclical and which are not.

Typically, the model will find it easier to interpret the engineered features and utilize them in modeling the dataset accurately.

Feature Discretization

During model development, one of the techniques that many don't experiment with is feature discretization. As the name suggests, the idea behind discretization is to transform a continuous feature into discrete features.



Why, when, and how would you do that? Let's understand in this chapter.

Motivation

My rationale for using feature discretization has almost always been simple: "It just makes sense to discretize a feature."

For instance, consider your dataset has an age feature:

$$\hat{y} = w_1 * \text{age} + w_2 * x_2 + w_3 * x_3 + \dots$$

In many use cases, like understanding spending behavior based on transaction history, such continuous variables are better understood when they are discretized into meaningful groups → youngsters, adults, and seniors.

For instance, say we model this transaction dataset without discretization. This would result in some coefficients for each feature, which would tell us the influence of each feature on the final prediction.

$$\hat{y} = w_1 * \text{age} + w_2 * x_2 + w_3 * x_3 + \dots$$

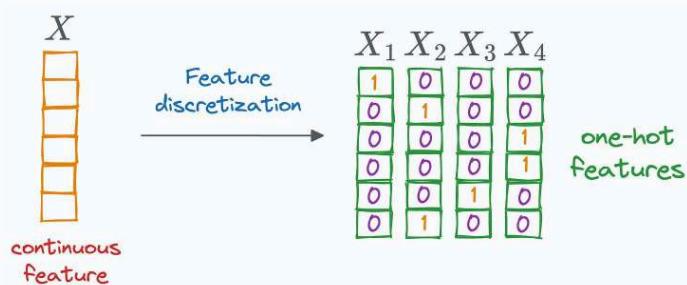
But if you think again, in our goal of understanding spending behavior, are we really interested in learning the correlation between exact age and spending behavior?

It makes very little sense to do that. Instead, it makes more sense to learn the correlation between different age groups and spending behavior.

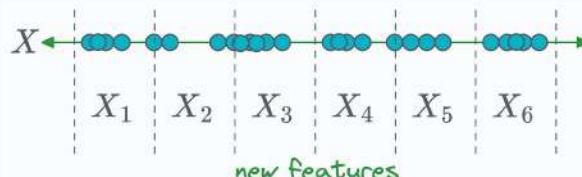
As a result, discretizing the age feature can potentially unveil much more valuable insights than using it as a raw feature.

2 common techniques

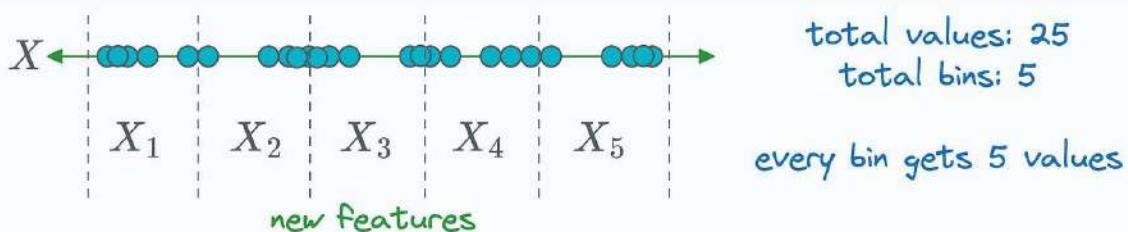
Now that we understand the rationale, there are 2 techniques that are widely preferred.



One way of discretizing features involves decomposing a feature into equally sized bins.

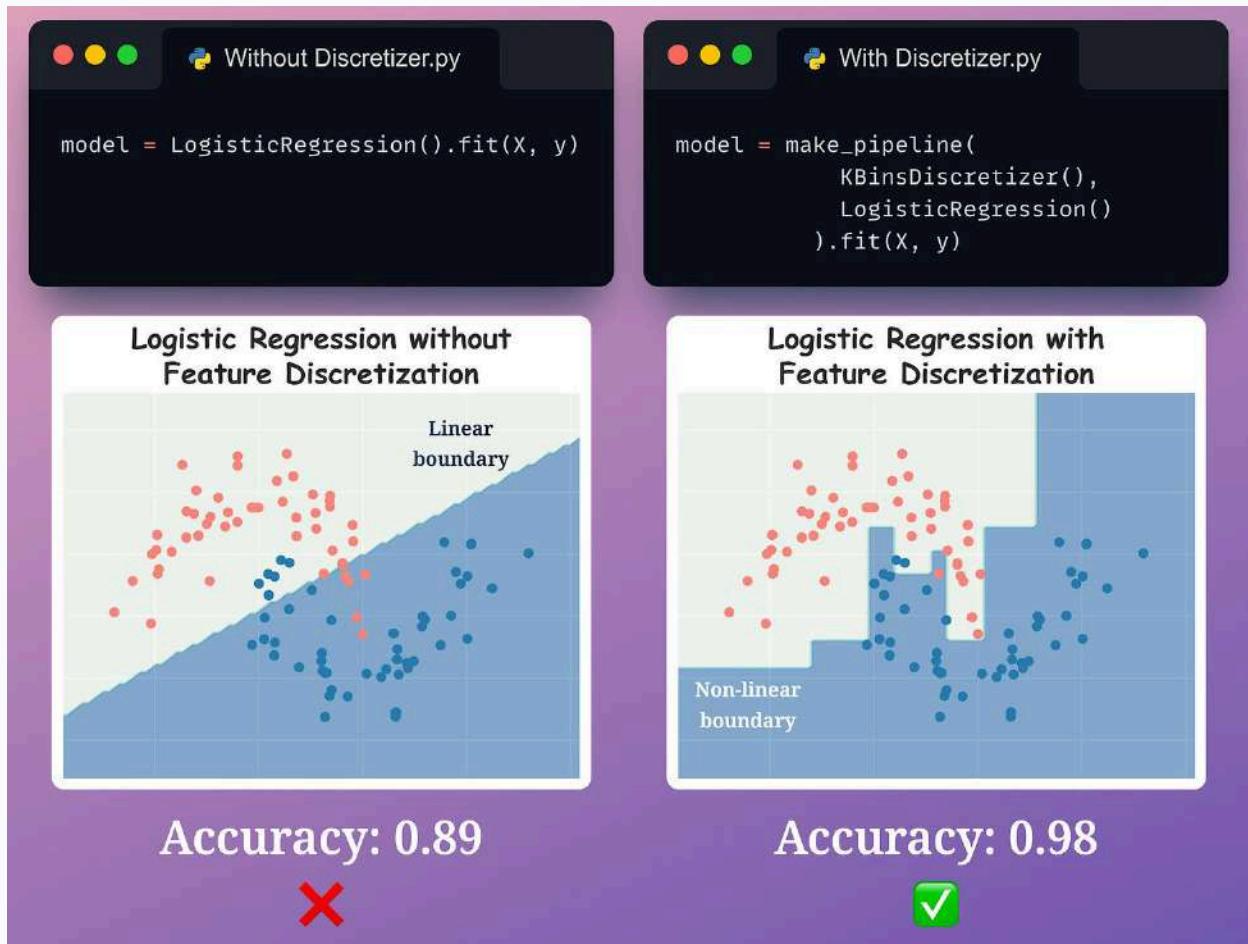


Another technique involves decomposing a feature into equal frequency bins:



After that, the discrete values are one-hot encoded.

One advantage of feature discretization is that it enables non-linear behavior even though the model is linear. This can potentially lead to better accuracy, which is also evident from the image below:



A linear model with feature discretization results in a:

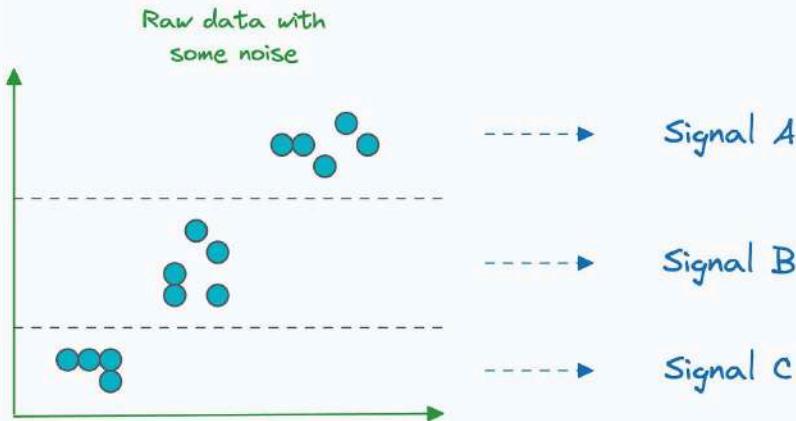
- non-linear decision boundary.
- better test accuracy.

So, in a way, we get to use a simple linear model but still get to learn non-linear patterns.

Another advantage of discretizing continuous features is that it helps us improve the signal-to-noise ratio.

Simply put, “signal” refers to the meaningful or valuable information in the data.

Binning a feature helps us mitigate the influence of minor fluctuations, which are often mere noise.



Each bin acts as a means of “smoothing” out the noise within specific data segments.

Before I conclude, do remember that feature discretization with one-hot encoding increases the number of features → thereby increasing the data dimensionality.

And typically, as we progress towards higher dimensions, data become more easily linearly separable. Thus, feature discretization can lead to overfitting.

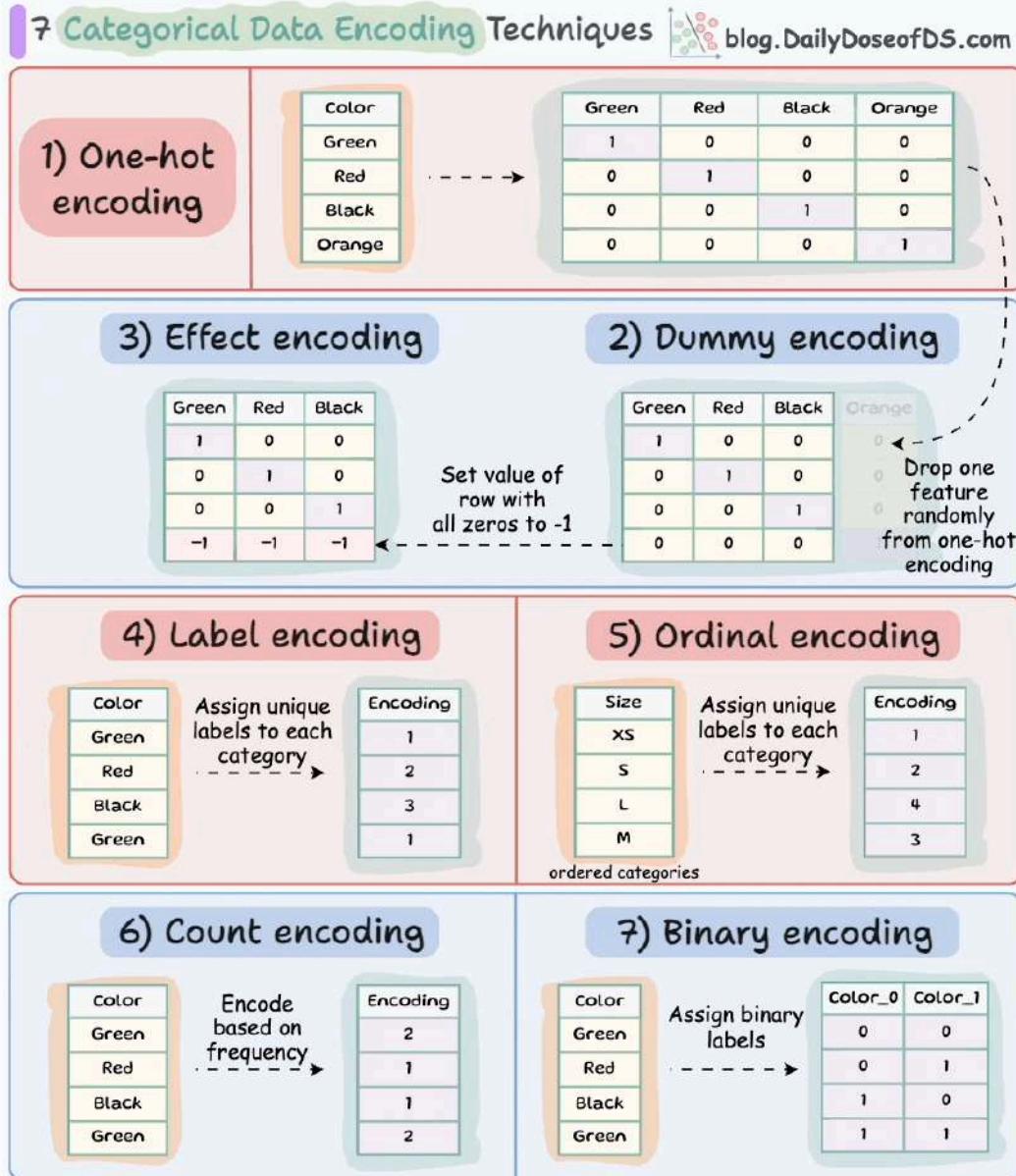
To avoid this, don’t overly discretize all features. Instead, use it when it makes intuitive sense, as we saw earlier.

Of course, its utility can vastly vary from one application to another, but at times, I have found that:

- Discretizing geospatial data like latitude and longitude can be useful.
- Discretizing age/weight-related data can be useful.
- Features that are typically constrained between a range makes sense, like savings/income (practically speaking), etc.

7 Categorical Data Encoding Techniques

Here are 7 ways to encode categorical features:



We covered them in detail here: <https://bit.ly/3LkfVq5>.

One-hot encoding:

- Each category is represented by a binary vector of 0s and 1s.
- Each category gets its own binary feature, and only one of them is "hot" (set to 1) at a time, indicating the presence of that category.
- Number of features = Number of unique categorical labels

Dummy encoding:

- Same as one-hot encoding but with one additional step.
- After one-hot encoding, we drop a feature randomly.
- We do this to avoid the dummy variable trap (discussed in [this chapter](#)).
- Number of features = Number of unique categorical labels - 1

Effect encoding:

- Similar to dummy encoding but with one additional step.
- Alter the row with all zeros to -1.
- This ensures that the resulting binary features represent not only the presence or absence of specific categories but also the contrast between the reference category and the absence of any category.
- Number of features = Number of unique categorical labels - 1.

Label encoding:

- Assign each category a unique label.
- Label encoding introduces an inherent ordering between categories, which may not be the case.
- Number of features = 1.

Ordinal encoding:

- Similar to label encoding — assign a unique integer value to each category.
- The assigned values have an inherent order, meaning that one category is considered greater or smaller than another.
- Number of features = 1.

Count encoding:

- Also known as frequency encoding.
- Encodes categorical features based on the frequency of each category.
- Thus, instead of replacing the categories with numerical values or binary representations, count encoding directly assigns each category with its corresponding count.
- Number of features = 1.

Binary encoding:

- Combination of one-hot encoding and ordinal encoding.
- It represents categories as binary code.
- Each category is first assigned an ordinal value, and then that value is converted to binary code.
- The binary code is then split into separate binary features.
- Useful when dealing with high-cardinality categorical features (or a high number of features) as it reduces the dimensionality compared to one-hot encoding.
- Number of features = $\log(n)$ (in base 2).

While these are some of the most popular techniques, do note that these are not the only techniques for encoding categorical data.

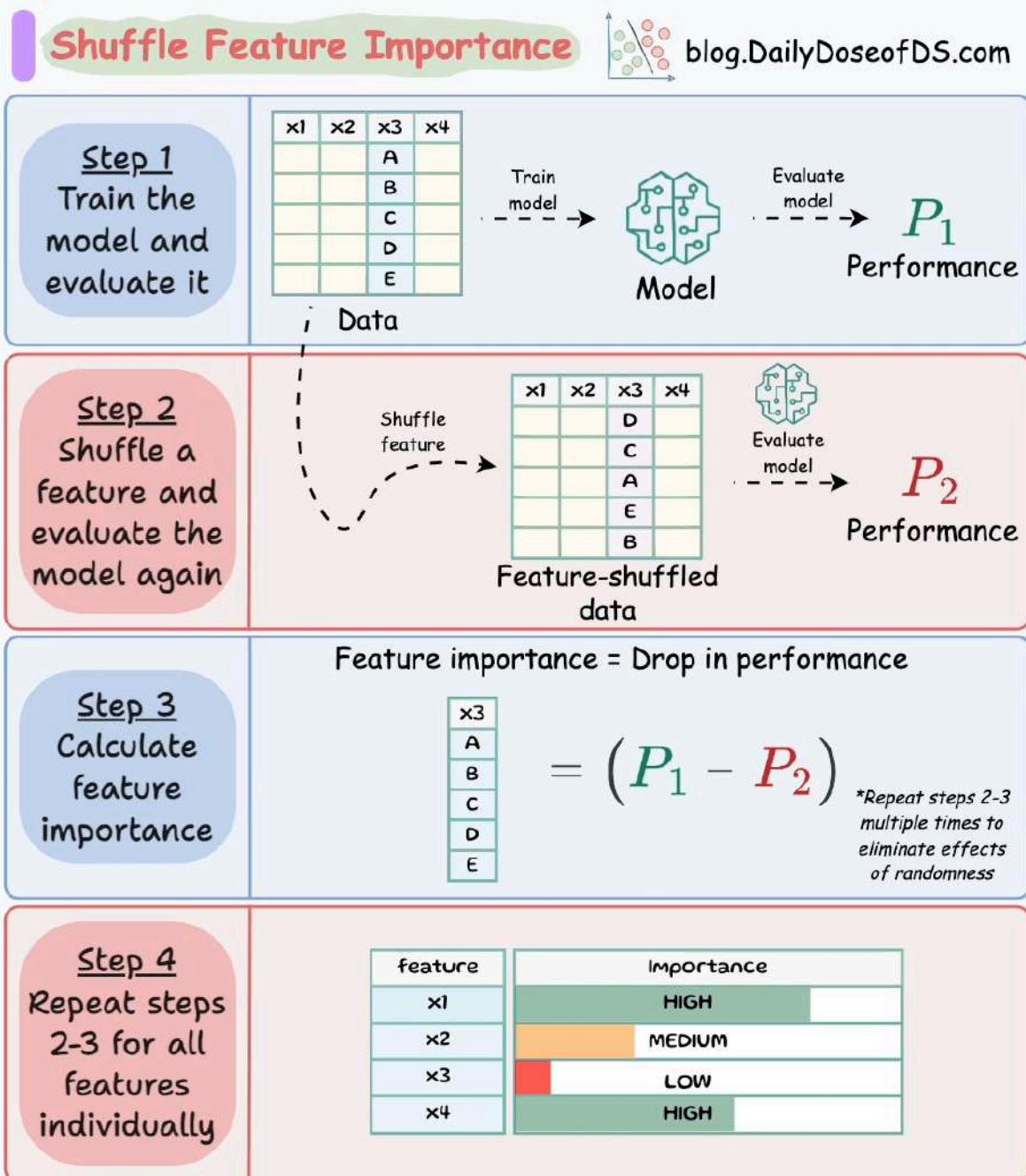
You can try plenty of techniques with the category-encoders library:

<https://pypi.org/project/category-encoders>.

Shuffle Feature Importance

I often find “Shuffle Feature Importance” to be a handy and intuitive technique to measure feature importance.

As the name suggests, it observes how shuffling a feature influences the model performance. The visual below illustrates this technique in four simple steps:

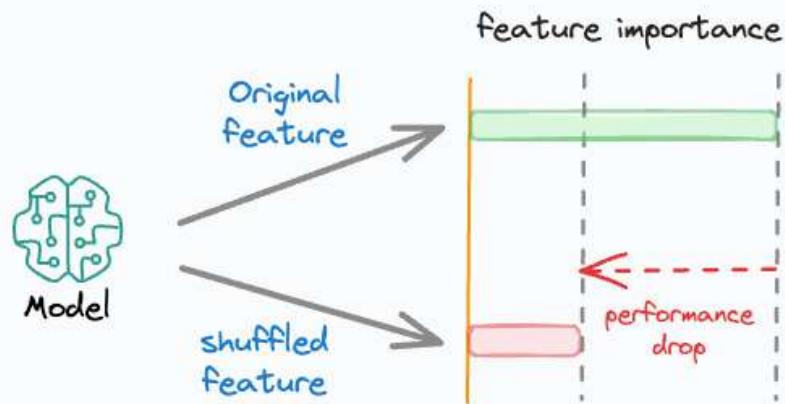


Here's how it works:

- Train the model and measure its performance $\rightarrow P_1$.
- Shuffle one feature randomly and measure performance again $\rightarrow P_2$ (model is NOT trained again).
- Measure feature importance using performance drop = $(P_1 - P_2)$.
- Repeat for all features.

This makes intuitive sense as well, doesn't it?

Simply put, if we randomly shuffle just one feature and everything else stays the same, then the performance drop will indicate how important that feature is.



- If the performance drop is low \rightarrow This means the feature has a very low influence on the model's predictions.
- If the performance drop is high \rightarrow This means that the feature has a very high influence on the model's predictions.

Do note that to eliminate any potential effects of randomness during feature shuffling, it is recommended to:

- Shuffle the same feature multiple times
- Measure average performance drop.

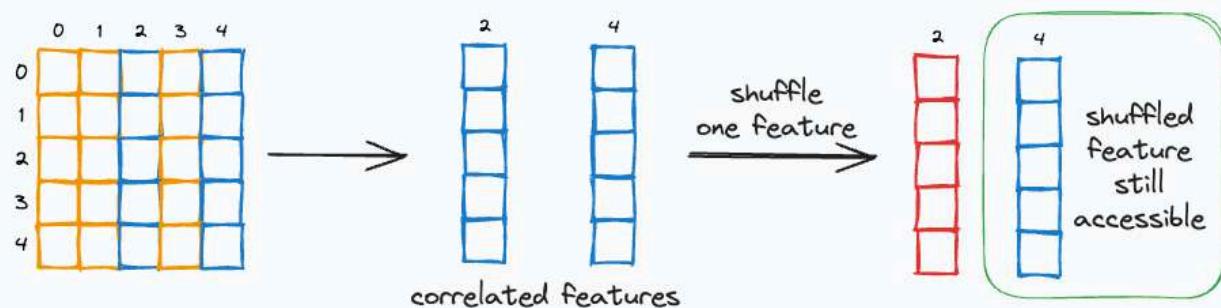
A few things that I love about this technique are:

- It requires no repetitive model training. Just train the model once and measure the feature importance.
- It is pretty simple to use and quite intuitive to interpret.
- This technique can be used for all ML models that can be evaluated.

Of course, there is one caveat as well.

Say two features are highly correlated, and one of them is permuted/shuffled.

In this case, the model will still have access to the feature through its correlated feature.



This will result in a lower importance value for both features.

One way to handle this is to cluster highly correlated features and only keep one feature from each cluster.

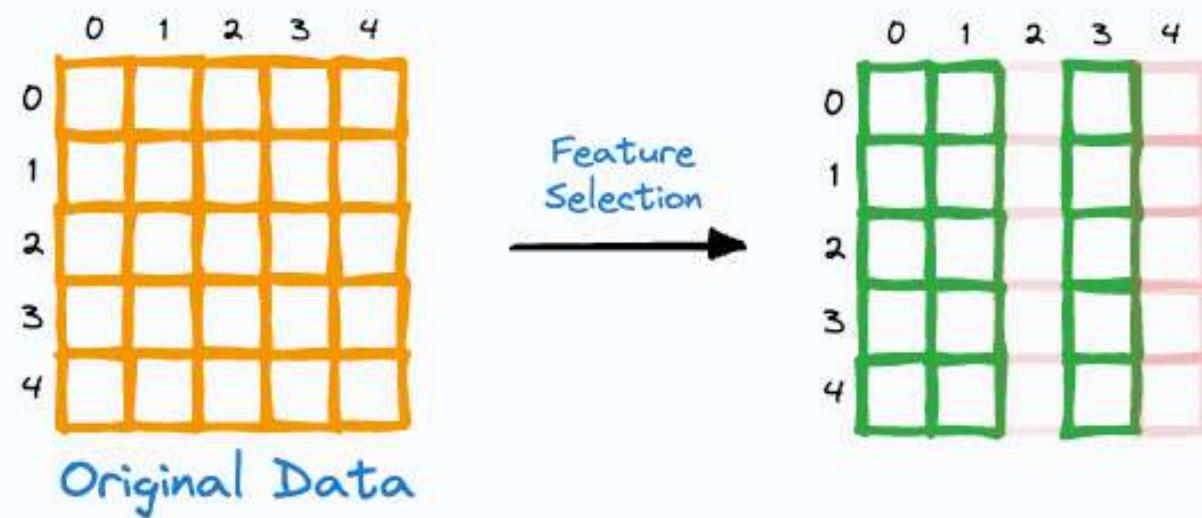
The Probe Method for Feature Selection

Real-world ML development is all about achieving a sweet balance between speed, model size, and performance.

One common way to:

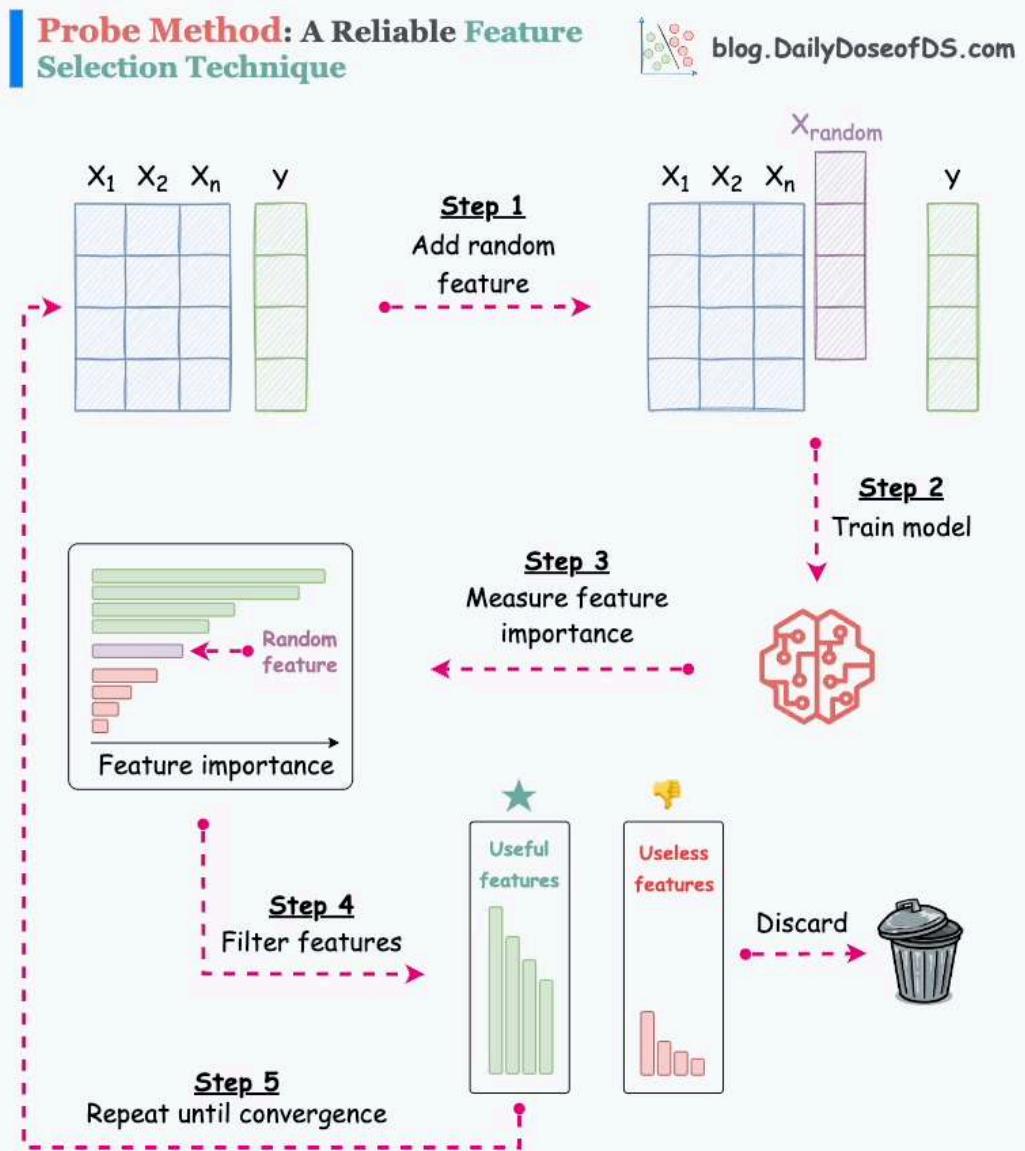
- Improve speed,
- Reduce size, and
- Maintain (or minimally degrade) performance...

...is by using feature selection. The idea is to select the most useful subset of features from the dataset.



While there are many many methods for feature selection, I have often found the “Probe Method” to be pretty reliable, practical and intuitive to use.

The image below depicts how it works:



- Step 1) Add a random feature (noise).
- Step 2) Train a model on the new dataset.
- Step 3) Measure feature importance (can use shuffle feature importance)
- Step 4) Discard original features that rank below the random feature.
- Step 5) Repeat until convergence.

This whole idea makes intuitive sense as well.

More specifically, if a feature's importance is ranked below a random feature, it is probably a useless feature for the model.

This can be especially useful in cases where we have plenty of features, and we wish to discard those that don't contribute to the model.

Of course, one shortcoming is that when using the Probe Method, we must train multiple models:

1. Train the first model with the random feature and discard useless features.
2. Keep training new models until the random feature is ranked as the least important feature (although typically, convergence does not result in plenty of models).
3. Train the final model without the random feature.

Nonetheless, the approach can be quite useful to reduce model complexity.

Regression

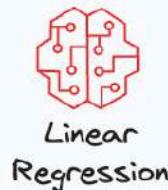
Why Mean Squared Error (MSE)?

Say you wish to train a linear regression model. We know that we train it by minimizing the squared error:

$$\text{Squared Error} = \sum_i^n \frac{(y_i - \theta^T x_i)^2}{n}$$

But have you ever wondered why we specifically use the squared error?

Train linear regression
using Mean Squared Error



But why?



See, many functions can potentially minimize the difference between observed and predicted values. But of all the possible choices, what is so special about the squared error?

In my experience, people often say:

- Squared error is differentiable. That is why we use it as a loss function.
WRONG.
- It is better than using absolute error as squared error penalizes large errors more. WRONG.

Sadly, each of these explanations are incorrect.

But approaching it from a probabilistic perspective helps us understand why the squared error is the ideal choice.

Let's understand in this chapter.

In linear regression, we predict our target variable y using the inputs X as follows:

$$y_i = \theta^T x_i + \epsilon_i$$

Here, epsilon is an error term that captures the random noise for a specific data point (i).

We assume the noise is drawn from a Gaussian distribution with zero mean based on the central limit theorem:

$$\epsilon \sim N(0, \sigma^2)$$

Thus, the probability of observing the error term can be written as:

$$p(\epsilon_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\epsilon_i^2}{2\sigma^2}\right)$$

Substituting the error term from the linear regression equation, we get:

$$\epsilon_i = \theta^T x_i - y_i$$

$$p(y_i|x_i; \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \theta^T x_i)^2}{2\sigma^2}\right)$$

For a specific set of parameters θ , the above tells us the probability of observing a data point (i).

Next, we can define the likelihood function as follows:

$$L(\theta) = p(y|X; \theta)$$

It means that by varying θ , we can fit a distribution to the observed data and quantify the likelihood of observing it.

We further write it as a product for individual data points because we assume all observations are independent.

$$L(\text{colorful dots}) = p(\text{blue}) \cdot p(\text{yellow}) \cdot p(\text{red}) \cdot p(\text{purple}) \cdot p(\text{cyan})$$

independent observations *product of observing individual observations*

Thus, we get:

$$\begin{aligned} L(\theta) &= \prod_i^n p(y_i|x_i; \theta) \\ &= \prod_i^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \theta^T x_i)^2}{2\sigma^2}\right) \end{aligned}$$

Since the log function is monotonic, we use the log-likelihood and maximize it. This is called maximum likelihood estimation (MLE).

$$\begin{aligned} \log(L(\theta)) &= \log\left(\prod_i^n p(y_i|x_i; \theta)\right) \\ &= \log\left(\prod_i^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \theta^T x_i)^2}{2\sigma^2}\right)\right) \end{aligned}$$

Simplifying, we get:

$$\begin{aligned}
 \log(L(\theta)) &= \sum_i^n \log\left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \theta^T x_i)^2}{2\sigma^2}\right)\right) \\
 &\quad \text{distribute the logarithm} \\
 &= \sum_i^n \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) - \frac{(y_i - \theta^T x_i)^2}{2\sigma^2} \\
 &\quad \text{distribute the summation} \\
 &= \boxed{n \cdot \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)} - \boxed{\sum_i^n \frac{(y_i - \theta^T x_i)^2}{2\sigma^2}}
 \end{aligned}$$

constant

To reiterate, the objective is to find the θ that maximizes the above expression. But the first term is independent of θ . Thus, maximizing the above expression is equivalent to minimizing the second term. And if you notice closely, it's precisely the squared error.

$$\frac{1}{2\sigma^2} \sum_i^n (y_i - \theta^T x_i)^2$$

Thus, you can maximize the log-likelihood by minimizing the squared error. And this is the origin of least-squares in linear regression. See, there's clear proof and reasoning behind using squared error as a loss function in linear regression.

Nothing comes from thin air in machine learning.

Sklearn Linear Regression Has No Hyperparameters

Almost all ML models we work with have some hyperparameters, such as:

- Learning rate
- Regularization
- Layer size (for neural network), etc.

But as shown in the image below, why don't we see any hyperparameter in Sklearn's Linear Regression implementation?

The screenshot shows the documentation for the `sklearn.linear_model.LinearRegression` class. At the top, there is a logo of a scatter plot with a regression line and the text "DailyDoseofDS.com". Below the title, the class definition is shown:

```
class sklearn.linear_model.LinearRegression(*, fit_intercept=True, copy_X=True, n_jobs=None, positive=False)
```

Next to the code, there is a link "[source]". A large red arrow points from the text "No Hyperparameters!" to the parameters section. The parameters are listed as follows:

- Parameters:**
 - fit_intercept : bool, default=True**
Whether to calculate the intercept for this model. If set to False, no intercept will be used in calculations (i.e. data is expected to be centered).
 - copy_X : bool, default=True**
If True, X will be copied; else, it may be overwritten.
 - n_jobs : int, default=None**
The number of jobs to use for the computation. This will only provide speedup in case of sufficiently large problems, that is if firstly `n_targets > 1` and secondly `X` is sparse or if `positive` is set to True. None means 1 unless in a `joblib.parallel_backend` context. -1 means using all processors. See `Glossary` for more details.
 - positive : bool, default=False**
When set to True, forces the coefficients to be positive. This option is only supported for dense arrays.

A note at the bottom states: "New in version 0.24."

It must have learning rate as a hyperparameter, right?

To understand the reason why it has no hyperparameters, we first need to learn that the Linear Regression can model data in two different ways:

1. Gradient Descent (which many other ML algorithms use for optimization):
 - It is a stochastic algorithm, i.e., involves some randomness.
 - It finds an approximate solution using optimization.
 - It has hyperparameters.
2. Ordinary Least Square (OLS):
 - It is a deterministic algorithm. Thus, if run multiple times, it will always converge to the same weights.
 - It always finds the optimal solution.
 - It has no hyperparameters.

Now, instead of the typical gradient descent approach, Sklearn's Linear Regression class implements the OLS method.

sklearn.linear_model.LinearRegression

```
class sklearn.linear_model.LinearRegression(*, fit_intercept=True, copy_X=True, n_jobs=None,
positive=False)
[source]
```

Ordinary least squares Linear Regression.

OLS

That is why it has no hyperparameters.

How does OLS work?

With OLS, the idea is to find the set of parameters (Θ) such that:

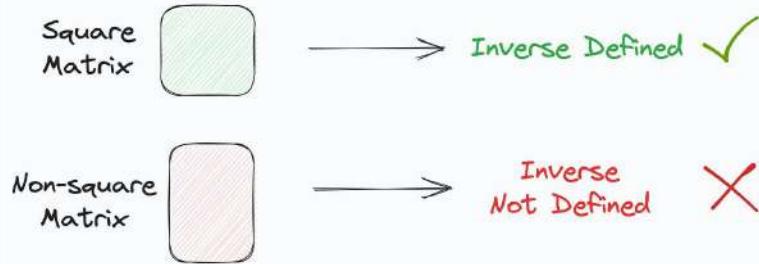
$$X\theta = Y$$

- X: input data with dimensions (n,m).
- Θ : parameters with dimensions (m,1).
- y: output data with dimensions (n,1).
- n: number of samples.
- m: number of features.

One way to determine the parameter matrix Θ is by multiplying both sides of the equation with the inverse of X , as shown below:

$$\theta = X^{-1}Y$$

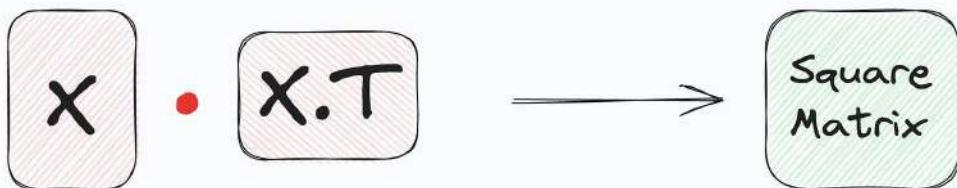
But because X might be a non-square matrix, its inverse may not be defined.



To resolve this, first, we multiply with the transpose of X on both sides, as shown below:

$$X^T X \theta = X^T y$$

This makes the product of X with its transpose a square matrix.



The obtained matrix, being square, can be inverted (provided it is non-singular).

Next, we take the collective inverse of the product to get the following:

$$\theta = (X^T X)^{-1} X^T y$$

It's clear that the above definition has:

- No hyperparameters.
- No randomness. Thus, it will always return the same solution, which is also optimal.

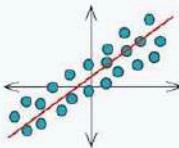
This is precisely what the Linear Regression class of Sklearn implements.

To summarize, it uses the OLS method instead of gradient descent.

That is why it has no hyperparameters.

Of course, do note that there is a significant tradeoff between run time and convenience when using OLS vs. gradient descent.

This is also clear from the time-complexity table we discussed in an earlier chapter.

Time Complexity of 10 Most Popular ML Algorithms		Training	Inference
	Linear Regression (OLS)	$O(nm^2 + m^3)$	$O(m)$
	Linear Regression (SGD)	$O(n_{epoch}nm)$	$O(m)$

As depicted above, the run-time of OLS is cubically related to the number of features (m).

Thus, when we have many features, it may not be a good idea to use the `LinearRegression()` class. Instead, use the `SGDRegressor()` class from Sklearn.

Of course, the good thing about `LinearRegression()` class is that it involves no hyperparameter tuning.

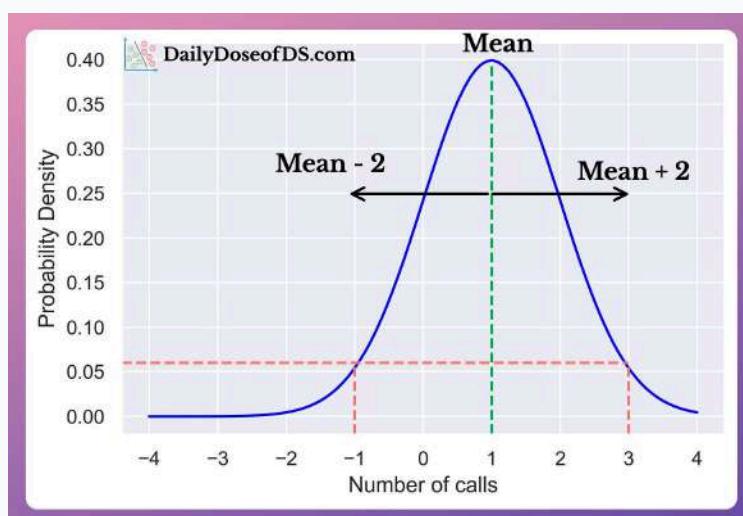
Thus, when we use OLS, we trade run-time for finding an optimal solution without hyperparameter tuning.

Poisson Regression vs. Linear Regression

Linear regression comes with its own set of challenges/assumptions. For instance, after modeling, the output can be negative for some inputs.

But this may not make sense at times — predicting the number of goals scored, number of calls received, etc. Thus, it is clear that it cannot model count (or discrete) data.

Furthermore, in linear regression, residuals are expected to be normally distributed around the mean. So the outcomes on either side of the mean ($m-x$, $m+x$) are equally likely.



For instance:

- if the expected number (mean) of calls received is 1...
- ...then, according to linear regression, receiving 3 calls (1+2) is just as likely as receiving -1 (1-2) calls. (This relates to the concept of prediction intervals, which we covered in the [earlier chapter](#).)
- But in this case, a negative prediction does not make any sense.

Thus, if the above assumptions do not hold, linear regression won't help.

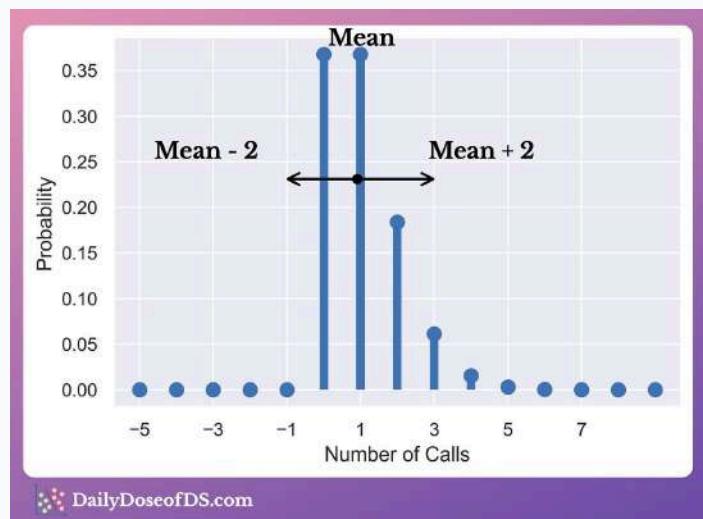
Instead, in this specific case, what you may need is Poisson regression.

Poisson regression is more suitable if your response (or outcome) is count-based. It assumes that the response comes from a Poisson distribution.

$$P(x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

It is a type of generalized linear model (GLM) that is used to model count data. It works by estimating a Poisson distribution parameter (λ), which is directly linked to the expected number of events in a given interval.

Contrary to linear regression, in Poisson regression, residuals may follow an asymmetric distribution around the mean (λ). Hence, outcomes on either side of the mean ($\lambda-x, \lambda+x$) are NOT equally likely.



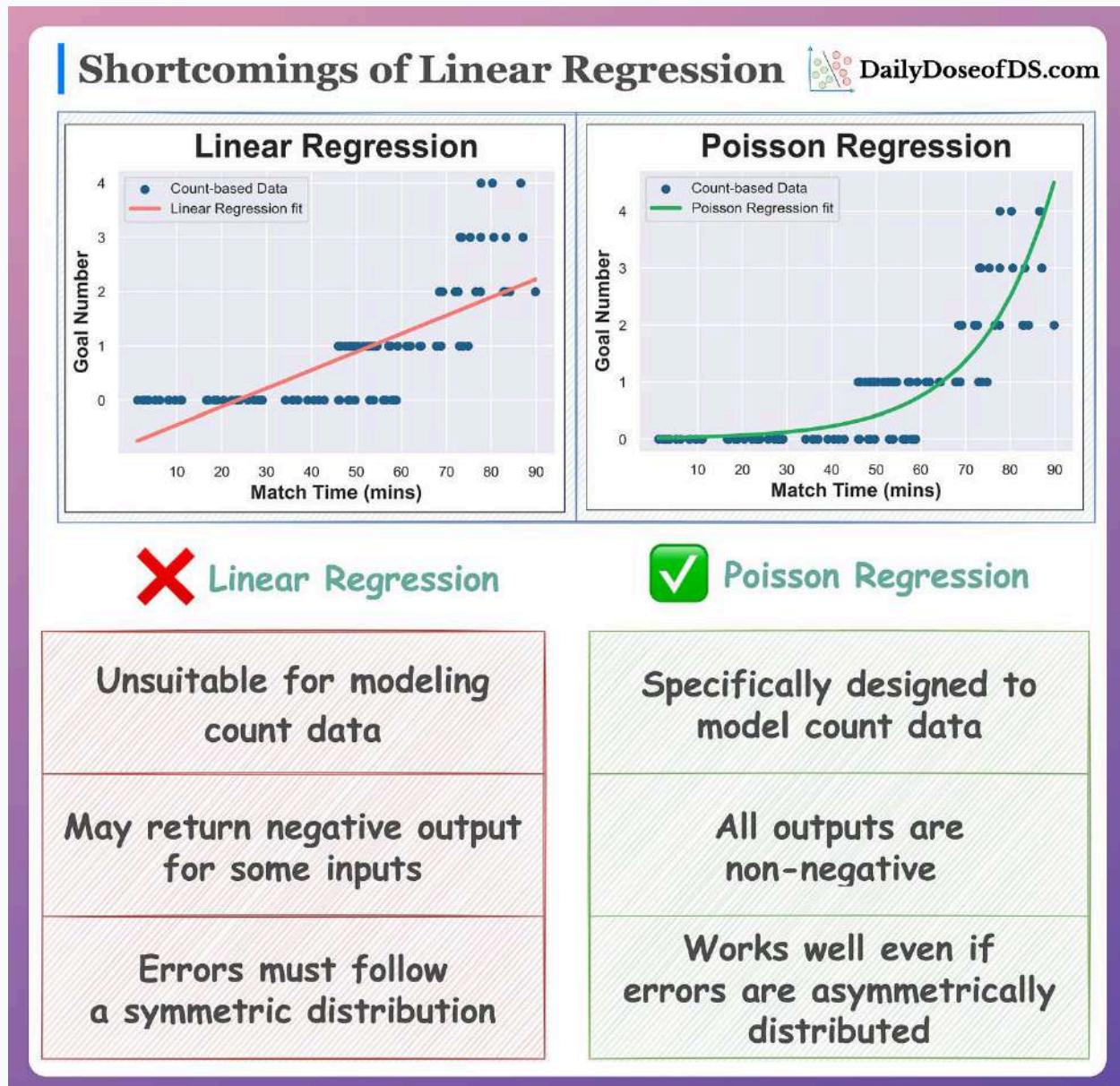
For instance:

- if the expected number (mean) of calls received is 1...
- ...then, according to Poisson regression, it is possible to receive 3 (1+2) calls, but it is impossible to receive -1 (1-2) calls.
- This is because its outcome is also non-negative.

The regression fit is mathematically defined as follows:

$$\ln(E(Y|X)) = w_1X_1 + w_2X_2 + \cdots + w_nX_n$$

The following visual neatly summarizes this post:



We shall continue this discussion in the next chapter.

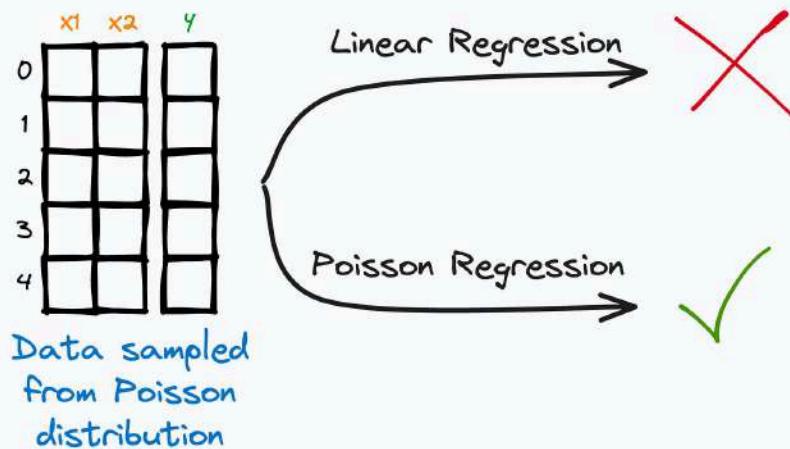
How to Build Linear Models?

In this chapter, I will help you cultivate what I think is one of the MOST overlooked and underappreciated skills in developing linear models.

I can guarantee that harnessing this skill will give you a lot of clarity and intuition in the modeling stages.

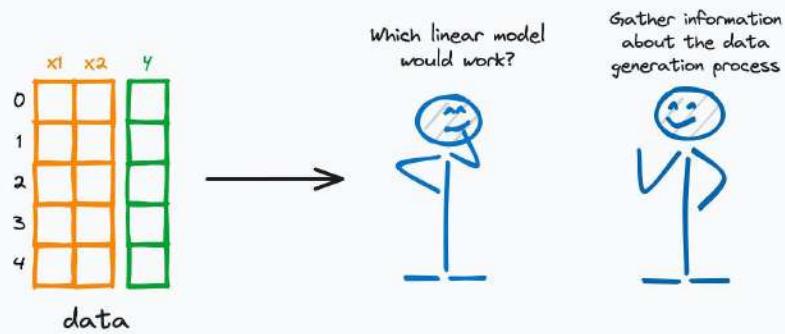
To begin, understand that the Poisson regression we discussed above is no magic.

It's just that, in our specific use case, the data generation process didn't perfectly align with what linear regression is designed to handle. In other words, earlier when we trained a linear regression model, we inherently assumed that the data was sampled from a normal distribution. But that was not true in this Poisson regression case.

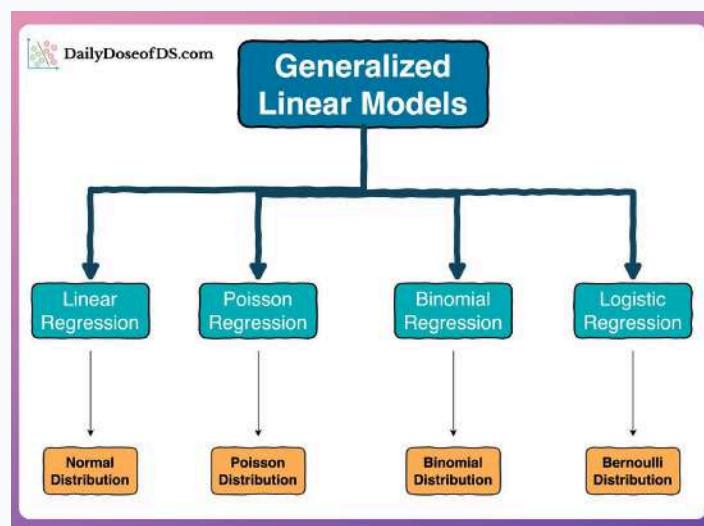


Instead, it came from a Poisson distribution, which is why Poisson regression worked better. Thus, the takeaway is that whenever you train linear models, always always and always think about the data generation process. It goes like this:

- Okay, I have this data.
- I want to fit a linear model through it.
- What information do I get about the data generation process that can help me select an appropriate linear model?



You'd start appreciating the importance of data generation when you realize that literally every member of the generalized linear model family stems from altering the data generation process.



For instance:

- If the data generation process involves a Normal distribution → you get linear regression.
- If the data has only positive integers in the response variable, maybe it came from a Poisson distribution → and this gives us Poisson regression. This is precisely what we discussed yesterday.
- If the data has only two targets — 0 and 1, maybe it was generated using Bernoulli distribution → and this gives rise to logistic regression.
- If the data has finite and fixed categories (0, 1, 2,...n), then this hints towards Binomial distribution → and we get Binomial regression.

See...

Every linear model makes an assumption and is then derived from an underlying data generation process.

Linear Regression \dashrightarrow Assumes Normal distribution

Logistic Regression \dashrightarrow Assumes Bernoulli distribution

Poisson Regression \dashrightarrow Assumes Poisson distribution

Binomial Regression \dashrightarrow Assumes Binomial distribution

Thus, developing a habit of holding for a second and thinking about the data generation process will give you so much clarity in the modeling stages.

I am confident this will help you get rid of that annoying and helpless habit of relentlessly using a specific sklearn algorithm without truly knowing why you are using it.

Consequently, you'd know which algorithm to use and, most importantly, why.

This improves your credibility as a data scientist and allows you to approach data science problems with intuition and clarity rather than hit-and-trial.

In fact, once you understand the data generation process, you will automatically get to know about most of the assumptions of that specific linear model.

Dummy Variable Trap

With one-hot encoding, we introduce a big problem in the data.

When we one-hot encode categorical data, we unknowingly introduce perfect multicollinearity.

Multicollinearity arises when two or more features can predict another feature. In this case, as the sum of one-hot encoded features is always 1, it leads to perfect multicollinearity.

	isSmall	isMedium	isLarge
Small	1	0	0
Medium	0	1	0
Large	0	0	1

	isSmall	isMedium	isLarge	Sum
Small	1	0	0	1
Medium	0	1	0	1
Large	0	0	1	1

$\text{isSmall} + \text{isMedium} + \text{isLarge} = 1$

This is often called the Dummy Variable Trap. It is bad because the model has redundant features. Moreover, the regression coefficients aren't reliable in the presence of multicollinearity.

So how to resolve this?

The solution is simple. Drop any arbitrary feature from the one-hot encoded features.

This instantly mitigates multicollinearity and breaks the linear relationship which existed before.

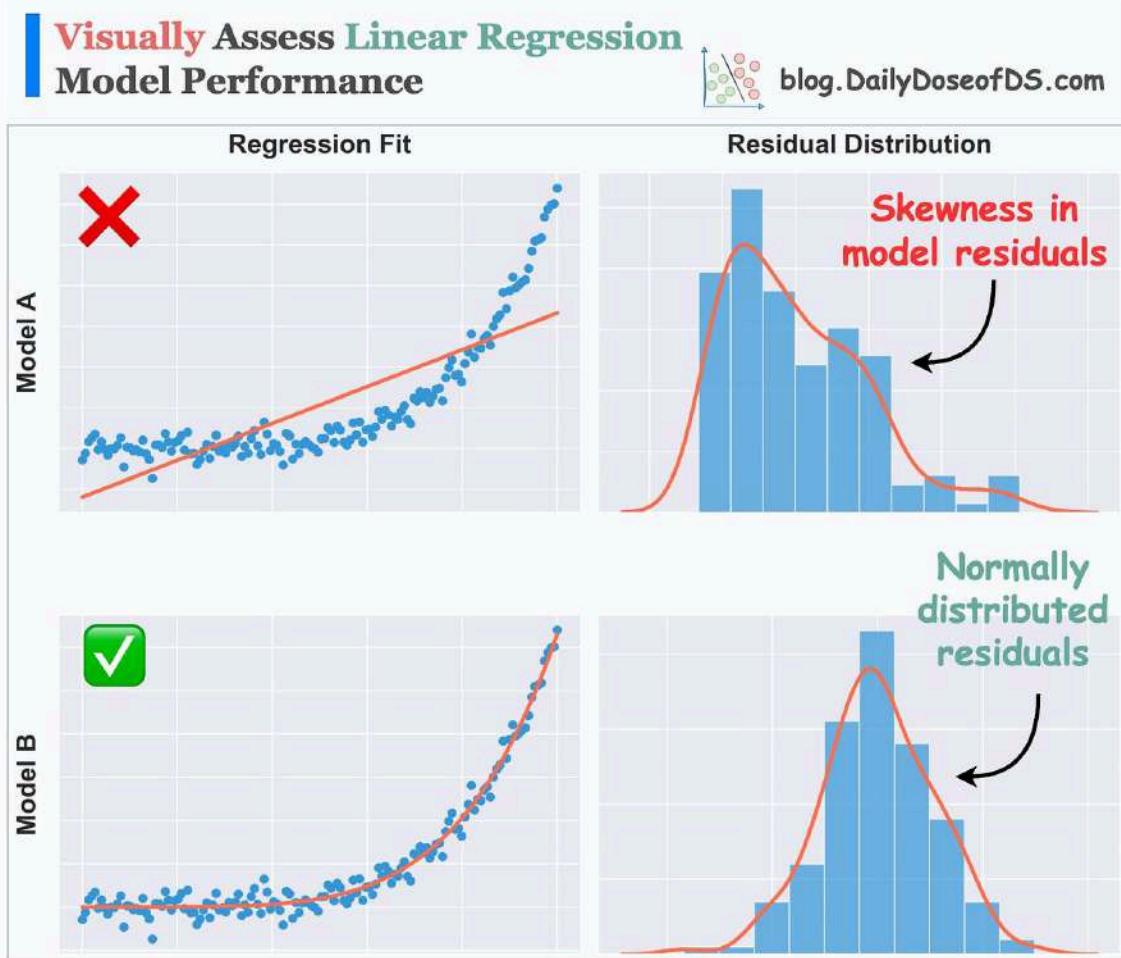
	isSmall	isMedium	Sum
Small	1	0	1
Medium	0	1	1
Large	0	0	0

$\text{isSmall} + \text{isMedium} \neq 1$

Visually Assess Linear Regression Performance

Linear regression assumes that the model residuals ($= \text{actual} - \text{predicted}$) are normally distributed. If the model is underperforming, it may be due to a violation of this assumption.

Here, I often use a residual distribution plot to verify this and determine the model's performance. As the name suggests, this plot depicts the distribution of residuals ($= \text{actual} - \text{predicted}$), as shown below:



A good residual plot will:

- Follow a normal distribution
- NOT reveal trends in residuals

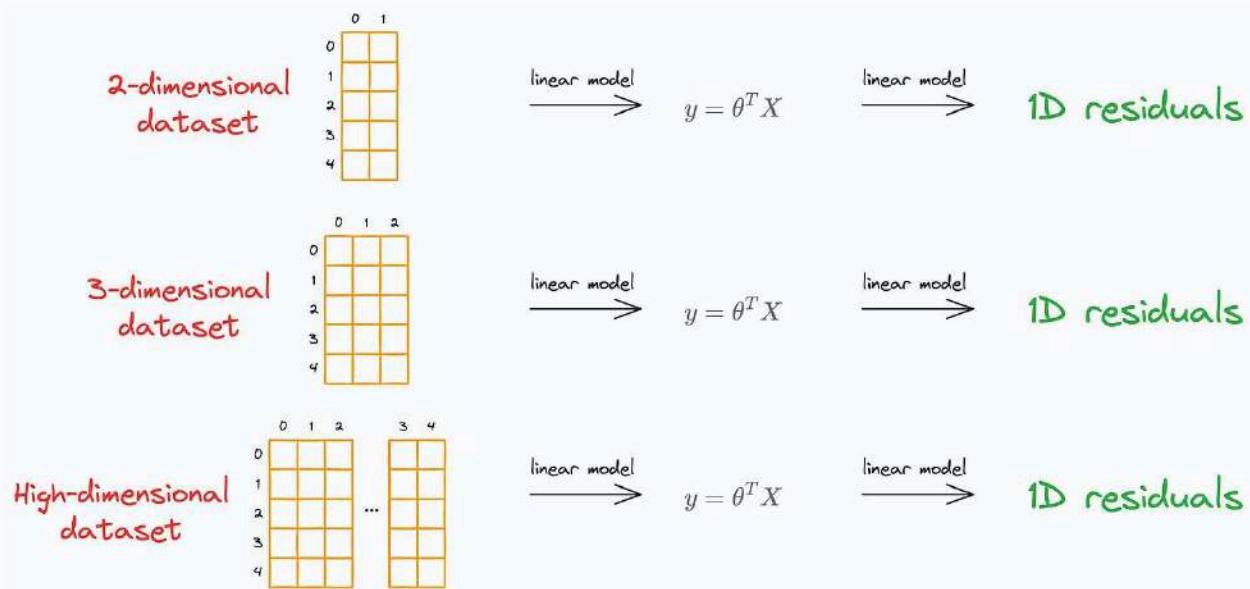
A bad residual plot will:

- Show skewness
- Reveal patterns in residuals

Thus, the more normally distributed the residual plot looks, the more confident we can be about our model. This is especially useful when the regression line is difficult to visualize, i.e., in a high-dimensional dataset.

Why?

Because a residual distribution plot depicts the distribution of residuals, which is always one-dimensional.



Thus, it can be plotted and visualized easily.

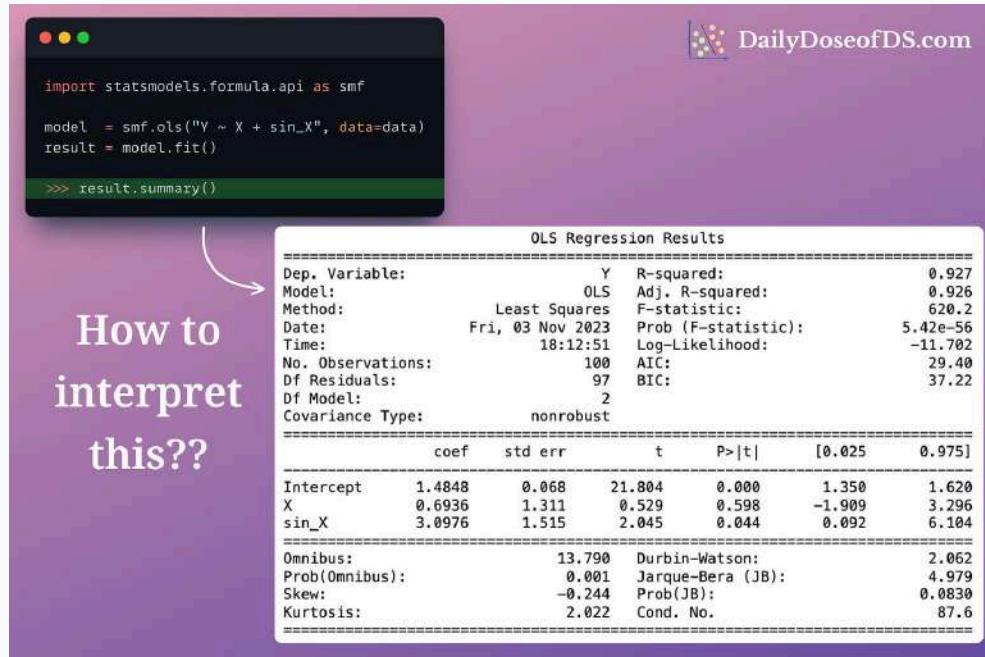
Of course, this was just about validating one assumption — the normality of residuals.

However, linear regression relies on many other assumptions, which must be tested as well.

Statsmodel provides a pretty comprehensive report for this, which we shall discuss in the next chapter.

Statsmodel Regression Summary

Statsmodel provides one of the most comprehensive summaries for regression analysis.



The screenshot shows a Jupyter Notebook cell with Python code to fit a linear regression model and print its summary. The output is an 'OLS Regression Results' table with various statistical measures and coefficient details. A large green arrow points from the text 'How to interpret this??' to the 'Dep. Variable' section of the table.

```

import statsmodels.formula.api as smf
model = smf.ols("Y ~ X + sin_X", data=data)
result = model.fit()
>>> result.summary()

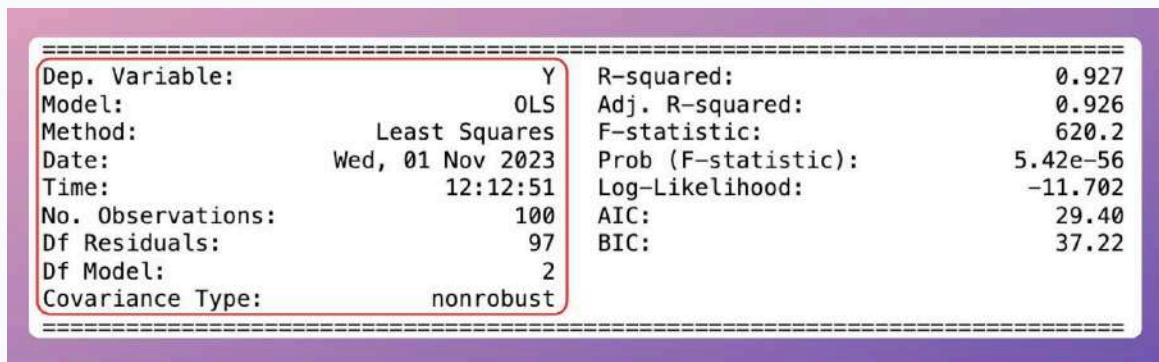
```

OLS Regression Results						
Dep. Variable:	Y	R-squared:	0.927			
Model:	OLS	Adj. R-squared:	0.926			
Method:	Least Squares	F-statistic:	620.2			
Date:	Fri, 03 Nov 2023	Prob (F-statistic):	5.42e-56			
Time:	18:12:51	Log-Likelihood:	-11.702			
No. Observations:	100	AIC:	29.40			
Df Residuals:	97	BIC:	37.22			
Df Model:	2					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
Intercept	1.4848	0.068	21.804	0.000	1.350	1.620
X	0.6936	1.311	0.529	0.598	-1.909	3.296
sin_X	3.0976	1.515	2.045	0.044	0.092	6.104
Omnibus:	13.790	Durbin-Watson:	2.062			
Prob(Omnibus):	0.001	Jarque-Bera (JB):	4.979			
Skew:	-0.244	Prob(JB):	0.0830			
Kurtosis:	2.022	Cond. No.	87.6			

Yet, I have seen so many people struggling to interpret the critical model details mentioned in this report. In this chapter, let's understand the entire summary support provided by statsmodel and why it is so important.

Section #1

The first column of the first section lists the model's settings (or config). This part has nothing to do with the model's performance.



The screenshot shows the first section of the regression summary, which contains the model's configuration parameters. These include the dependent variable (Y), model type (OLS), method (Least Squares), date, time, number of observations, degrees of freedom for residuals and model, and covariance type (nonrobust). The entire section is highlighted with a red box.

Dep. Variable:	Y	R-squared:	0.927
Model:	OLS	Adj. R-squared:	0.926
Method:	Least Squares	F-statistic:	620.2
Date:	Wed, 01 Nov 2023	Prob (F-statistic):	5.42e-56
Time:	12:12:51	Log-Likelihood:	-11.702
No. Observations:	100	AIC:	29.40
Df Residuals:	97	BIC:	37.22
Df Model:	2		
Covariance Type:	nonrobust		

- Dependent variable: The variable we are predicting.
- Model and Method: We are using OLS to fit a linear model.
- Date and time: You know it.
- No. of observations: The dataset's size.
- Df residuals: The degrees of freedom associated with the residuals. It is essentially the number of data points minus the number of parameters estimated in the model (including intercept term).
- Df Model: This represents the degrees of freedom associated with the model. It is the number of predictors, 2 in our case — X and \sin_X .

If your data has categorical features, statsmodel will one-hot encode them. But in that process, it will drop one of the one-hot encoded features.



This is done to avoid the dummy variable trap, which we discussed in an earlier chapter ([this chapter](#)).

- Covariance type: This is related to the assumptions about the distribution of the residual.
 - In linear regression, we assume that the residuals have a constant variance (homoscedasticity).
 - We use “nonrobust” to train a model under that assumption.

In the second column, statsmodel provides the overall performance-related details:

Dep. Variable:	Y	R-squared:	0.927
Model:	OLS	Adj. R-squared:	0.926
Method:	Least Squares	F-statistic:	620.2
Date:	Wed, 01 Nov 2023	Prob (F-statistic):	5.42e-56
Time:	12:12:51	Log-Likelihood:	-11.702
No. Observations:	100	AIC:	29.48
Df Residuals:	97	BIC:	37.22
Df Model:	2		
Covariance Type:	nonrobust		

- R-squared: The fraction of original data variability captured by the model.

$$R^2 = \frac{\text{Variability captured by the Model}}{\text{Variability in the original dataset}}$$

- For instance, in this case, 0.927 means that the current model captures 92.7% of the original variability in the training data.
- Statsmodel reports R2 on the input data, so you must not overly optimize for it. If you do, it will lead to overfitting.
- Adj. R-squared:
 - It is somewhat similar to R-squared, but it also accounts for the number of predictors (features) in the model.
 - The problem is that R-squared always increases as we add more features.
 - So even if we add totally irrelevant features, R-squared will never decrease. Adj. R-squared penalizes this behavior of R-squared.
- F-statistic and Prob (F-statistic):
 - These assess the overall significance of a regression model.
 - They compare the estimated coefficients by OLS with a model whose all coefficients (except for the intercept) are zero.



- F-statistic tests whether the independent variables collectively have any effect on the dependent variable or not.
- Prob (F-statistic) is the associated p-value with the F-statistic.
- A small p-value (typically less than 0.05) indicates that the model as a whole is statistically significant.
- This means that at least one independent variable has a significant effect on the dependent variable.
- Log-Likelihood:

- This tells us the log-likelihood that the given data was generated by the estimated model.
- The higher the value, the more likely the data was generated by this model.
- AIC and BIC:
 - Like adjusted R-squared, these are performance metrics to determine goodness of fit while penalizing complexity.
 - Lower AIC and BIC values indicate a better fit.

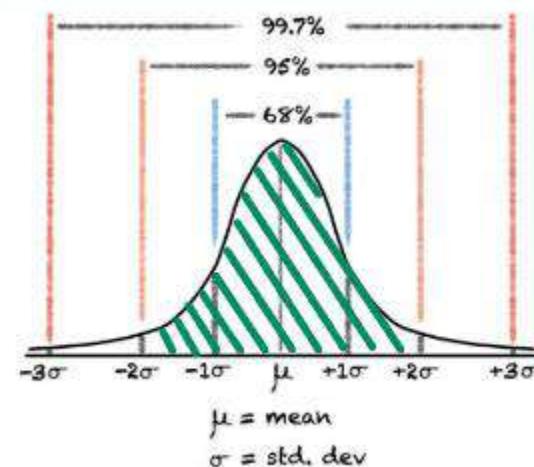
Section #2

The second section provides details related to the features:

	coef	std err	t	P> t	[0.025	0.975]
Intercept	1.4848	0.068	21.804	0.000	1.350	1.620
X	0.6936	1.311	0.529	0.598	-1.909	3.296
sin_X	3.0976	1.515	2.045	0.044	0.092	6.104

- coef: The estimated coefficient for a feature.
- t and P>|t|:
 - Earlier, we used F-statistic to determine the statistical significance of the model as a whole.
 - t-statistic is more granular on that front as it determines the significance of every individual feature.
 - P>|t| is the associated p-value with the t-statistic.
 - A small p-value (typically less than 0.05) indicates that the feature is statistically significant.
 - For instance, the feature “X” has a p-value of ~0.6. This suggests that there is a 60% chance that the feature “X” has no effect on “Y”.
- [0.025, 0.975] and std err:

- See, the coefficients we have obtained from the model are just estimates. They may not be absolute true coefficients of the process that generated the data.
- Thus, the estimated parameters are subject to uncertainty, aren't they?
- Note: The width of the interval [0.025, 0.975] is 0.95 → or 95%. This constitutes the area between 2 standard deviations from the mean in a normal distribution.



- A 95% confidence interval provides a range of values within which you can be 95% confident that the true value of the parameter lies.

	coef	std err	t	P> t	[0.025	0.975]
Intercept	1.4848	0.068	21.804	0.000	1.350	1.620
X	0.6936	1.311	0.529	0.598	-1.909	3.296
sin_X	3.0976	1.515	2.045	0.044	0.092	6.104

- For instance, the interval for sin_X is (0.092, 6.104). So although the estimated coefficient is 3.09, we can be 95% confident that the true coefficient lies in the range (0.092, 6.104).

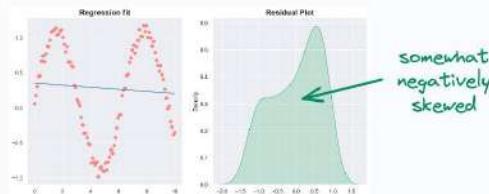
Section #3

Details in the last section of the report test the assumptions of linear regression.

Omnibus:	13.790	Durbin-Watson:	2.062
Prob(Omnibus):	0.001	Jarque-Bera (JB):	4.979
Skew:	-0.244	Prob(JB):	0.0830
Kurtosis:	2.022	Cond. No.	87.6

- Omnibus and Prob(Omnibus):
 - They test the normality of the residuals.
 - Omnibus value of zero means residuals are perfectly normal.
 - Prob(Omnibus) is the corresponding p-value.
 - In this case, Prob(Omnibus) is 0.001. This means there is a 0.1% chance that the residuals are normally distributed.
- Skew and Kurtosis:
 - They also provide information about the distribution of the residuals.
 - Skewness measures the asymmetry of the distribution of residuals.
 - Zero skewness means perfect symmetry.
 - Positive skewness indicates a distribution with a long right tail. This indicates a concentration of residuals on lower values. Good to check for outliers in this case.

Negative skewness indicates a distribution with a long left tail. This is mostly indicative of poor features. For instance, consider fitting a sin curve with a linear feature (X). Most residuals will be high, resulting in negative skewness.



- Durbin-Watson:
 - This measures autocorrelation between residuals.
 - Autocorrelation occurs when the residuals are correlated, indicating that the error terms are not independent.
 - But linear regression assumes that residuals are not correlated.
 - The Durbin-Watson statistic ranges between 0 and 4.
 - A value close to 2 indicates no autocorrelation.
 - Values closer to 0 indicate positive autocorrelation.
 - Values closer to 4 indicate negative autocorrelation.
- Jarque-Bera (JB) and Prob(JB):
 - They solve the same purpose as Omnibus and Prob(Omnibus) — measuring the normality of residuals.
- Condition Number:
 - This tests multicollinearity.
 - Multicollinearity occurs when two features are correlated, or two or more features determine the value of another feature.
 - A standalone value for Condition Number can be difficult to interpret so here's how I use it:
 - Add features one by one to the regression model and notice any spikes in the Condition Number.

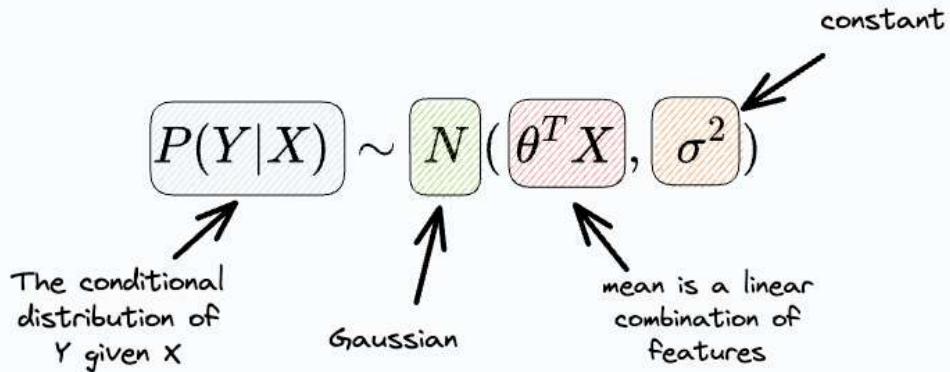
As discussed above, every section of this report has its importance:

- The first section tells us about the model's config, the overall performance of the model, and its statistical significance.
- The second section tells us about the statistical significance of individual features, the model's confidence in finding the true coefficient, etc.
- The last section lets us validate the model's assumptions, which are immensely critical to linear regression's performance.

Now you know how to interpret the entire regression summary from statsmodel.

Generalized Linear Models (GLMs)

A linear regression model is undeniably an extremely powerful model, in my opinion. However, it makes some strict assumptions about the type of data it can model, as depicted below.

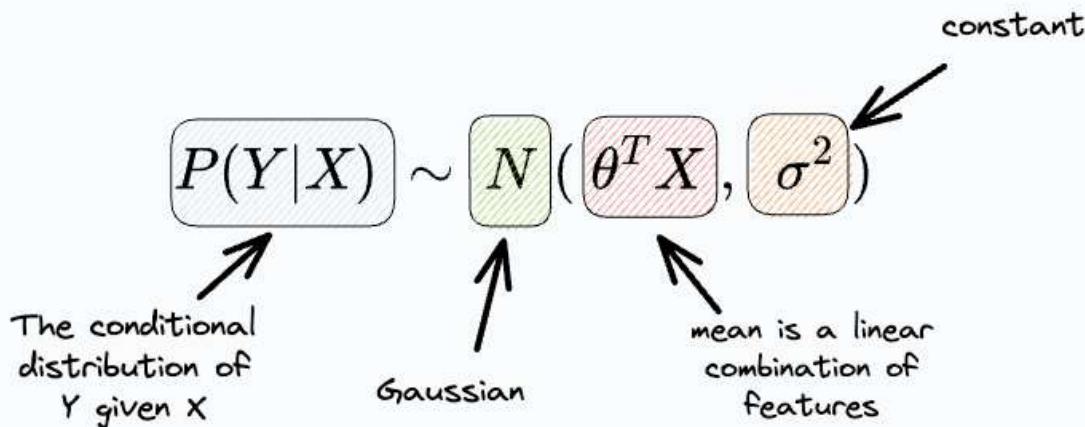


These conditions often restrict its applicability to data situations that do not obey the above assumptions. That is why being aware of its extensions is immensely important.

Generalized linear models (GLMs) precisely do that. They relax the assumptions of linear regression to make linear models more adaptable to real-world datasets.

Why GLMs?

Linear regression is pretty restricted in terms of the kind of data it can model. For instance, its assumed data generation process looks like this:



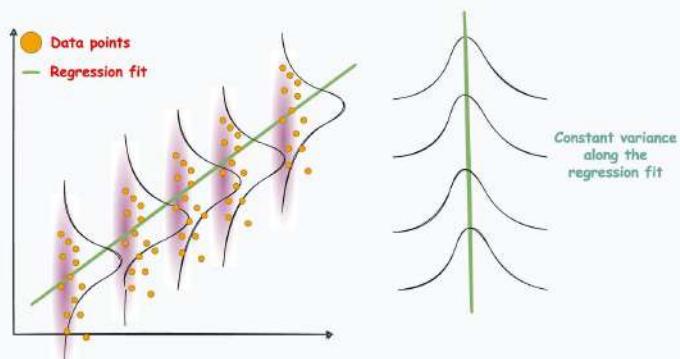
The assumed data generation process of linear regression

- Firstly, it assumes that the conditional distribution of Y given X is a Gaussian.
- Next, it assumes a very specific form for the mean of the above Gaussian. It says that the mean should always be a linear combination of the features (or predictors).

$$\mu(x) = \sum_{j=1}^K \theta_j \cdot x_j$$

mean is a linear combination of features

- Lastly, it assumes a constant variance for the conditional distribution $P(Y|X)$ across all levels of X. A graphical way of illustrating this is as follows:



These conditions often restrict its applicability to data situations that do not obey the above assumptions.

In other words, nothing stops real-world datasets from violating these assumptions.

In fact, in many scenarios, the data might exhibit complex relationships, heteroscedasticity (varying variance), or even follow entirely different distributions altogether.

Yet, if we intend to build linear models, we should formulate better algorithms that can handle these peculiarities.

Generalized linear models (GLMs) precisely do that.

They relax the assumptions of linear regression to make linear models more adaptable to real-world datasets.

More specifically, they consider the following:

- What if the distribution isn't normal but some other distribution?

$$P(Y|X) \sim N(\theta^T X, \sigma^2)$$

What if this is not normal distribution?

- What if X has a more sophisticated relationship with the mean?

$$\mu(x) = \sum_{i=1}^K \theta_j \cdot x_j$$

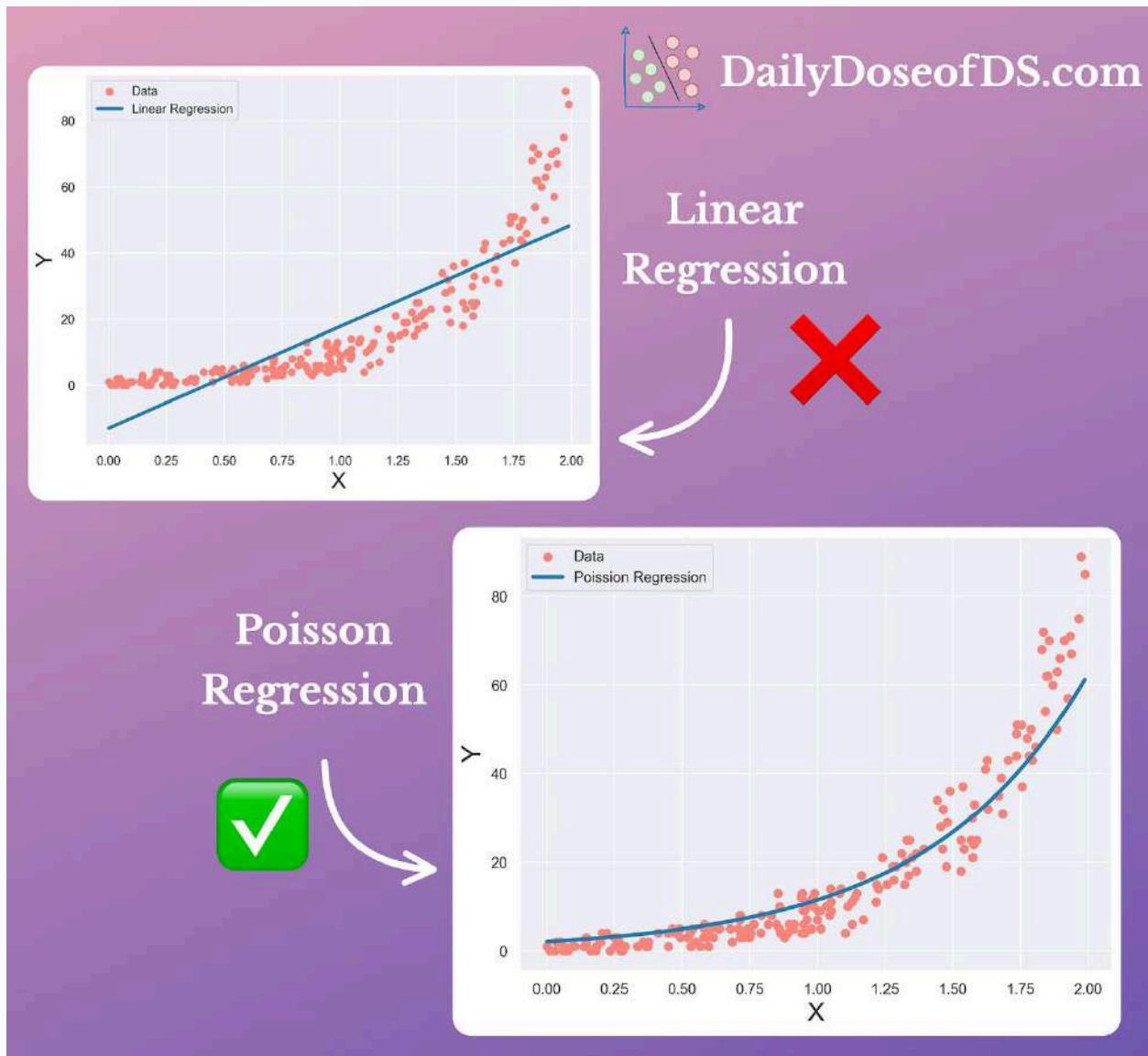
What if this does not hold?

- What if the variance varies with X?

$$P(Y|X) \sim N(\theta^T X, \sigma^2)$$

What if this is not constant?

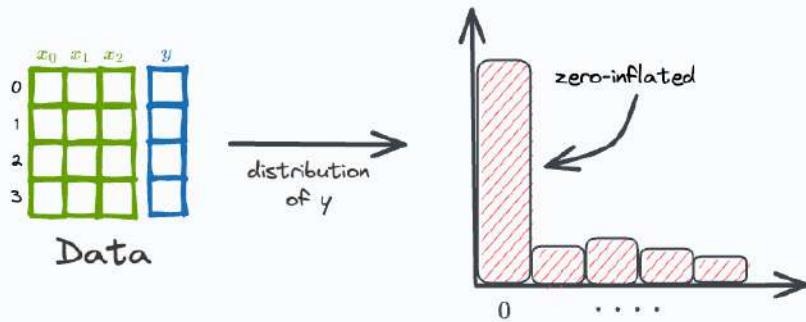
The effectiveness of a specific GLM — Poisson regression (which we discussed in an [earlier chapter](#)) over linear regression is evident from the image below:



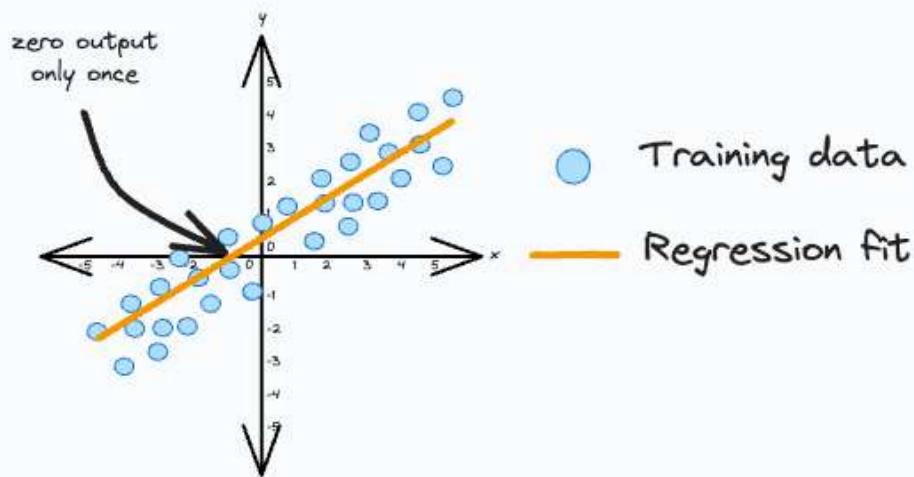
- Linear regression assumes the data is drawn from a Gaussian, when in reality, it isn't. Hence, it underperforms.
- Poisson regression adapts its regression fit to a non-Gaussian distribution. Hence, it performs significantly better.

Zero-inflated Regression

The target variable of typical regression datasets is somewhat evenly distributed. But, at times, the target variable may have plenty of zeros. Such datasets are called zero-inflated datasets.



They may raise many problems during regression modeling. This is because a regression model can not always predict exact “zero” values when, ideally, it should. For instance, consider simple linear regression. The regression line will output exactly “zero” only once (if it has a non-zero slope).

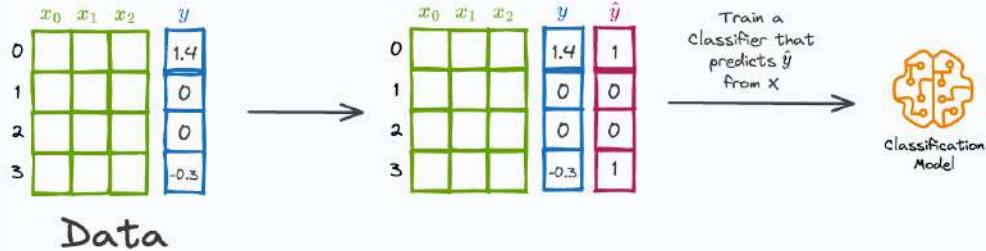


This issue persists not only in higher dimensions but also in complex models like neural nets for regression.

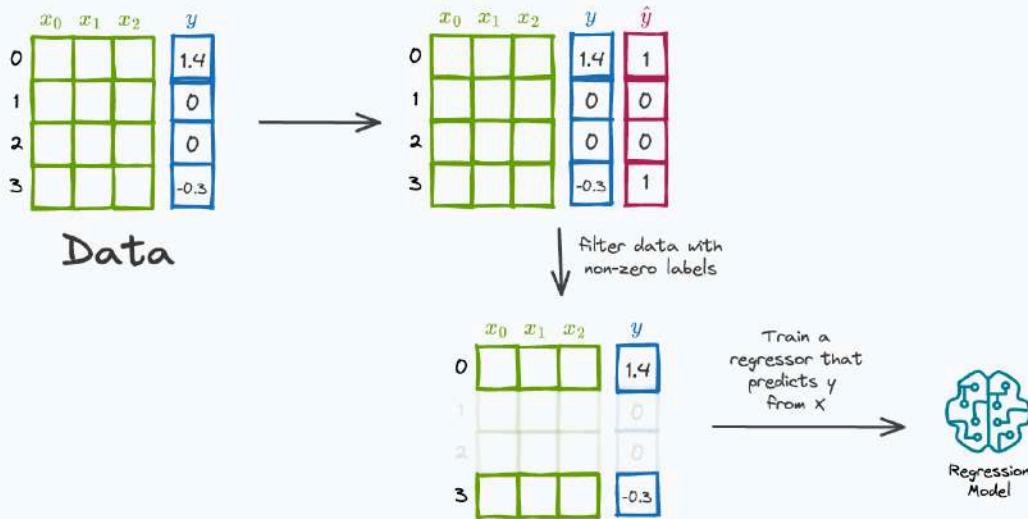
One great way to solve this is by training a combination of a classification and a regression model.

This goes as follows:

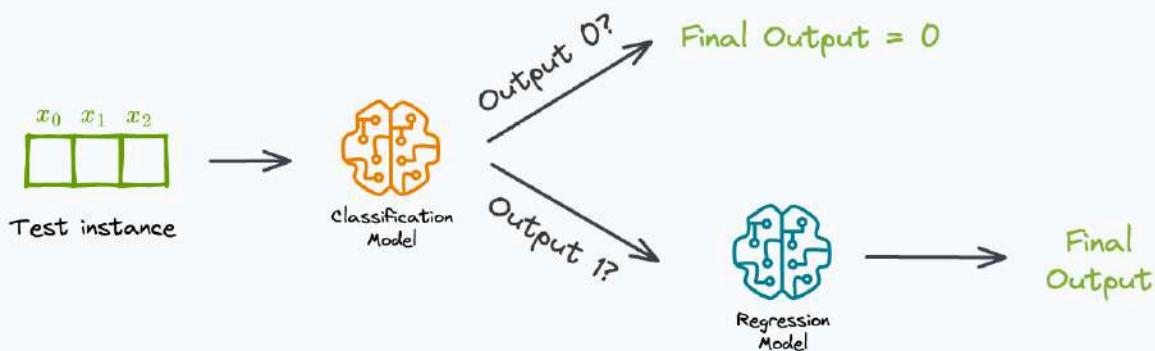
- Mark all non-zero targets as “1” and the rest as “0”.
- Train a binary classifier on this dataset.



- Next, train a regression model only on those data points with a non-zero true target.



During prediction:



- If the classifier's output is “0”, the final output is also zero.
- If the classifier's output is “1”, use the regression model to predict the final output.

Its effectiveness over the regular regression model is evident from the image below:

Always Plot the Response Distribution Before Training a Model

 blog.DailyDoseofDS.com



Regression vs. Regression + Classification results

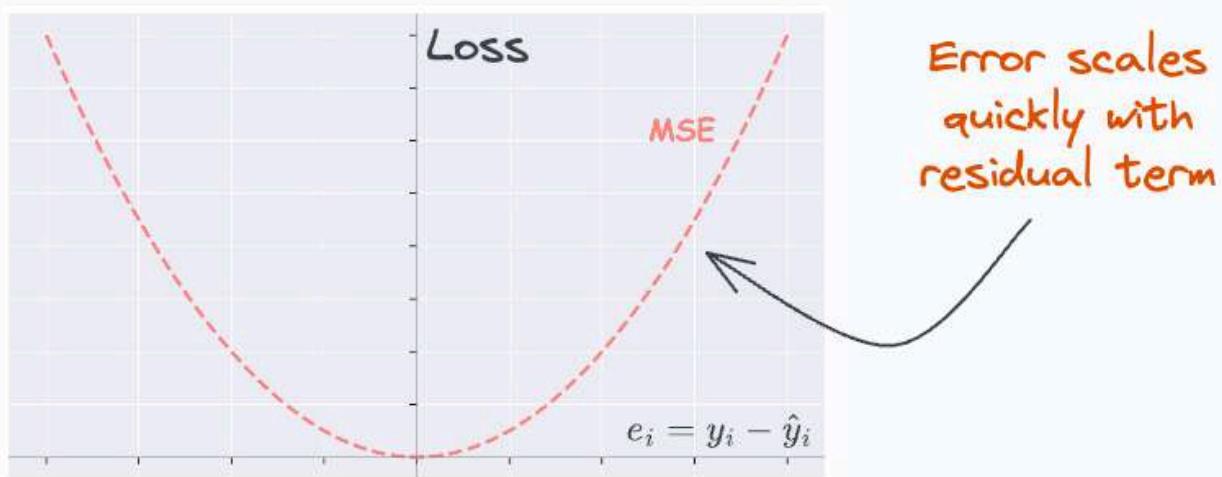
- Linear regression alone underfits the data.
- Linear regression with a classifier performs as expected.

Huber Regression

One big problem with regression models is that they are sensitive to outliers. Consider linear regression. Even a few outliers can significantly impact Linear Regression performance, as shown below:



And it isn't hard to identify the cause of this problem. Essentially, the loss function (MSE) scales quickly with the residual term (true-predicted).



Thus, even a few data points with a large residual can impact parameter estimation.

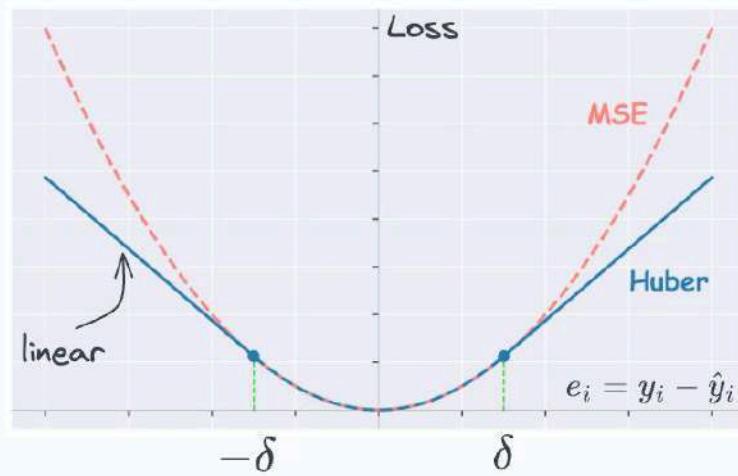
Huber loss (or Huber Regression) precisely addresses this problem. In a gist, it attempts to reduce the error contribution of data points with large residuals.

One simple, intuitive, and obvious way to do this is by applying a threshold (δ) on the residual term:

- If the residual is smaller than the threshold, use MSE (no change here).
- Otherwise, use a loss function which has a smaller output than MSE — linear, for instance.

This is depicted below:

- For residuals smaller than the threshold (δ) → we use MSE.
- Otherwise, we use a linear loss function which has a smaller output than MSE.



Mathematically, Huber loss is defined as follows:

$$e_i = y_i - \hat{y}_i$$

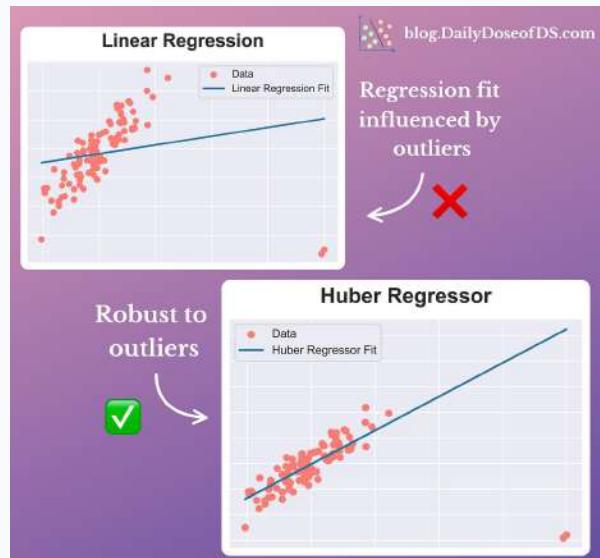
$$L_{\text{Huber}}(e_i) = \begin{cases} \frac{1}{2}e_i^2 & : e_i \leq \delta \\ \delta \cdot (|e_i| - \frac{1}{2}\delta) & : \text{otherwise} \end{cases}$$

MSE

Linear

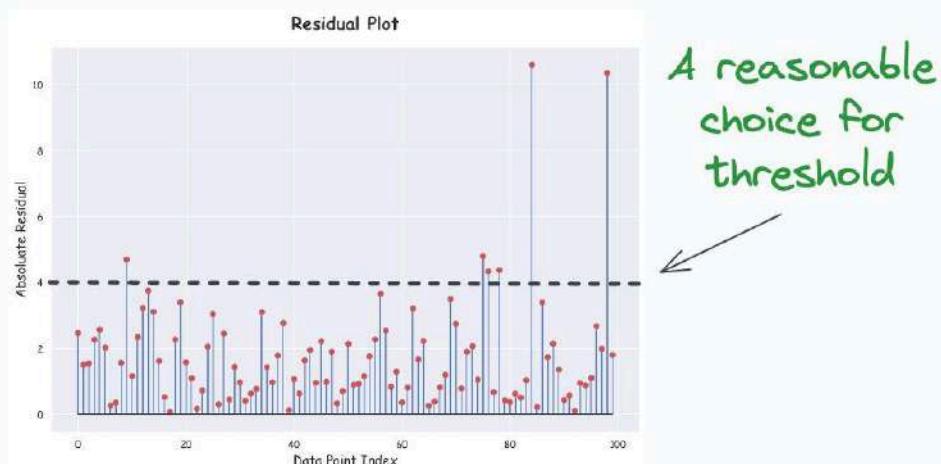
Its effectiveness is evident from the following image:

- Linear Regression is affected by outliers
- Huber Regression is more robust.



How do we determine the threshold (δ)?

While trial and error is one way, I often like to create a residual plot. This is depicted below: The below plot is generally called a lollipop plot because of its appearance.



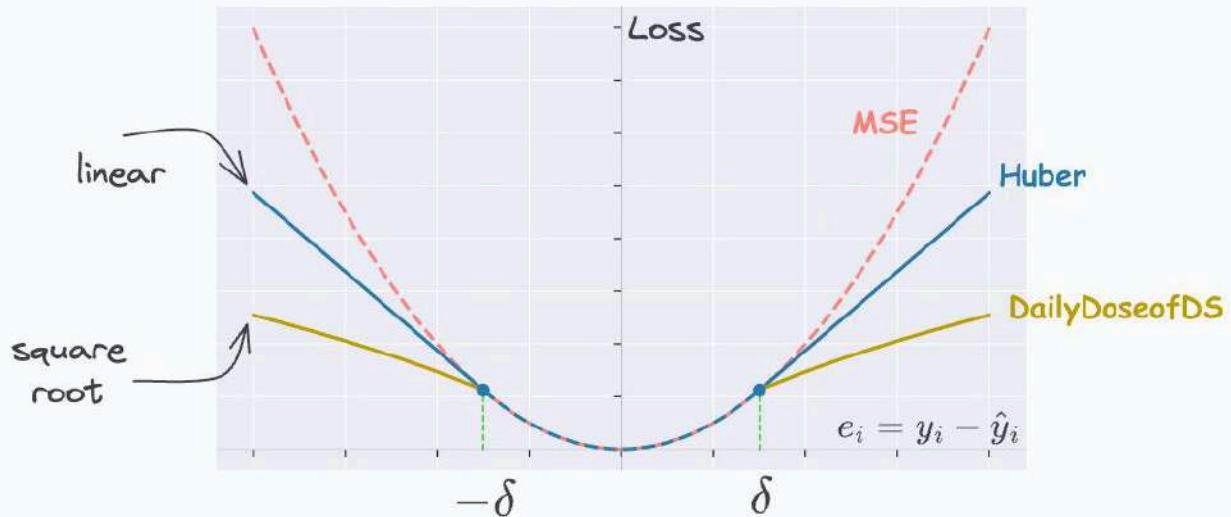
- Train a linear regression model as you usually would.
- Compute the residuals (=true-predicted) on the training data.
- Plot the absolute residuals for every data point.

One good thing is that we can create this plot for any dimensional dataset. The objective is just to plot (true-predicted) values, which will always be 1D.

Next, you can subjectively decide a reasonable threshold value δ .

Here's another interesting idea.

By using a linear loss function in Huber regressor, we intended to reduce the large error contributions that would have happened otherwise by using MSE. Thus, we can further reduce that error contribution by using, say, a square root loss function, as shown below:



I am unsure if this has been proposed before, so I decided to call it the DailyDoseofDataScience Regressor 😊.

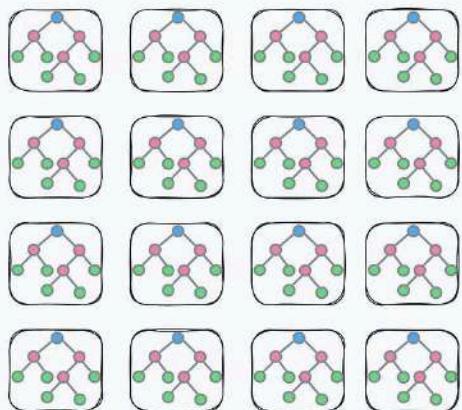
It is clear that the error contribution of the square root loss function is the lowest for all residuals above the threshold δ .

Decision Trees and Ensemble Methods

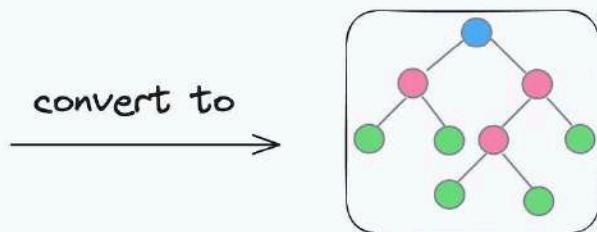
Condense Random Forest into a Decision Tree

There's an interesting technique, using which, we can condense an entire random forest model into a single decision tree.

Random Forest



Decision Tree



The benefits?

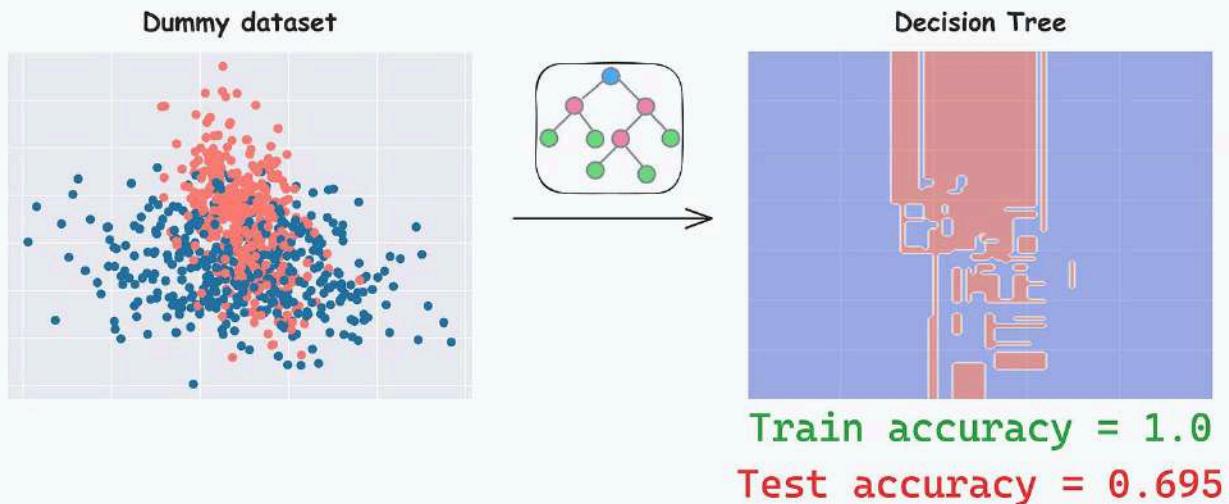
This technique can:

- Decrease the prediction run-time.
- Improve interpretability.
- Reduce the memory footprint.
- Simplify the model.
- Preserve the generalization power of the random forest model.

Let's understand in this chapter.

Technique Walkthrough

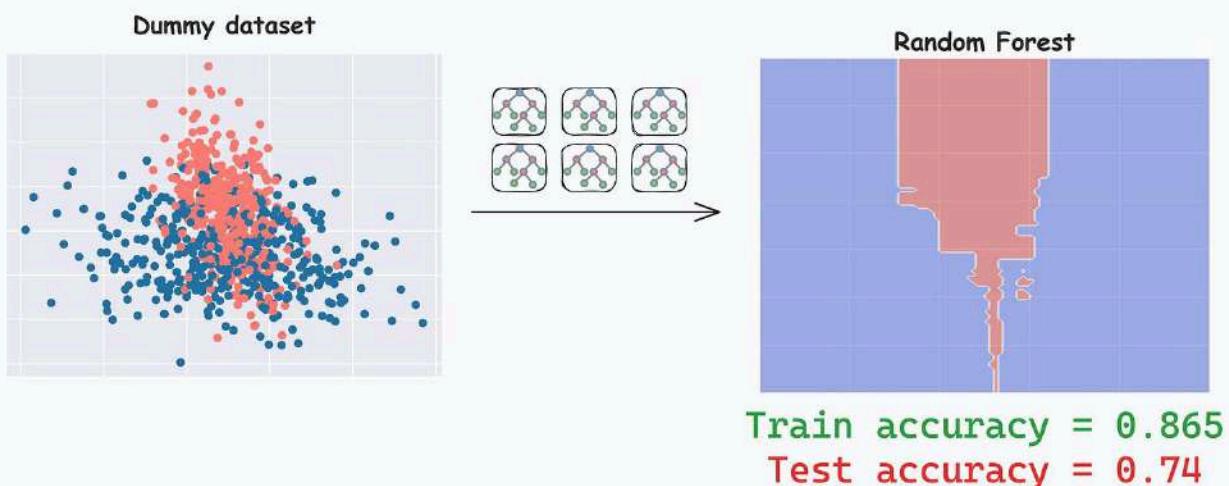
Let's fit a decision tree model on the following dummy dataset. It produces a decision region plot shown on the right.



It's clear that there is high overfitting.

In fact, we must note that, by default, a decision tree can always 100% overfit any dataset (we will use this information shortly). This is because it is always allowed to grow until all samples have been classified correctly.

This overfitting problem is resolved by a random forest model, as depicted below:

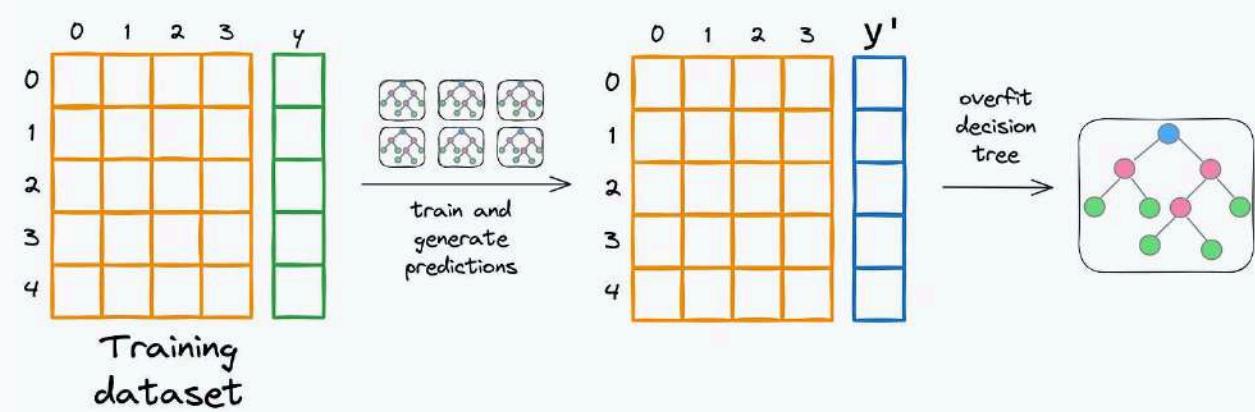


This time, the decision region plot suggests that we don't have a complex decision boundary. The test accuracy has also improved (69.5% to 74%).

Now, here's an interesting thing we can do.

We know that the random forest model has learned some rules that generalize on unseen data.

So, how about we train a decision tree on the predictions generated by the random forest model on the training set?



More specifically, given a dataset (X, y) :

- Train a random forest model. This will learn some rules from the training set which are expected to generalize on unseen data (due to Bagging).
- Generate predictions on X , which produces the output y' . These predictions will capture the essence of the rules learned by the random forest model.
- Finally, train a decision tree model on (X, y') . Here, we want to intentionally overfit this mapping as this mapping from (X) to (y') is a proxy for the rules learned by the random forest model.

This idea is implemented below:



```

from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier

X, y = ... # classification dataset

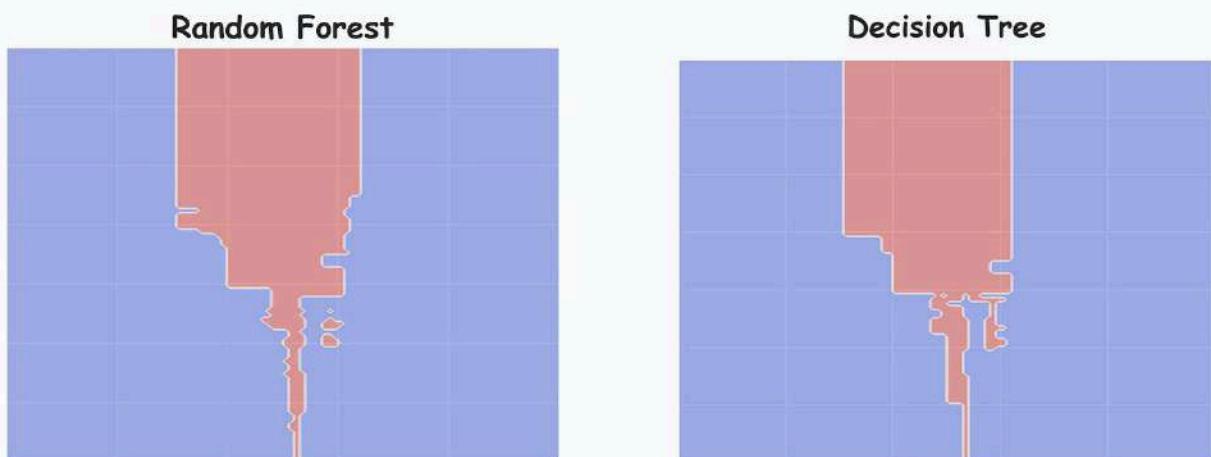
# train random forest model
rfmodel = RandomForestClassifier(...).fit(X, y)

# generate predictions on training set
y1 = rfmodel.predict(X)

# train decision tree model (overfit it intentionally)
dtree = DecisionTreeClassifier().fit(X, y1) ←
    
```

Train on output of Random Forest model

The decision region plot we get with the new decision tree is pretty similar to what we saw with the random forest earlier:



Measuring the test accuracy of the decision tree and random forest model, we notice them to be similar too:

Similar test accuracy

In fact, this approach also significantly reduces the run-time, as depicted below:

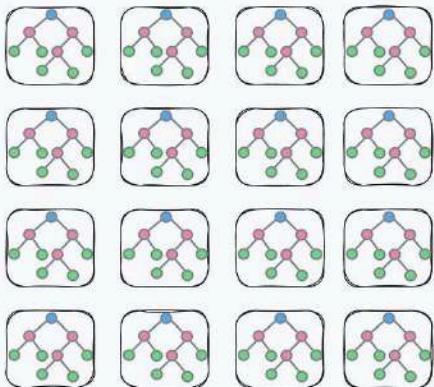
Decision tree is over 150x faster

Isn't that cool?

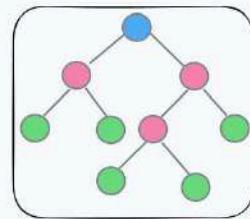
Another rationale for considering doing this is that it adds interpretability.

This is because if we have 100 trees in a random forest, there's no way we can interpret them.

However, if we have condensed it to a decision tree, now we can inspect it.



Low interpretability



High interpretability

A departing note

I devised this very recently. I also tested this approach on a couple of datasets, and they produced promising results.

But it won't be fair to make any conclusions based on just two instances.

While the idea makes intuitive sense, I understand there could be some potential flaws that are not evident right now.

So, I'm not saying that you should adopt this technique right away.

Instead, it is advised to test this approach on your random forest use cases. Considering reverting back to me with what you discovered.

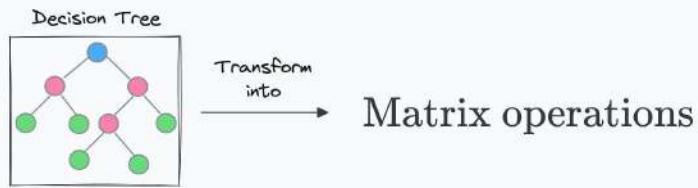
The code for this chapter is available here: <https://bit.ly/3XSPejD>.

In this next chapter, let's understand a technique to transform a decision tree into matrix operations which can run on GPUs.

Transform Decision Tree into Matrix Operations

Inference using a decision tree is an iterative process. We traverse a decision tree by evaluating the condition at a specific node in a layer until we reach a leaf node.

In this chapter, let's learn a superb technique that to represent inferences from a decision tree in the form of matrix operations.



As a result:

1. It makes inference much faster as matrix operations can be radically parallelized.
2. These operations can be loaded on a GPU for even faster inference, making them more deployment-friendly.

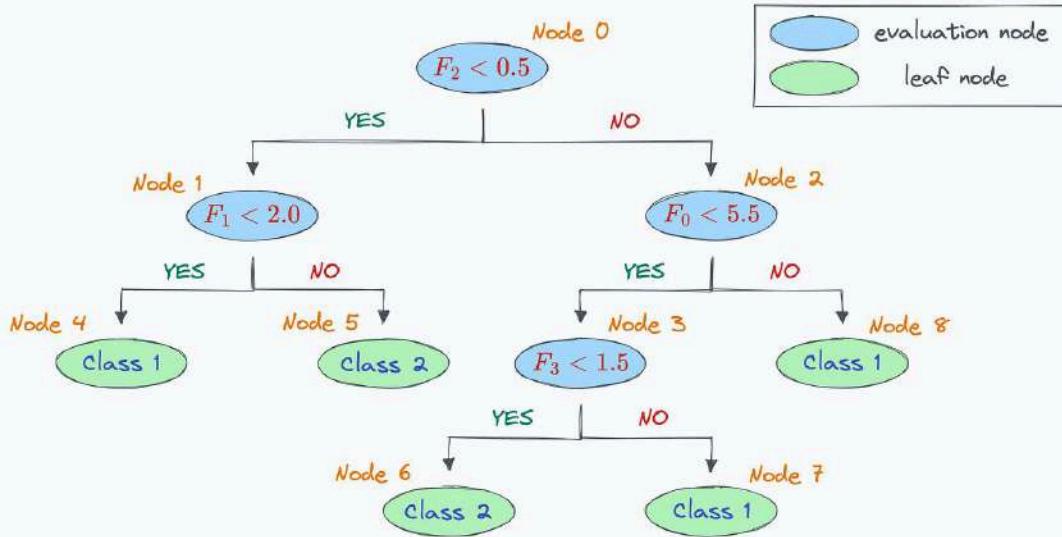
Setup

Consider a binary classification dataset with 5 features.

Binary
classification
dataset

F_1	F_2	F_3	F_4	F_5	y
					1
					0
					0
					1

Let's say we get the following tree structure after fitting a decision tree on the above dataset:



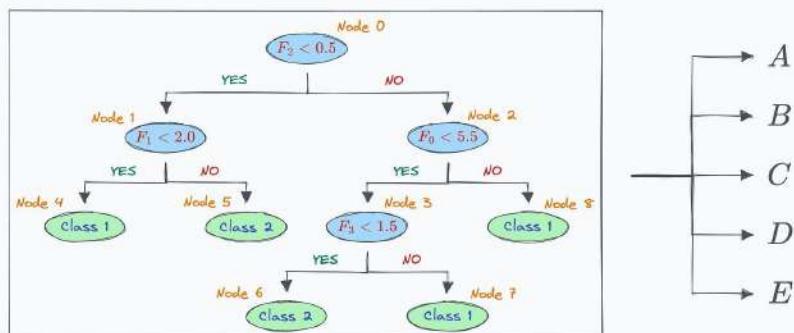
Notations

Before proceeding ahead, let's assume that:

- $m \rightarrow$ Total features in the dataset (5 in the dataset above).
- $e \rightarrow$ Total evaluation nodes in the tree (4 blue nodes in the tree above).
- $l \rightarrow$ Total leaf nodes in the tree (5 green nodes in the tree).
- $c \rightarrow$ Total classes in the dataset (2 in the dataset above).

Tree to Matrices

The core idea in this conversion is to derive five matrices (A, B, C, D, E) that capture the structure of the decision tree.



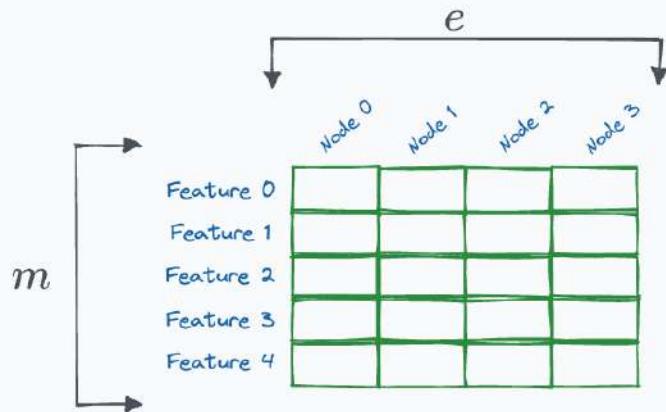
Let's understand them one by one!

Note: The steps we shall be performing below might not make immediate sense so please keep reading. You'll understand everything once we have derived all 5 matrices.

#1) Matrix A

This matrix captures the relationship between input features and evaluation nodes (blue nodes above).

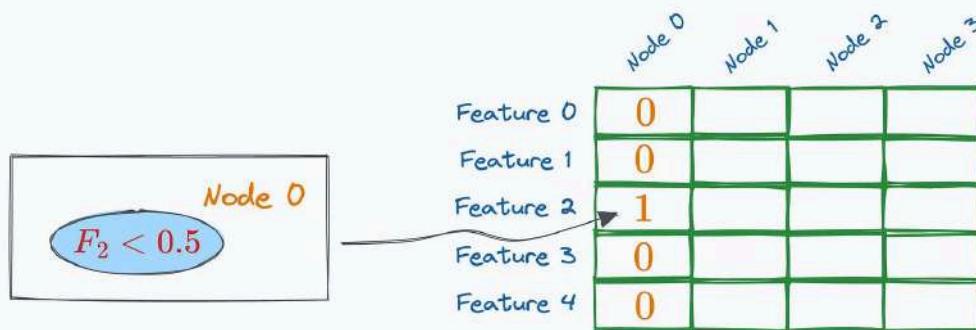
So it's an $(m \times e)$ shaped matrix.



A specific entry is set to “1” if the corresponding node in the column evaluates the corresponding feature in the row. For instance, in our decision tree, “Node 0” evaluates “Feature 2”.

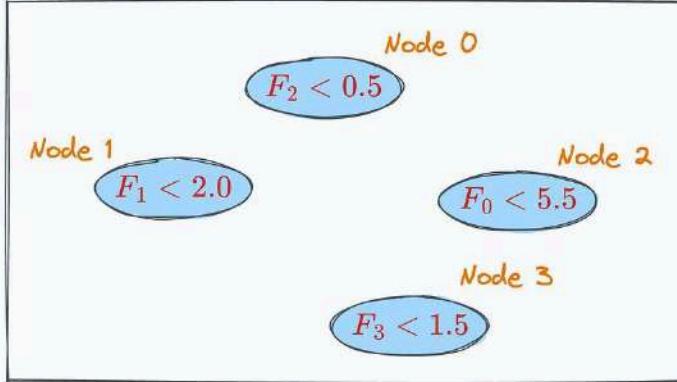


Thus, the corresponding entry will be “1” and all other entries will be “0.”



Filling out the entire matrix this way, we get:

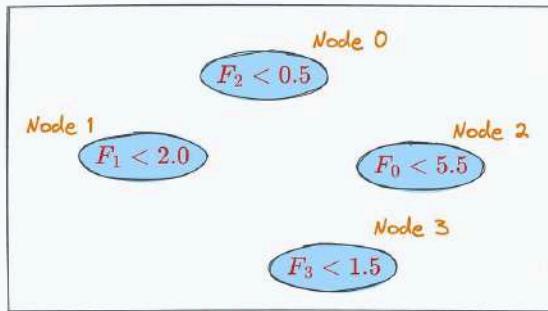
Matrix A



	Node 0	Node 1	Node 2	Node 3
Feature 0	0	0	1	0
Feature 1	0	1	0	0
Feature 2	1	0	0	0
Feature 3	0	0	0	1
Feature 4	0	0	0	0

#2) Matrix B

The entries of matrix B are the threshold value at each node. Thus, its shape is $1 \times e$.



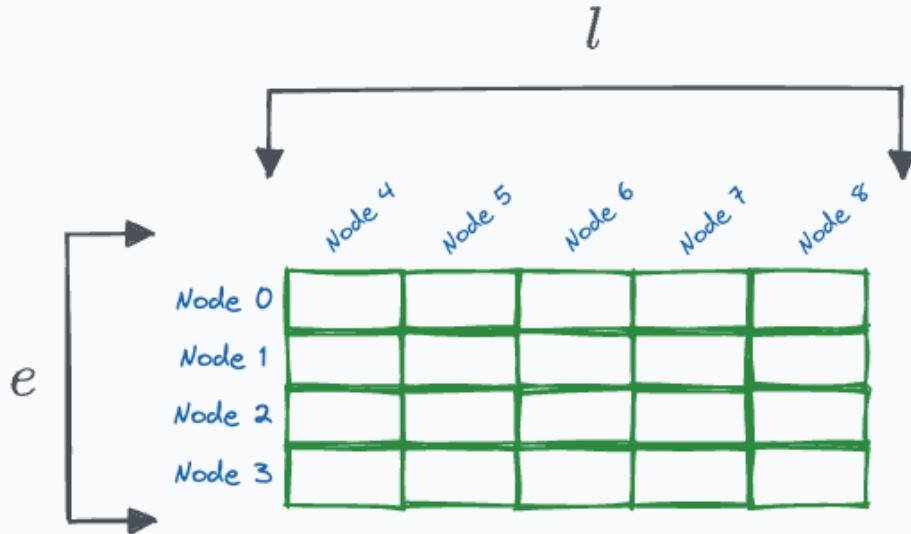
Matrix B

	Node 0	Node 1	Node 2	Node 3
	0.5	2.0	5.5	1.5

This is vector though, but the terminology is not important here.

#3) Matrix C

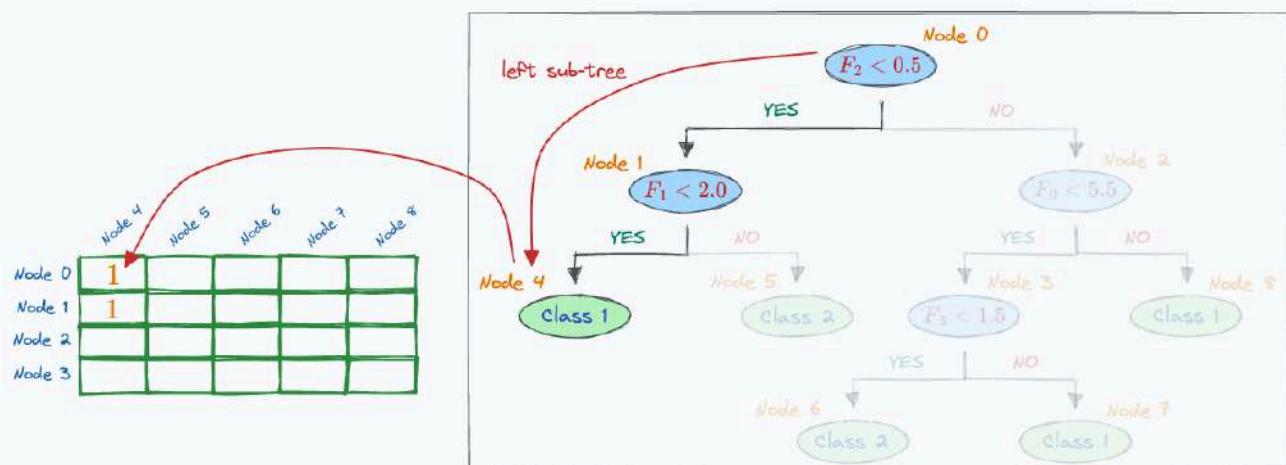
This is a matrix between every pair of leaf nodes and evaluation nodes. Thus, its dimensions are $e \times l$.



A specific entry is set to:

- “1” if the corresponding leaf node in the column lies in the left sub-tree of the corresponding evaluation node in the row.
- “-1” if the corresponding leaf node in the column lies in the right sub-tree of the corresponding evaluation node in the row.
- “0” if the corresponding leaf node and evaluation node have no link.

For instance, in our decision tree, the “leaf node 4” lies on the left sub-tree of both “evaluation node 0” and “evaluation node 1”. Thus, the corresponding values will be 1.



Filling out the entire matrix this way, we get:

Matrix C

	Node 4	Node 5	Node 6	Node 7	Node 8
Node 0	1	1	-1	-1	-1
Node 1	1	-1	0	0	0
Node 2	0	0	1	1	-1
Node 3	0	0	1	-1	0

#4) Matrix (or Vector) D

The entries of vector D are the sum of non-negative entries in every column of Matrix C:

Matrix C

	Node 4	Node 5	Node 6	Node 7	Node 8
Node 0	1	1	-1	-1	-1
Node 1	1	-1	0	0	0
Node 2	0	0	1	1	-1
Node 3	0	0	1	-1	0



Matrix D

	Node 4	Node 5	Node 6	Node 7	Node 8
	2	1	2	1	0

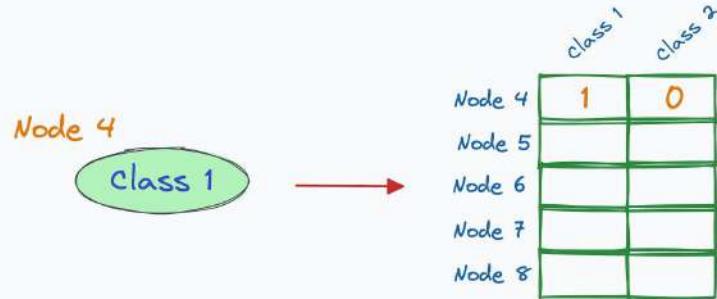
#5) Matrix E

Finally, this matrix holds the mapping between leaf nodes and their corresponding output labels. Thus, its dimensions are $l \times c$.

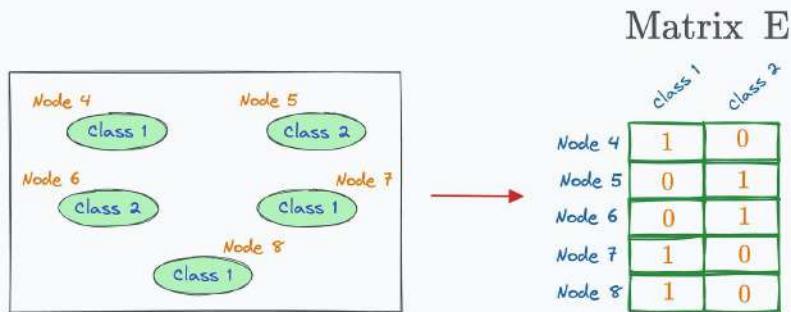
	class 1	class 2
Node 4		
Node 5		
Node 6		
Node 7		
Node 8		

If a leaf node classifies a sample to “Class 1”, the corresponding entry will be 1, and the other cell entry will be 0.

For instance, “lead node 4” outputs “Class 1”, thus the corresponding entries for the first row will be (1,0):



We repeat this for all other leaf nodes to get the following matrix as Matrix E:



With this, we have compiled our decision tree into matrices. To recall, these are the five matrices we have created so far:

	Node 0	Node 1	Node 2	Node 3
Feature 0	0	0	1	0
Feature 1	0	1	0	0
Feature 2	1	0	0	0
Feature 3	0	0	0	1
Feature 4	0	0	0	0

	Node 0	Node 1	Node 2	Node 3
	0.5	2.0	5.5	1.5

	Node 4	Node 5	Node 6	Node 7	Node 8
	class 1	class 2			
Node 4	1	0			
Node 5	0	1			
Node 6	0	1			
Node 7	1	0			
Node 8	1	0			

	Node 4	Node 5	Node 6	Node 7	Node 8
	1	1	-1	-1	-1
Node 0	1	-1	0	0	0
Node 1	0	0	1	1	-1
Node 2	0	0	1	-1	0
Node 3	0	0	1	-1	0

	Node 4	Node 5	Node 6	Node 7	Node 8
	2	1	2	1	0

Matrix D

- Matrix A captures which input feature was used at each evaluation node.
- Matrix B stores the threshold of each evaluation node.
- Matrix C captures whether a leaf node lies on the left or right sub-tree of a specific evaluation node or has no relation to it.
- Matrix D stores the sum of non-negative entries in every column of Matrix C.
- Finally, Matrix E maps from leaf nodes to their class labels.

Inference using matrices

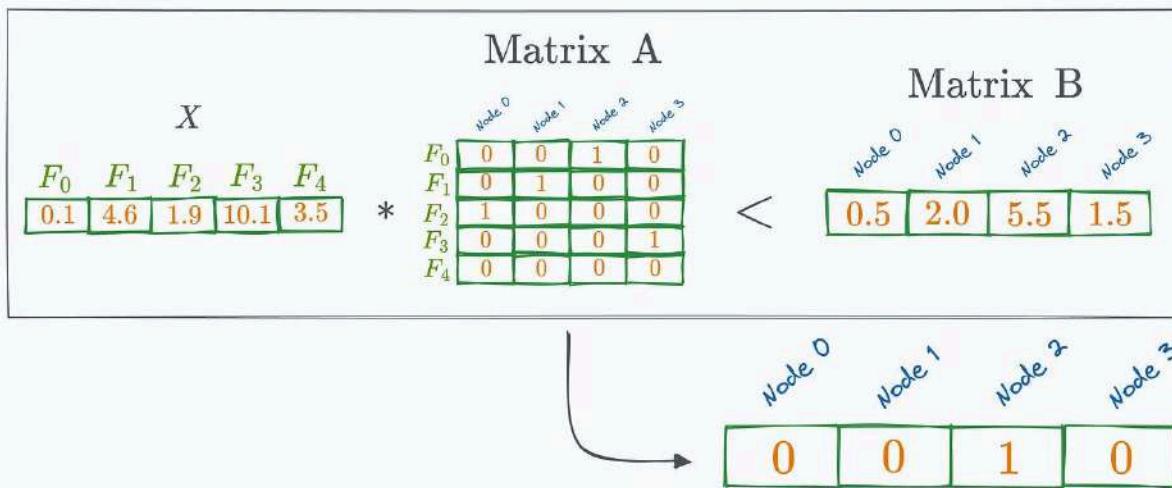
Say this is our input feature vector X (5 dimensions):

F_0	F_1	F_2	F_3	F_4
0.1	4.6	1.9	10.1	3.5

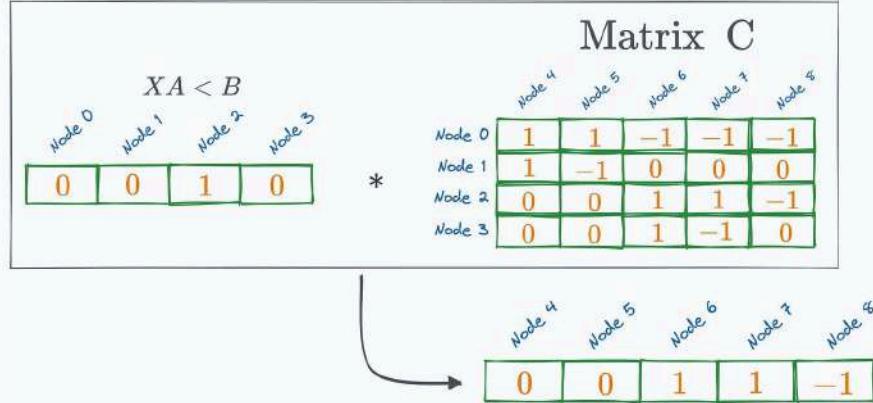
The whole inference can now be done using just these matrix operations:

$$(((X \cdot A) < B) * C) == D * E$$

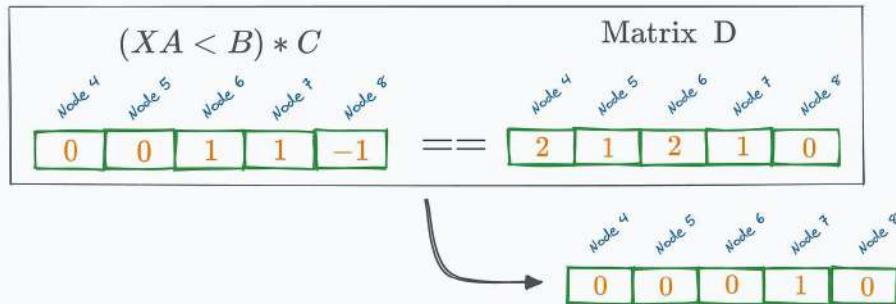
- $X \cdot A < B$ gives:



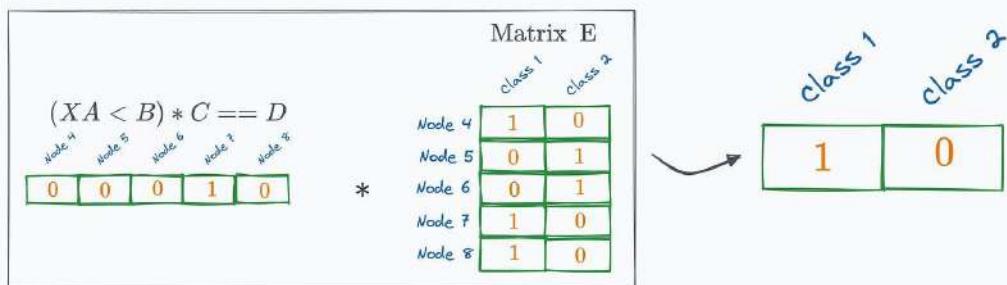
- The above result multiplied by C gives:



- The above result, when matched to D, gives:



- Finally, multiplying with E, we get:



The final prediction comes out to be “Class 1,” which is indeed correct! Notice that we carried out the entire inference process using only matrix operations:

$$(((XA < B) * C) == D) * E$$

As a result, the inference operation can largely benefit from parallelization and GPU capabilities.

The run-time efficacy of this technique is evident from the image below:

Sklearn RF Model	Tensor CPU Model	Tensor GPU Model
<code>RF_model.score(X, y)</code>	<code>RF_tensor_CPU.score(X, y)</code>	<code>RF_tensor_GPU.score(X, y)</code>
0.7262	0.7262	0.7262
<code>%%timeit</code>	<code>%%timeit</code>	<code>%%timeit</code>
<code>RF_model.predict(X)</code>	<code>RF_tensor_CPU.predict(X)</code>	<code>RF_tensor_GPU.predict(X)</code>
113 ms ± 14 ms per loop	51.4 ms ± 499 µs per loop	2.81 ms ± 18.5 µs per loop

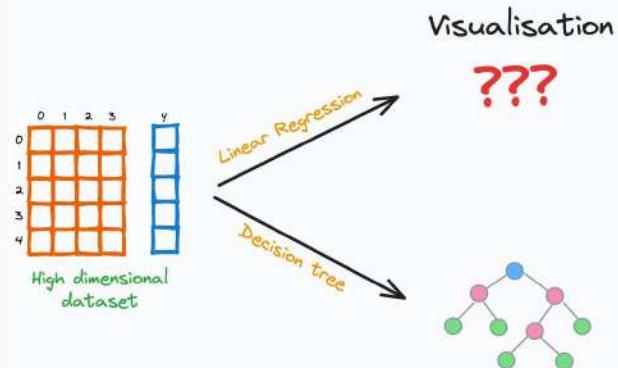
CPU Model is 2x Faster

GPU Model is 40x Faster

- Here, we have trained a random forest model.
- The compiled model runs:
 - Over twice as fast on a CPU (Tensor CPU Model).
 - ~40 times faster on a GPU, which is huge (Tensor GPU Model).
- All models have the same accuracy — indicating no loss of information during compilation.

Interactively Prune a Decision Tree

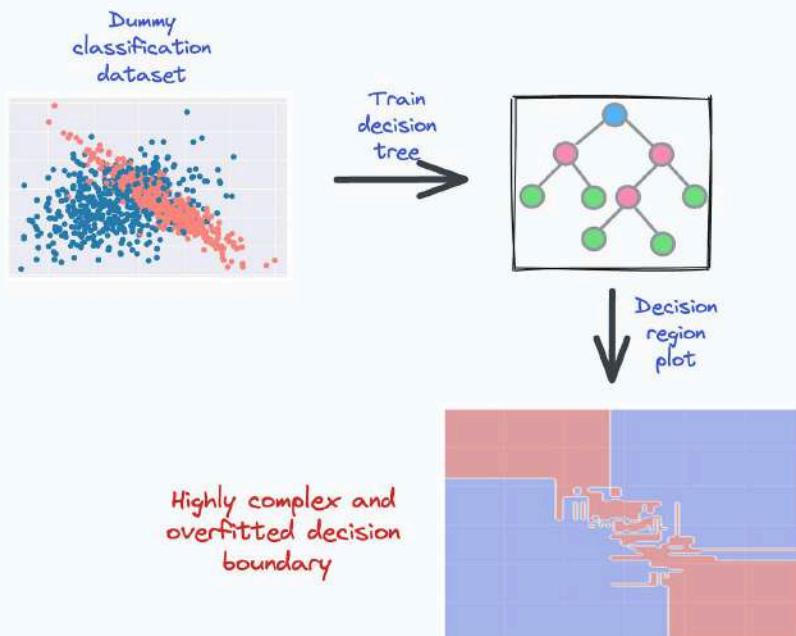
One thing I always appreciate about decision trees is their ease of visual interpretability. No matter how many features our dataset has, we can **ALWAYS** visualize and interpret a decision tree.



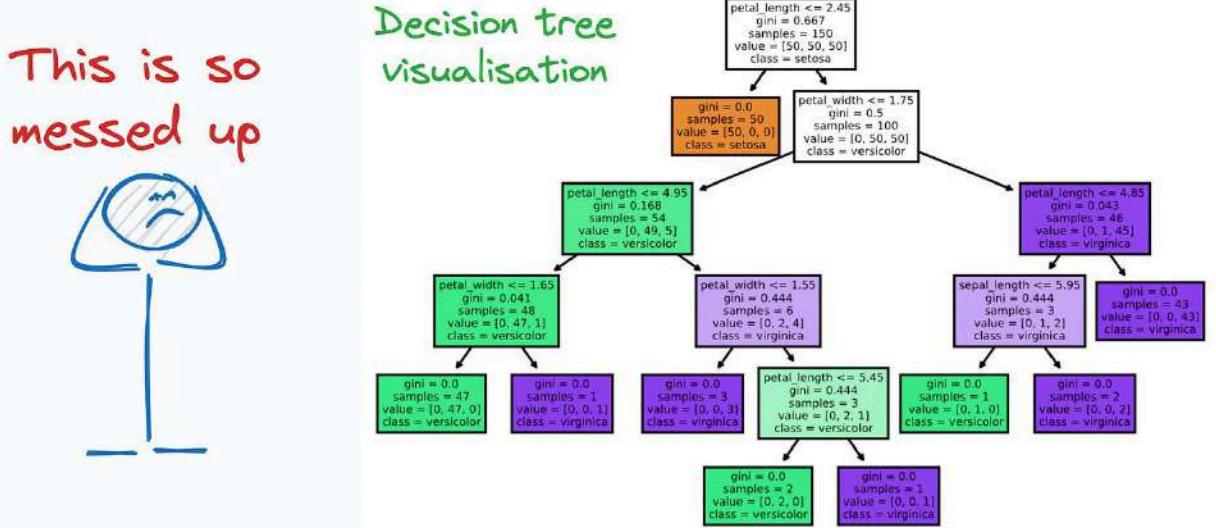
This is not always possible with other intuitive and simple models like linear regression. But decision trees stand out in this respect. Nonetheless, one thing I often find a bit time-consuming and somewhat hit-and-trial-driven is pruning a decision tree.

Why prune?

The problem is that under default conditions, decision trees **ALWAYS** 100% overfit the dataset, as depicted in this image:



Thus, pruning is ALWAYS necessary to reduce model variance. Scikit-learn already provides a method to visualize them as shown below:

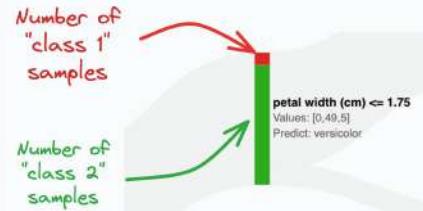


But the above visualisation is pretty non-elegant, tedious, messy, and static (or non-interactive). I recommend using an interactive Sankey diagram to prune decision trees. This is depicted below:

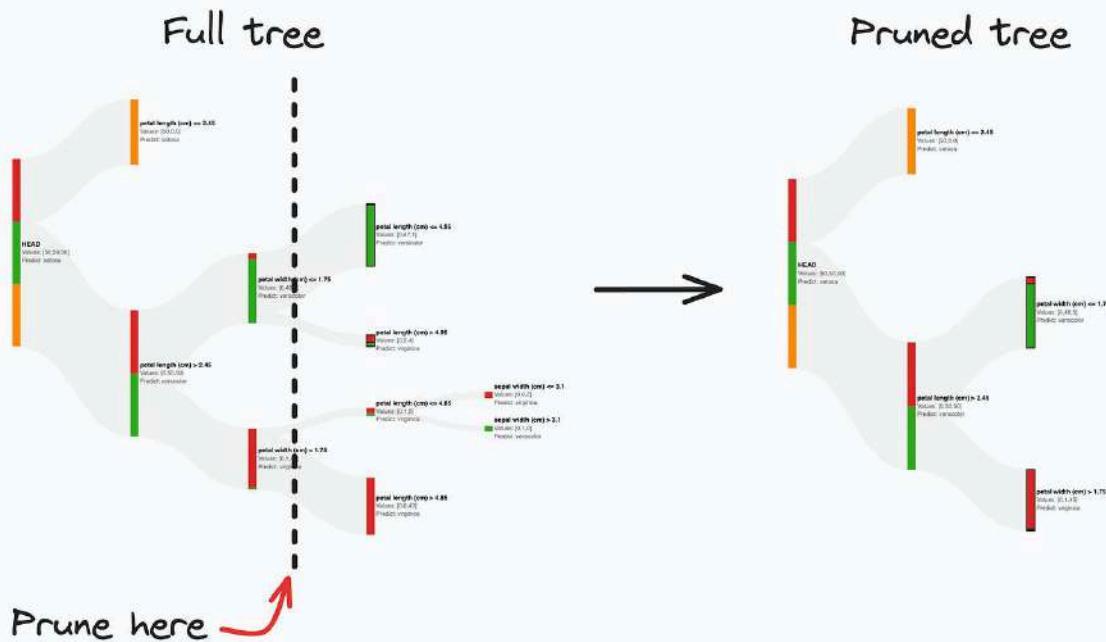


As shown above, the Sankey diagram allows us to interactively visualize and prune a decision tree by collapsing its nodes.

Also, the number of data points from each class is size and color-encoded in each node, as shown below.



This instantly gives an estimate of the node's impurity, based on which, we can visually and interactively prune the tree in seconds. For instance, in the full decision tree shown below, pruning the tree at a depth of two appears reasonable:



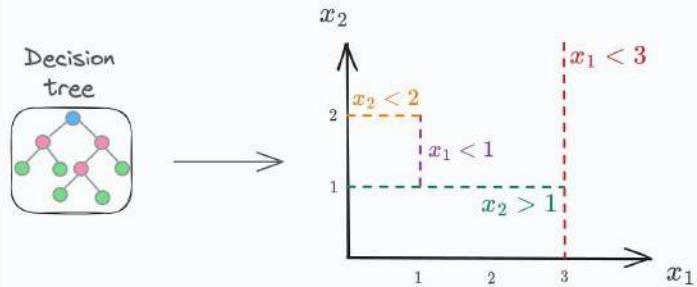
Next, we can train a new decision tree after obtaining an estimate for hyperparameter values. This will help us reduce the variance of the decision tree.

You can download the code notebook for the interactive decision tree here: <https://bit.ly/4bBwY1p>. Instructions are available in the notebook.

Next, let's understand a point of caution when using decision trees.

Why Decision Trees Must Be Thoroughly Inspected After Training

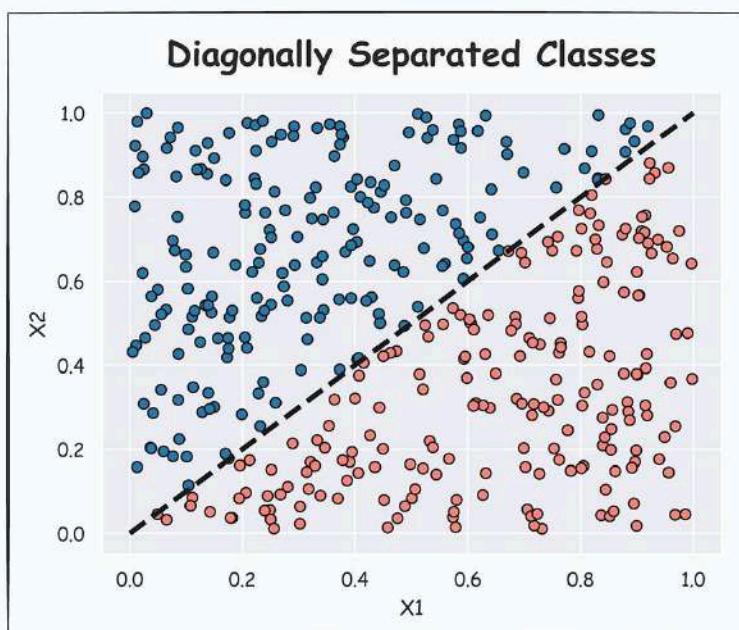
If we were to visualize the decision rules (the conditions evaluated at every node) of ANY decision tree, we would ALWAYS find them to be perpendicular to the feature axes, as depicted in the image.



In other words, every decision tree progressively segregates feature space based on such perpendicular boundaries to split the data.

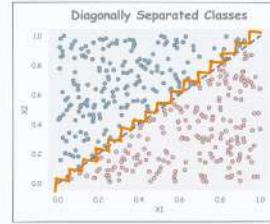
Of course, this is not a “problem” per se.

In fact, this perpendicular splitting is what makes it so powerful to perfectly overfit any dataset. However, this also brings up a pretty interesting point that is often overlooked when fitting decision trees. More specifically, what would happen if our dataset had a diagonal decision boundary, as depicted below:

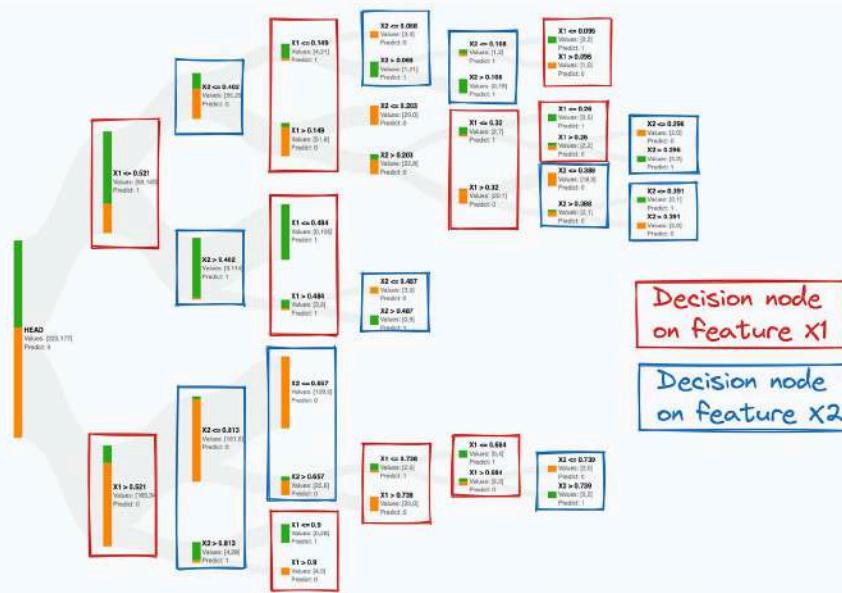


Dummy dataset
with diagonal
decision boundary

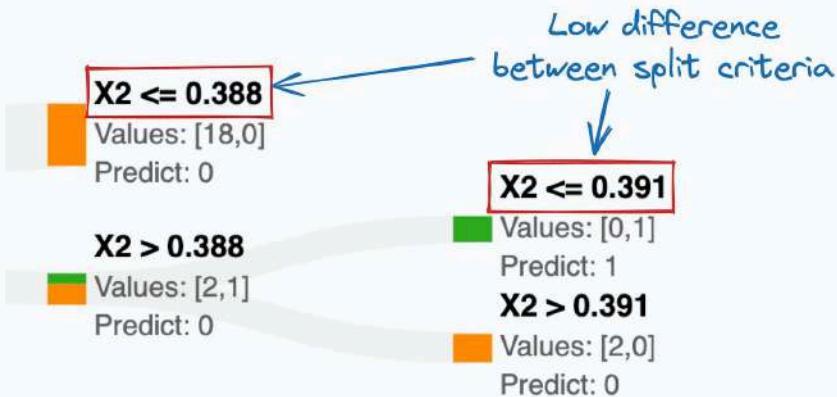
It is easy to guess that in such a case, the decision boundary learned by a decision tree is expected to appear as follows:



In fact, if we plot this decision tree, we notice that it creates so many splits just to fit this easily separable dataset, which a model like logistic regression, support vector machine (SVM), or even a small neural network can easily handle:



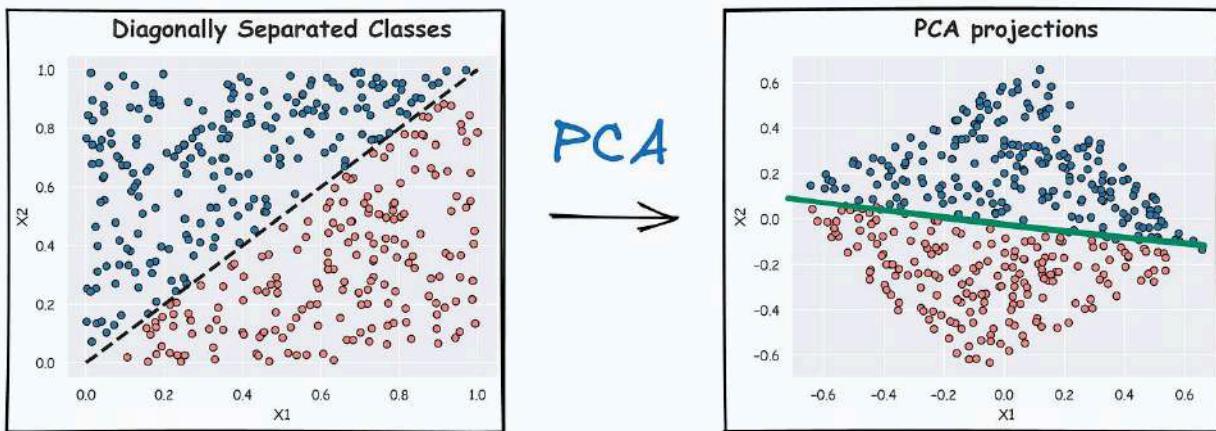
It becomes more evident if we zoom into this decision tree and notice how close the thresholds of its split conditions are:



This is a bit concerning because it clearly shows that the decision tree is meticulously trying to mimic a diagonal decision boundary, which hints that it might not be the best model to proceed with. To double-check this, I often do the following:

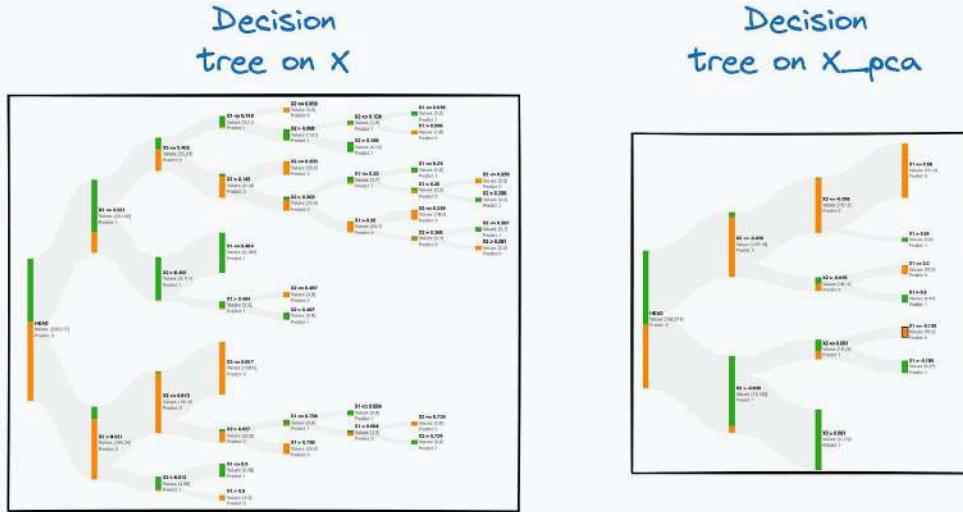
- Take the training data (X, y);
 - Shape of X : (n, m) .
 - Shape of y : $(n, 1)$.
- Run PCA on X to project data into an orthogonal space of m dimensions. This will give X_{pca} , whose shape will also be (n, m) .
- Fit a decision tree on X_{pca} and visualize it (thankfully, decision trees are always visualizable).
- If the decision tree depth is significantly smaller in this case, it validates that there is a diagonal separation.

For instance, the PCA projections on the above dataset are shown below:



It is clear that the decision boundary on PCA projections is almost perpendicular to the X_2 feature (the 2nd principal component).

Fitting a decision tree on this X_{pca} drastically reduces its depth, as depicted below:



This lets us determine that we might be better off using some other algorithm instead.

Or, we can spend some time engineering better features that the decision tree model can easily work with using its perpendicular data splits.

At this point, if you are thinking, why can't we use the decision tree trained on X_{pca} ?

While nothing stops us from doing that, do note that PCA components are not interpretable, and maintaining feature interpretability can be important at times. Thus, whenever you train your next decision tree model, consider spending some time inspecting what it's doing.

Before ending this chapter

I don't intend to discourage the use of decision trees. They are the building blocks of some of the most powerful ensemble models we use today.

My point is to bring forward the structural formulation of decision trees and why/when they might not be an ideal algorithm to work with.

Decision Trees **ALWAYS** Overfit!

In addition to the above inspection, there's one more thing you need to be careful of when using decision trees. This is about overfitting.

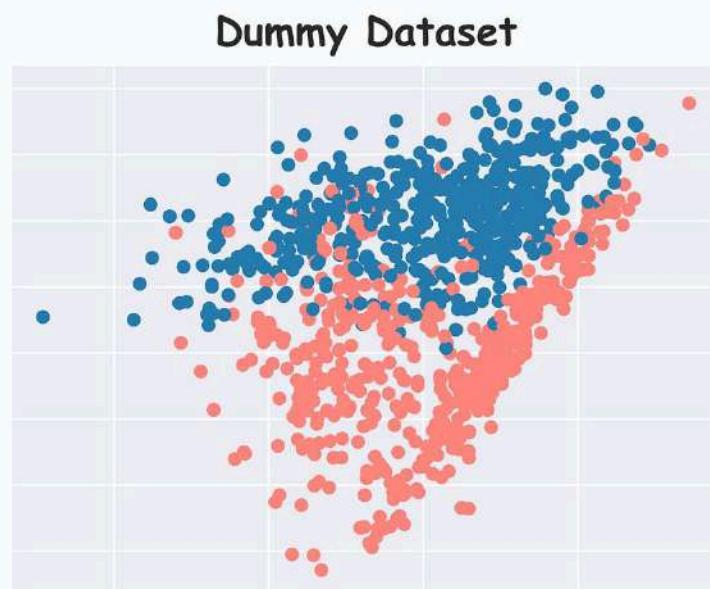
The thing is that, by default, a decision tree (in sklearn's implementation, for instance), is allowed to grow until all leaves are pure. This happens because a standard decision tree algorithm greedily selects the best split at each node.



This makes its nodes more and more pure as we traverse down the tree. As the model correctly classifies ALL training instances, it leads to 100% overfitting, and poor generalization.

For instance, consider this dummy dataset:

Dummy dataset
with intermixed
classes



Fitting a decision tree on this dataset gives us the following decision region plot:



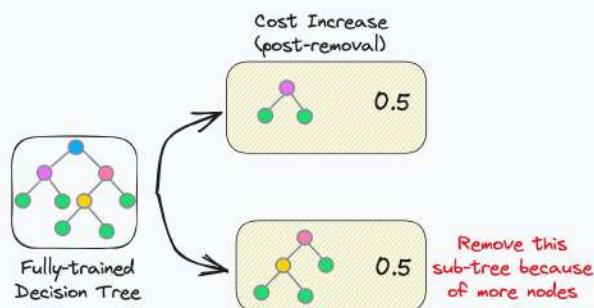
It is pretty evident from the decision region plot, the training and test accuracy that the model has entirely overfitted our dataset.

Cost-complexity-pruning (CCP) is an effective technique to prevent this.

CCP considers a combination of two factors for pruning a decision tree:

- Cost (C): Number of misclassifications
- Complexity (C): Number of nodes

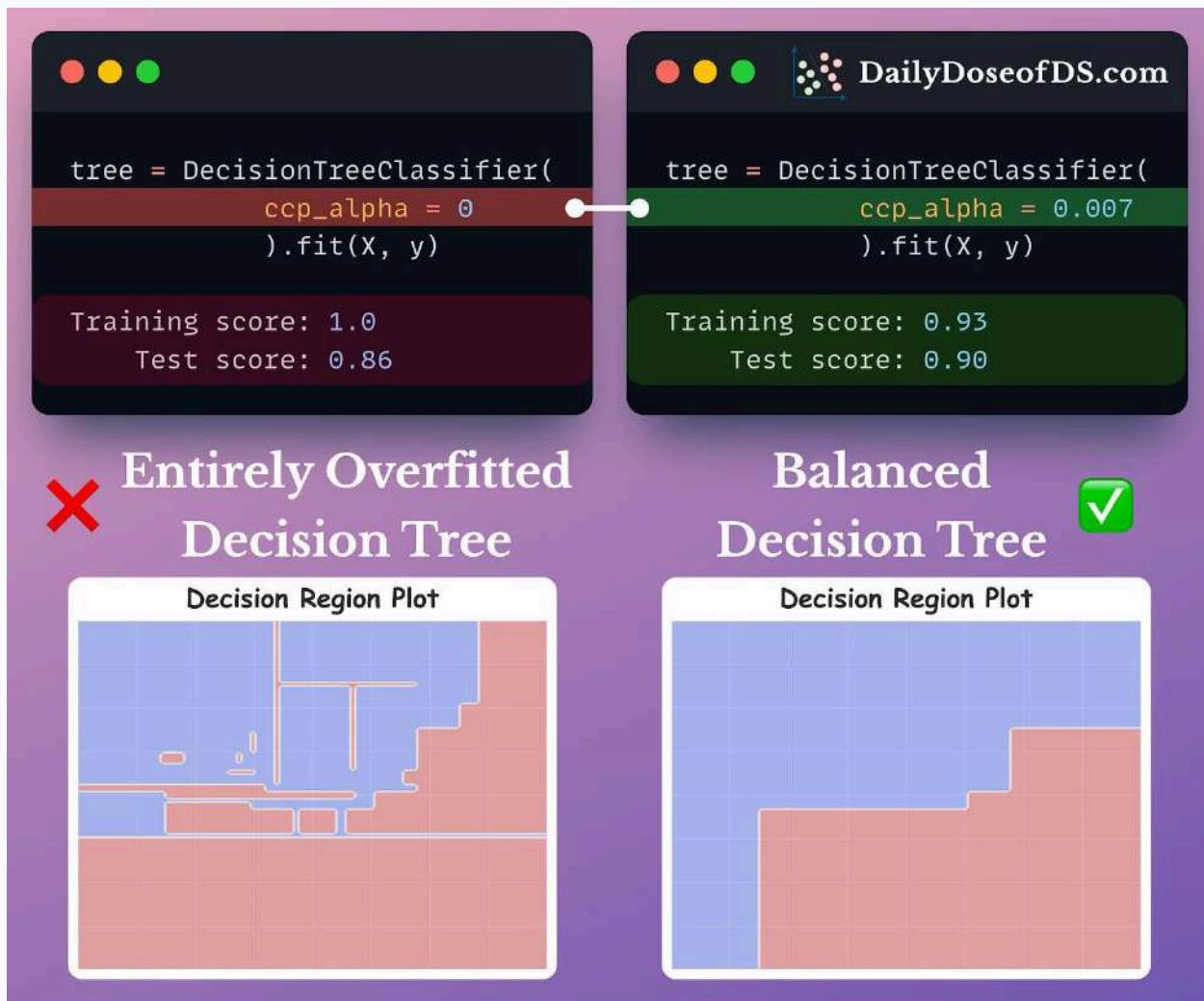
The core idea is to iteratively drop sub-trees, which, after removal, lead to a minimal increase in classification cost AND a maximum reduction of complexity (or nodes). In other words, if two sub-trees lead to a similar increase in classification cost, then it is wise to remove the sub-tree with more nodes.



In sklearn, you can control cost-complexity-pruning using the `ccp_alpha` parameter:

- large value of `ccp_alpha` → results in underfitting
- small value of `ccp_alpha` → results in overfitting

The objective is to determine the optimal value of `ccp_alpha`, which gives a better model. The effectiveness of cost-complexity-pruning is evident from the image below:



As depicted above, CCP results in a much simpler and acceptable decision region plot.

OOB Validation in Random Forest

After training an ML model on a training set, we always keep a held-out validation/test set for evaluation. I am sure you already know the purpose, so we won't discuss that.

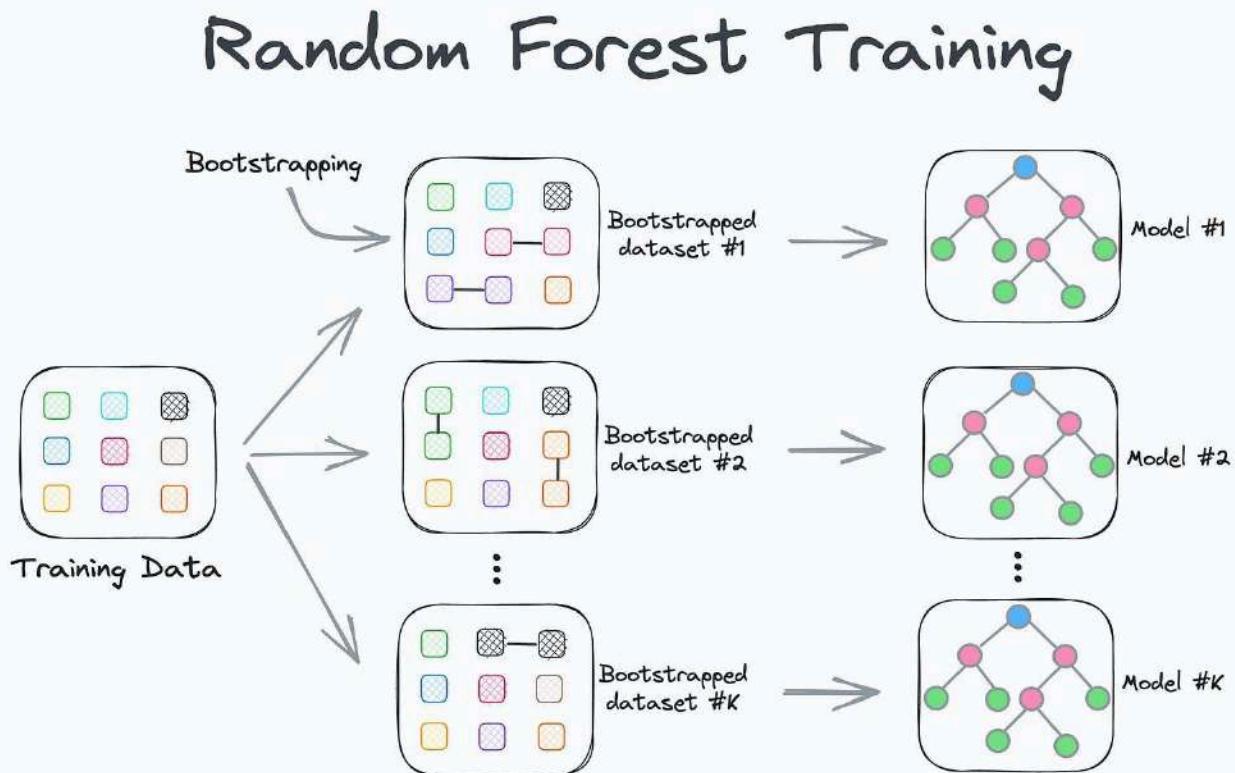
But do you know that random forests are an exception to that? In other words, one can somewhat "evaluate" a random forest using the training set itself.

Let's understand how.

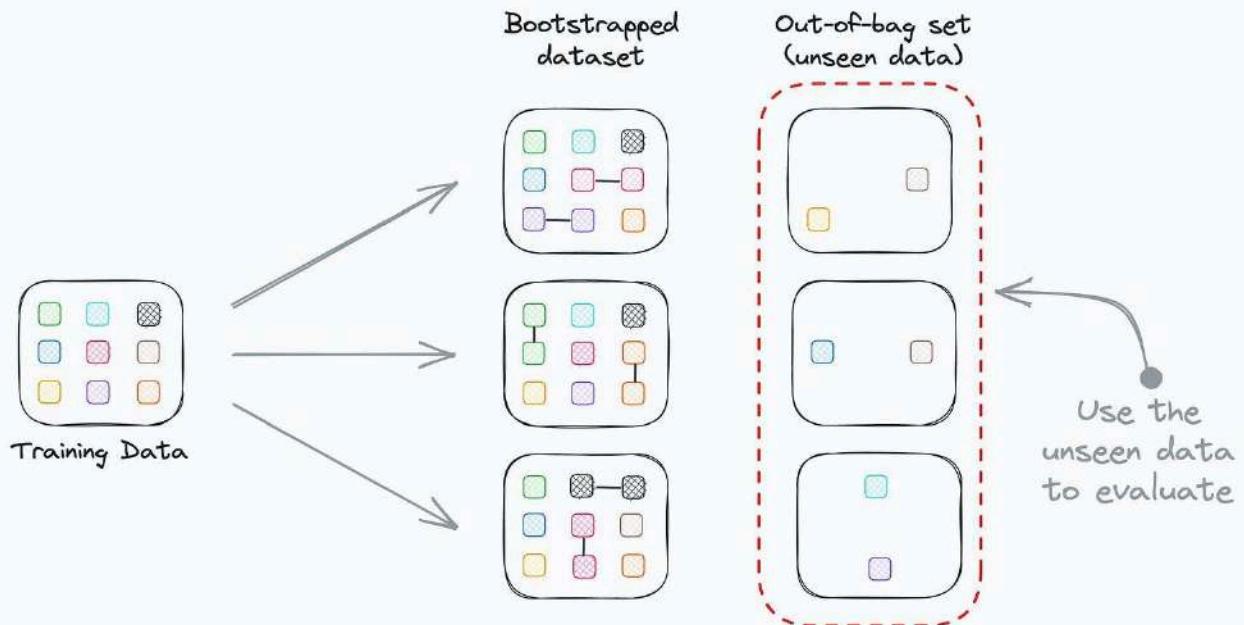
To recap, a random forest is trained as follows:

- First, we create different subsets of data with replacement (this process is called bootstrapping).
- Next, we train one decision tree per subset.
- Finally, we aggregate all predictions to get the final prediction.

This process is depicted below:

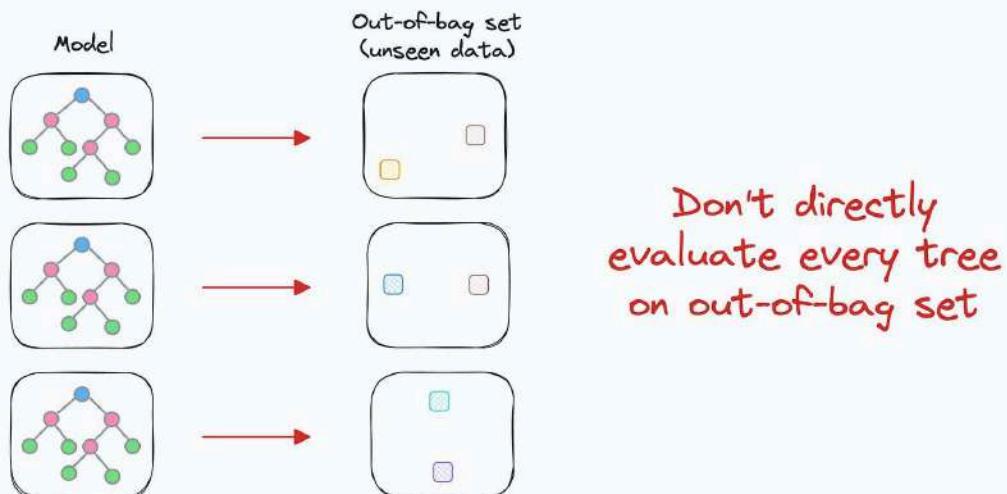


If we look closely above, every subset has some missing data points from the original training set.



We can use these observations to validate the model. This is also called out-of-bag validation. Calculating the out-of-bag score for the whole random forest is simple too.

But one thing to remember is that we CAN NOT evaluate individual decision trees on their specific out-of-bag sample and generate some sort of “aggregated score” for the entire random forest model.

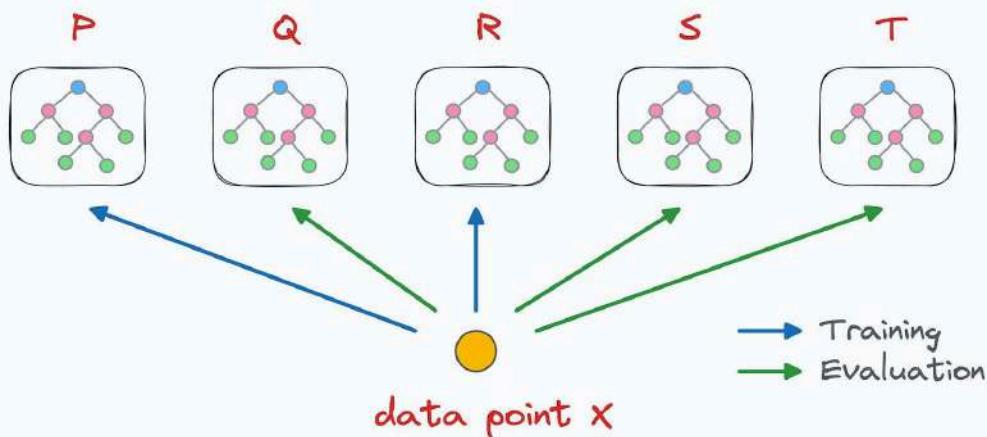


This is because a random forest is not about what a decision tree says individually. Instead, it's about what all decision trees say collectively.

So here's how we can generate the out-of-bag score for the random forest model. For every data point in the training set:

- Gather predictions from all decision trees that did not use it as a training data point.
- Aggregate predictions to get the final prediction.

For instance, consider a RF model with 5 decision trees $\rightarrow (P, Q, R, S, T)$. Say a specific data point X was used as a training data in decision trees P and R .



So we shall gather the out-of-bag prediction for data point X from decision trees Q, S and T.

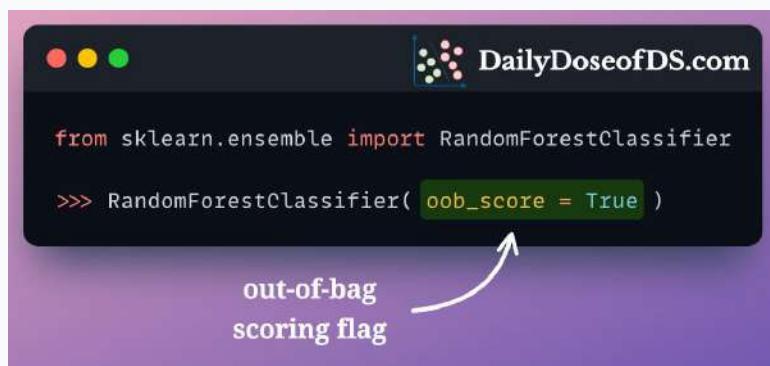
After obtaining out-of-bag predictions for all samples, we score them to get the out-of-bag score.

See...this technique allowed us to evaluate a random forest model on the training set. Of course, I don't want you to blindly adopt out-of-bag validation without understanding some of its advantages and considerations.

I have found out-of-bag validation to be particularly useful in the following situations:

- In low-data situations, out-of-bag validation prevents data splitting whilst obtaining a good proxy for model validation.
- In large-data situations, traditional cross-validation techniques are computationally expensive. Here, out-of-bag validation provides an efficient alternative. This is because, by its very nature, even cross-validation provides an out-of-fold metric. Out-of-bag validation is also based on a similar principle.

And, of course, an inherent advantage of out-of-bag validation is that it guarantees no data leakage. Luckily, out-of-bag validation is also neatly tied in sklearn's random forest implementation.



```

from sklearn.ensemble import RandomForestClassifier
>>> RandomForestClassifier( oob_score = True )

```

out-of-bag scoring flag

The most significant consideration about out-of-bag score is to use it with caution for model selection, model improvement, etc.

This is because if we do, we typically tend to overfit the out-of-bag score as the model is essentially being tuned to perform well on the data points that were left out during its training. And if we consistently improve the model based on the out-of-bag score, we obtain an overly optimistic evaluation of its generalization performance.

If I were to share just one lesson here based on my experience, it would be that if we don't have a true (and entirely different) held-out set for validation, we will overfit to some extent.

The decisions made may be too specific to the out-of-bag sample and may not generalize well to new data.

Train Random Forest on Large Datasets

Most classical ML algorithms cannot be trained with a batch implementation. This limits their usage to only small/intermediate datasets. For instance, this is the list of sklearn implementations that support a batch API:

Here is a list of incremental estimators for different tasks:

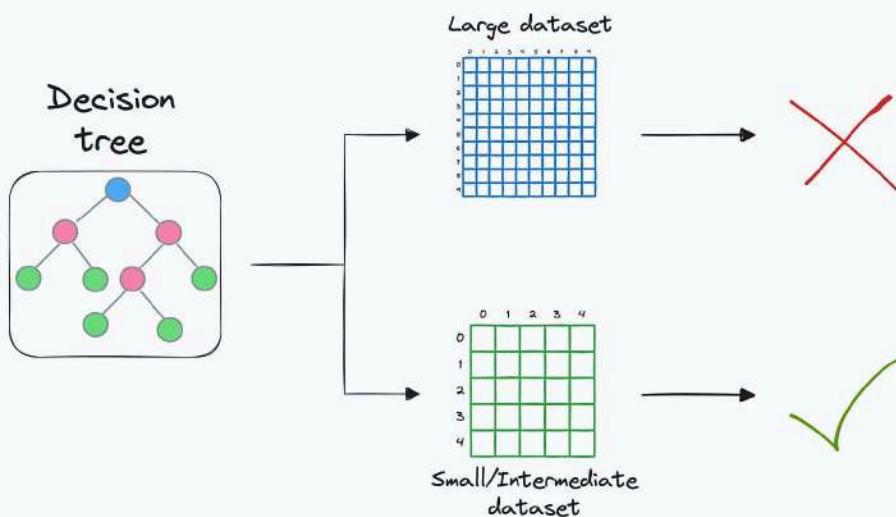
- Classification
 - `sklearn.naive_bayes.MultinomialNB`
 - `sklearn.naive_bayes.BernoulliNB`
 - `sklearn.linear_model.Perceptron`
 - `sklearn.linear_model.SGDClassifier`
 - `sklearn.linear_model.PassiveAggressiveClassifier`
- Regression
 - `sklearn.linear_model.SGDRegressor`
 - `sklearn.linear_model.PassiveAggressiveRegressor`
- Clustering
 - `sklearn.cluster.MiniBatchKMeans`
- Decomposition / feature Extraction
 - `sklearn.decomposition.MiniBatchDictionaryLearning`
 - `sklearn.cluster.MiniBatchKMeans`

Sklearn
implementations
that support
batch
implementation

It's pretty small, isn't it?

This is concerning because, in the enterprise space, the data is primarily tabular. Classical ML techniques, such as tree-based ensemble methods, are frequently used for modeling.

However, typical implementations of these models are not “big-data-friendly” because they require the entire dataset to be present in memory.



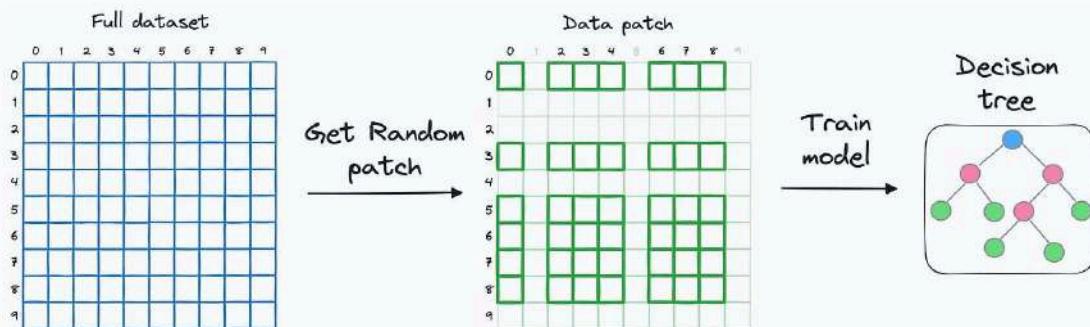
There are two ways to approach this:

1. Use big-data frameworks like Spark MLlib to train them.
2. There's one more way, which [Dr. Gilles Louppe](#) discussed in his PhD thesis — [Understanding Random Forests](#).

Here's what he proposed.

Random Patches

Before explaining, note that this approach will only work in an ensemble setting. So, you would have to train multiple models. The idea is to sample random data patches (rows and columns) and train a tree model on them.



Repeat this step multiple times by generating different patches of data randomly to obtain the entire random forest model.

The efficacy?

The thesis presents many benchmarks (check pages 174 and 178 if you need more details) on 13 datasets, and the results are shown below:

dataset	→													
Random patches	<table border="1"> <tr><td>46.13</td><td>98.87</td><td>98.91</td><td>96.57</td><td>95.20</td><td>80.0</td><td>71.62</td><td>81.36</td><td>88.20</td><td>98.57</td><td>93.45</td><td>92.06</td><td>97.52</td></tr> </table>	46.13	98.87	98.91	96.57	95.20	80.0	71.62	81.36	88.20	98.57	93.45	92.06	97.52
46.13	98.87	98.91	96.57	95.20	80.0	71.62	81.36	88.20	98.57	93.45	92.06	97.52		
Typical random forest	<table border="1"> <tr><td>46.17</td><td>98.77</td><td>98.82</td><td>96.56</td><td>95.07</td><td>77.76</td><td>69.91</td><td>81.41</td><td>87.92</td><td>98.46</td><td>93.36</td><td>91.53</td><td>97.44</td></tr> </table>	46.17	98.77	98.82	96.56	95.07	77.76	69.91	81.41	87.92	98.46	93.36	91.53	97.44
46.17	98.77	98.82	96.56	95.07	77.76	69.91	81.41	87.92	98.46	93.36	91.53	97.44		

From left to right → Cifar10, mnist3v8, mnist4v9, mnist, isolet, arcene, breast2, madelon, marti, reged, second, this, and sido.

From the above image, it is clear that in most cases, the random patches approach performs better than the traditional random forest.

In other cases, there wasn't a significant difference in performance.

And this is how we can train a random forest model on large datasets that do not fit into memory.

Why does it work?

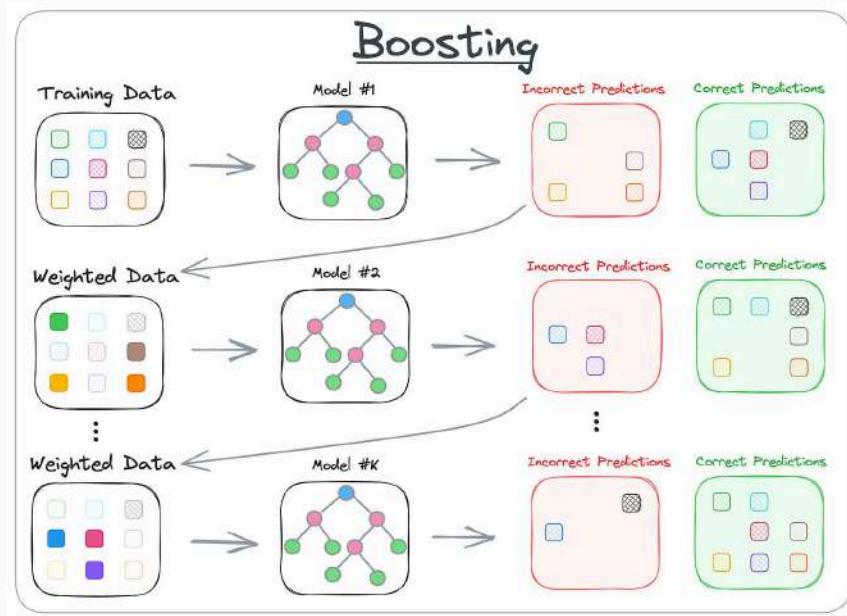
While the thesis did not provide a clear intuition behind this, I can understand why such an approach would still be as effective as random forest.

In a gist, building trees that are as different as possible guarantees a greater reduction in variance.

In this case, the dataset overlap between any two trees is NOT expected to be huge compared to the typical random forest. This aids in the Bagging objective.

A Visual Guide to AdaBoost

The following visual summarizes how Boosting models work:



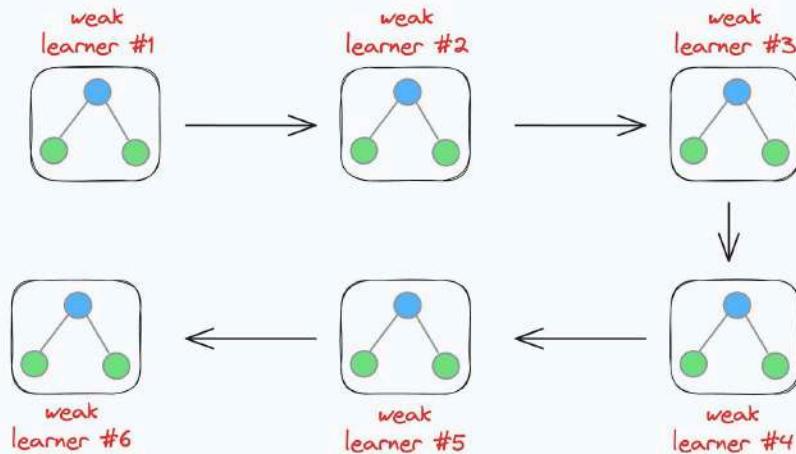
- Boosting is an iterative training process.
- The subsequent model puts more focus on misclassified samples from the previous model.
- The final prediction is a weighted combination of all predictions

However, many find it difficult to understand how this model is precisely trained and how instances are reweighed for subsequent models. AdaBoost is a common Boosting model, so in this chapter, let's understand how it works.

AdaBoost internal working

The core idea behind Adaboost is to train many weak learners to build a more powerful model. This technique is also called ensembling.

Specifically talking about Adaboost, the weak classifiers progressively learn from the previous model's mistakes, creating a powerful model when considered as a whole.



These weak learners are usually decision trees.

Let me make it more clear by implementing AdaBoost using the `DecisionTreeClassifier` class from sklearn.

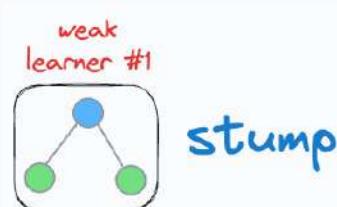
Consider we have the following classification dataset:

feature 1	feature 2	feature 3	target	weight
1.2	2.3	True	class A	0.2
-0.3	0.6	False	class B	0.2
0.7	1.8	False	class A	0.2
4.5	-3.5	True	class B	0.2
-0.4	0.4	True	class A	0.2

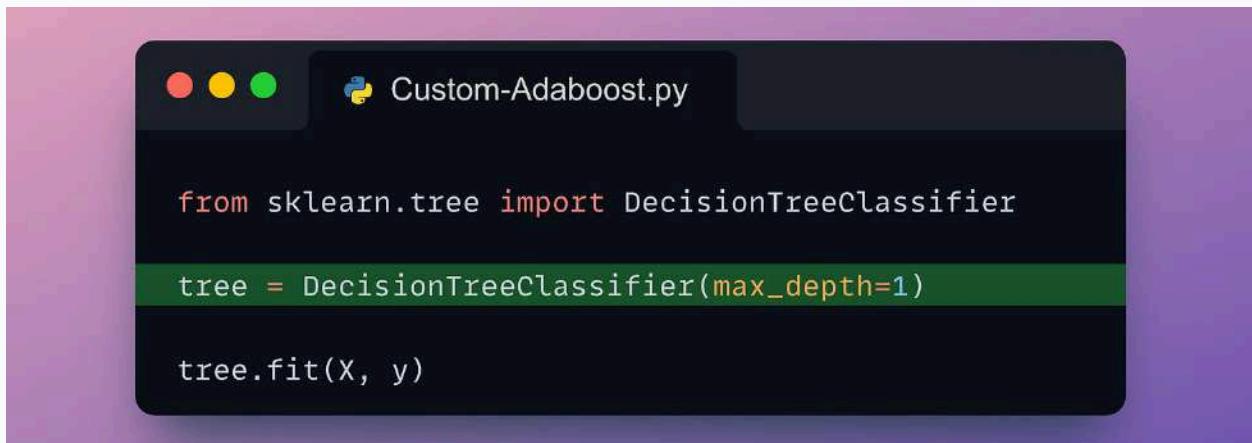
To begin, every row has an equal weight, and it is equal to $(1/n)$, where n is the number of training instances.

Step 1: Train a weak learner

In Adaboost, every decision tree has a unit depth, and they are also called stumps.



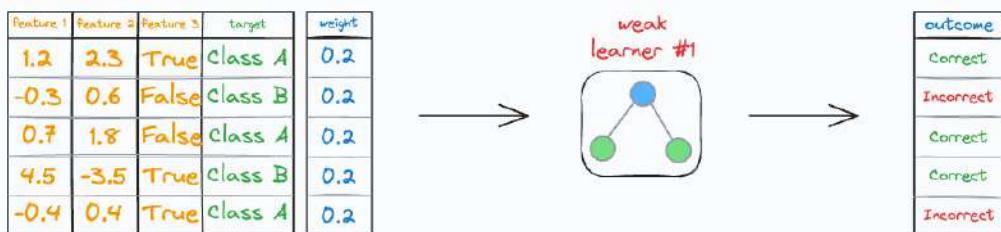
Thus, we define `DecisionTreeClassifier` with a `max_depth` of 1, and train it on the above dataset.



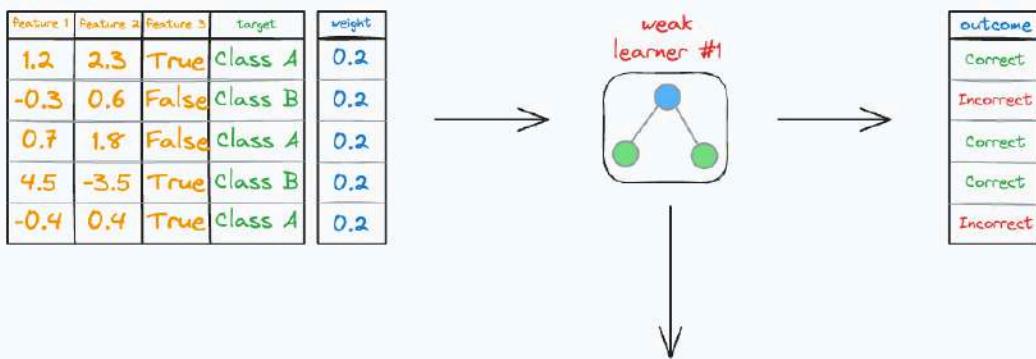
```
from sklearn.tree import DecisionTreeClassifier
tree = DecisionTreeClassifier(max_depth=1)
tree.fit(X, y)
```

Step 2: Calculate the learner's cost

Of course, there will be some correct and incorrect predictions.



The total cost (or error/loss) of this specific weak learner is the sum of the weights of the incorrect predictions. In our case, we have two errors, so the total error is:



Now, as discussed above, the idea is to let the weak learner progressively learn from previous learner's mistakes. So, going ahead, we want the subsequent model to focus more on the incorrect predictions produced earlier.

Here's how we do this:

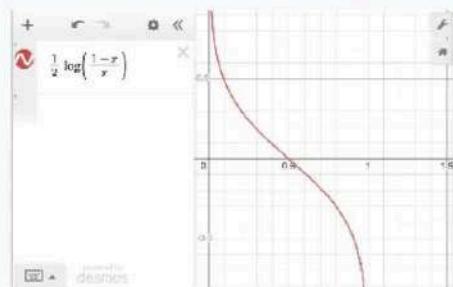
Step 3: Calculate the learner's importance

First, we determine the importance of the weak learner. Quite intuitively, we want the importance to be inversely related to the above error.

- If the weak learner has a high error, it must have a lower importance.
- If the weak learner has a low error, it must have a higher importance.

One choice is the following function:

$$\text{Importance} = \frac{1}{2} \cdot \log\left(\frac{1 - \text{Total error}}{\text{Total error}}\right)$$



- This function is only defined between [0,1].
- When the error is high (~1), this means there were no correct predictions → This gives a negative importance to the weak learner.
- When the error is low (~0), this means there were no incorrect predictions → This gives a positive importance to the weak learner.

If you feel there could be a better function to use here, you are free to use that and call it your own Boosting algorithm.

Now, we have the importance of the learner.

$$\text{Importance} = \frac{1}{2} \cdot \log\left(\frac{1 - \text{Total error}}{\text{Total error}}\right)$$

The importance value is used during model inference to weigh the predictions from the weak learners. So the next step is to...

Step 4: Reweigh the training instances

All the correct predictions are weighed down as follows:

$$w := w * e^{-\text{Importance}}$$

And all the incorrect predictions are weighed up as follows:

$$w := w * e^{\text{Importance}}$$

Once done, we normalize the new weights to add up to one.

Feature 1	Feature 2	Feature 3	target	weight	outcome	
1.2	2.3	True	Class A	0.2	Correct	
-0.3	0.6	False	Class B	0.2	Incorrect	
0.7	1.8	False	Class A	0.2	Correct	
4.5	-3.5	True	Class B	0.2	Correct	
-0.4	0.4	True	Class A	0.2	Incorrect	

reweigh and normalize →

weight
0.13
0.3
0.13
0.13
0.3

That's it!

Step 5: Sample from reweighed dataset

From step 4, we have the reweighed dataset.

We sample instances from this dataset in proportion to the new weights to create a new dataset.

feature 1	feature 2	feature 3	target	weight
1.2	2.3	True	Class A	0.13
-0.3	0.6	False	Class B	0.3
-0.3	0.6	False	Class B	0.3
-0.4	0.4	True	Class A	0.3
-0.4	0.4	True	Class A	0.3

Next, go back to step 1 — Train the next weak learner.

feature 1	feature 2	feature 3	target	weight
1.2	2.3	True	Class A	0.13
-0.3	0.6	False	Class B	0.3
-0.3	0.6	False	Class B	0.3
-0.4	0.4	True	Class A	0.3
-0.4	0.4	True	Class A	0.3

weak learner #2

And repeat the above process over and over for some pre-defined max iterations.

That's how we build the AdaBoost model.

All we have to do is consider the errors from the previous model, reweigh and sample the training instances for the next model, and repeat.

Dimensionality Reduction

The Utility of 'Variance' in PCA

The core objective of PCA is to retain the maximum variance of the original data while reducing the dimensionality. The rationale is that if we retain variance, we will retain maximum information.

But why?

Many people struggle to intuitively understand the motivation for using “variance” here. In other words:

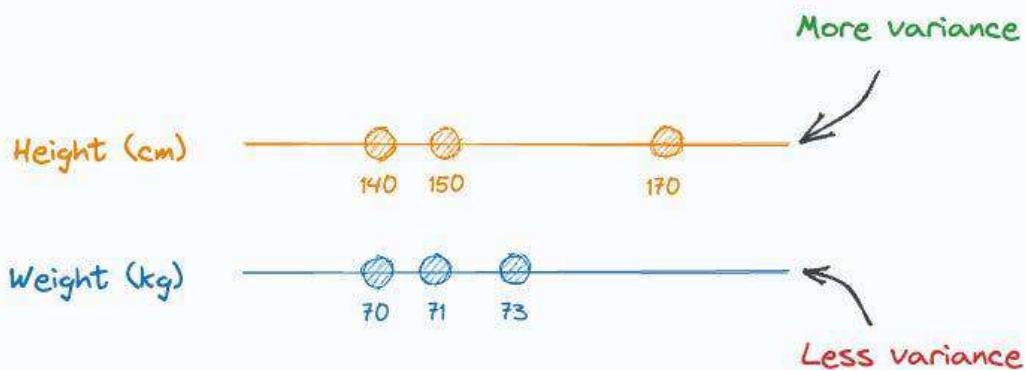
Why retaining maximum variance is an indicator for retaining maximum information?

This chapter provides an intuitive explanation of this.

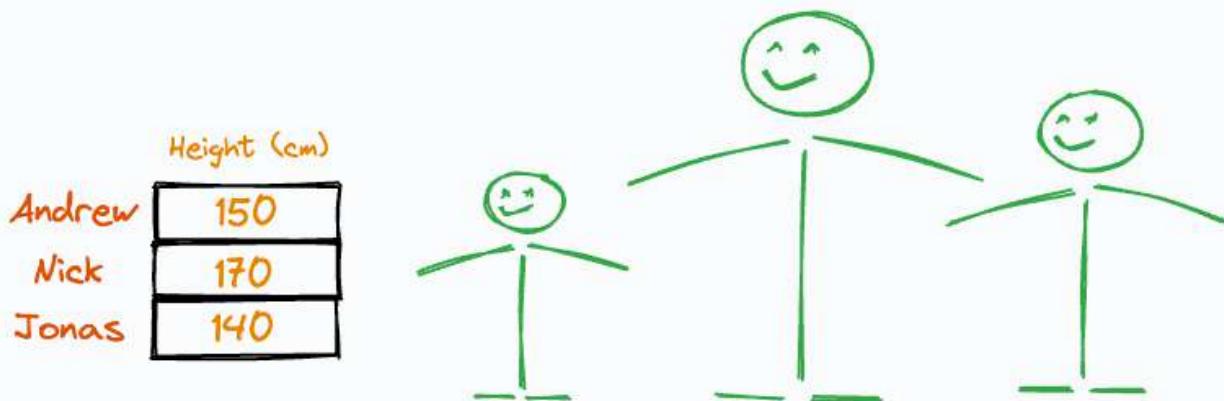
Imagine someone gave us the following weight and height information about three individuals:

	Height (cm)	Weight (kg)
Andrew	150	71
Nick	170	73
Jonas	140	70

It's clear that the height column has more variation than weight.



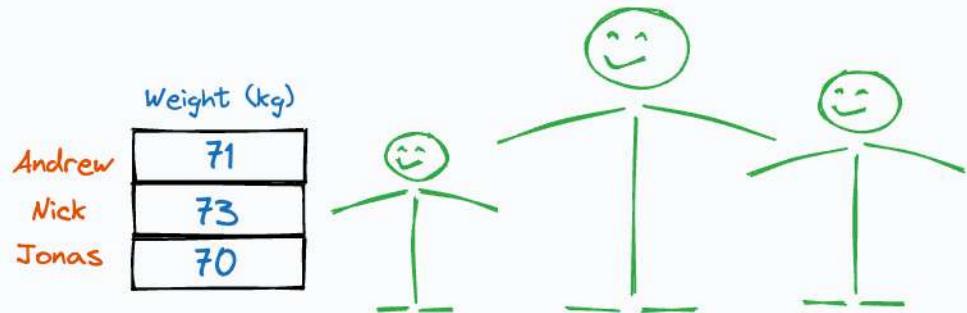
Thus, even if we discard the weight column, we can still identify these people solely based on their heights.



Dropping the weight column

- The one in the middle is Nick.
- The leftmost person is Jonas.
- The rightmost one is Andrew.

That was super simple. But what if we discarded the height column instead?



Can you identify them now?

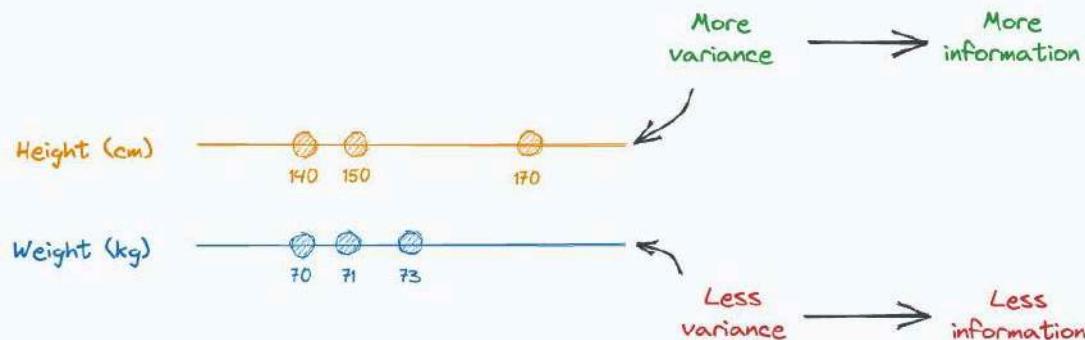
No, right?

And why?

This is because their heights have more variation than their weights.

And it's clear from the above example that, typically, if a column has more variation, it holds more information.

That is the core idea PCA is built around, and that is why it tries to retain maximum data variance. Simply put, PCA is devised on the premise that more variance means more information.

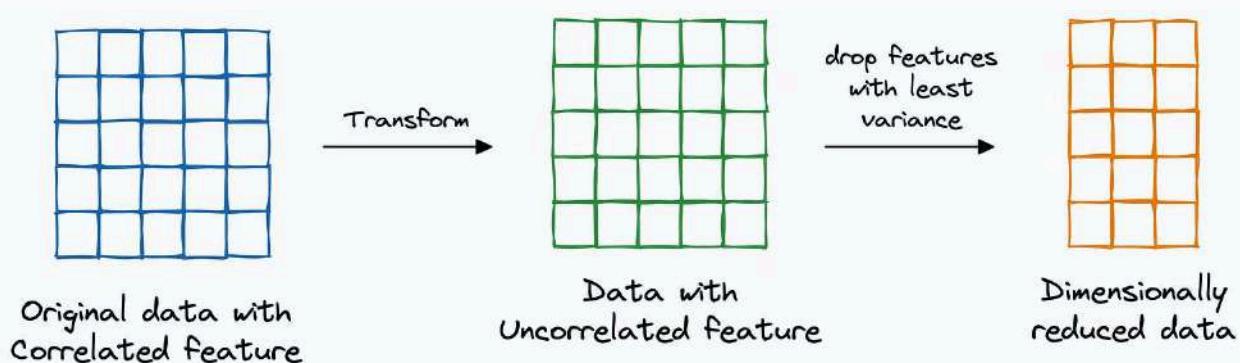


More variance means more information and less variance means less information. Thus, during dimensionality reduction, we can (somewhat) say that we are retaining maximum information if we retain maximum variance.

Of course, as we are using variance, this also means that it can be easily influenced by outliers. That is why we say that PCA is influenced by outliers.

As a concluding note, always remember that when using PCA, we don't just measure column-wise variance and drop the columns with the least variance.

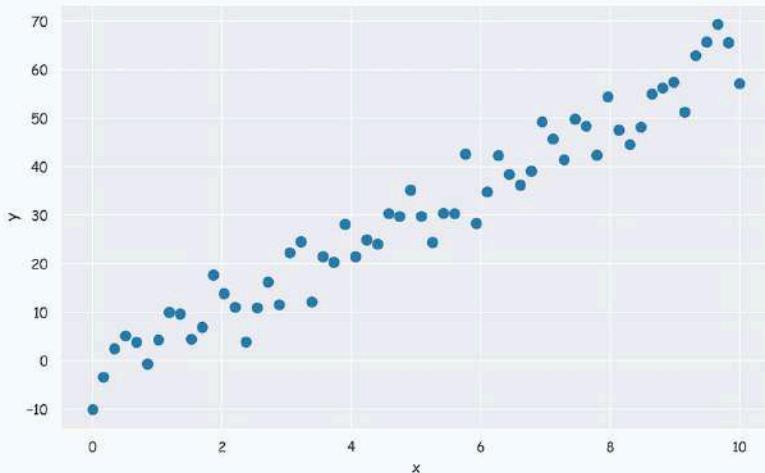
Instead, we must first transform the data to create uncorrelated features. After that, we drop the new features based on their variance.



KernelPCA vs. PCA

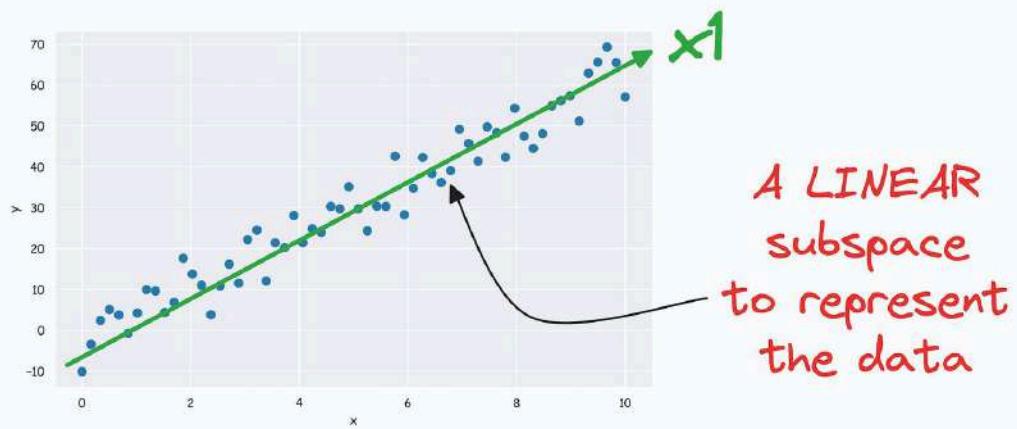
During dimensionality reduction, principal component analysis (PCA) tries to find a low-dimensional linear subspace that the given data conforms to.

For instance, consider the following dummy dataset:



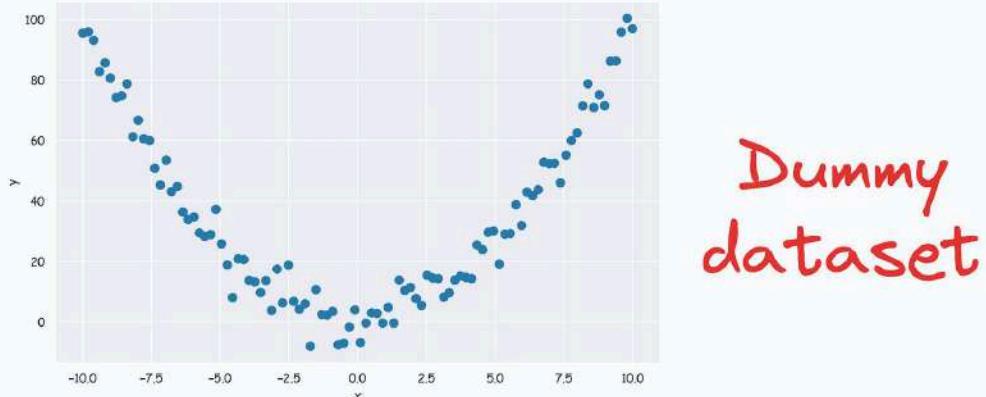
Dummy
dataset

It's pretty clear from the above visual that there is a linear subspace along which the data could be represented while retaining maximum data variance. This is shown below:



But what if our data conforms to a low-dimensional yet non-linear subspace.

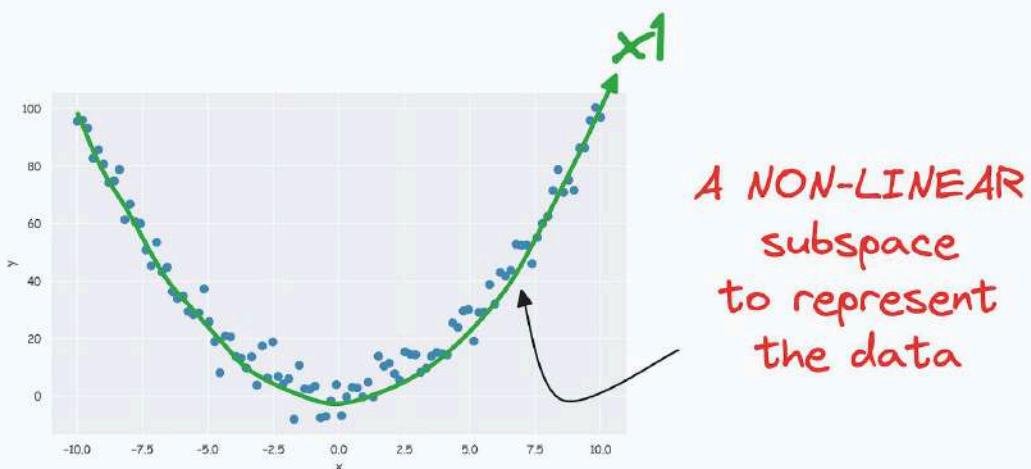
For instance, consider the following dataset:



Do you see a low-dimensional non-linear subspace along which our data could be represented?

No?

Don't worry. Let me show you!



The above curve is a continuous non-linear and low-dimensional subspace that we could represent our data given along.

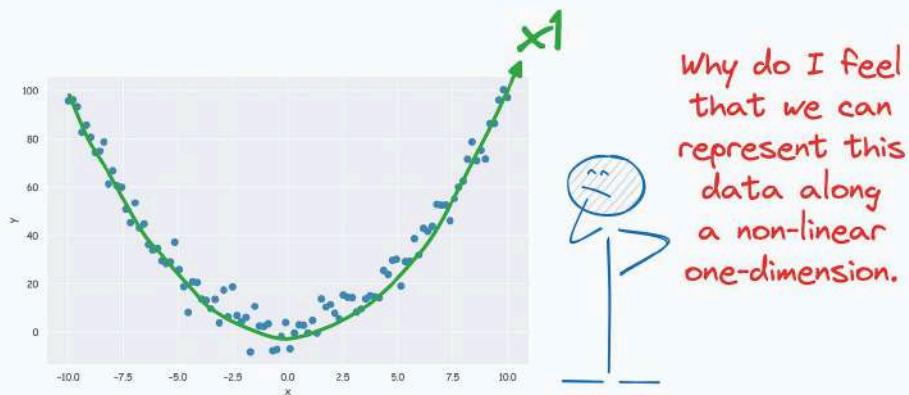
Okay...so why don't we do it then?

The problem is that PCA cannot determine this subspace because the data points are non-aligned along a straight line.

In other words, PCA is a linear dimensionality reduction technique.

Thus, it falls short in such situations.

Nonetheless, if we consider the above non-linear data, don't you think there's still some intuition telling us that this dataset can be reduced to one dimension if we can capture this non-linear curve.



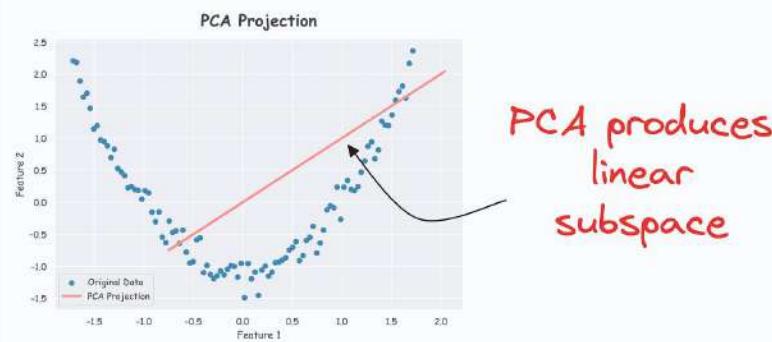
KernelPCA (which uses the kernel trick) precisely addresses this limitation of PCA.

The idea is pretty simple:

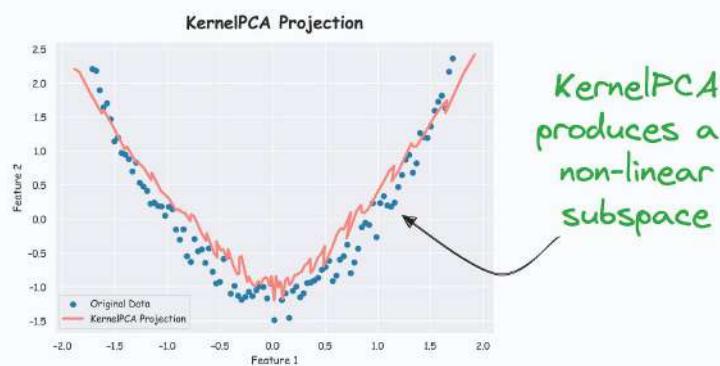
- Project the data to another high-dimensional space using a kernel function, where the data becomes linearly representable. Sklearn provides a KernelPCA wrapper, supporting many popularly used kernel functions.
- Apply the standard PCA algorithm to the transformed data.

The efficacy of KernelPCA over PCA is evident from the demo below.

As shown below, even though the data is non-linear, PCA still produces a linear subspace for projection:



However, KernelPCA produces a non-linear subspace:



What's the catch, you might be wondering?

The catch is the run time.

Please note that the run time of PCA is already cubically related to the number of dimensions.

$$O(nm^2 + m^3)$$

n → data points
 m → dimensions

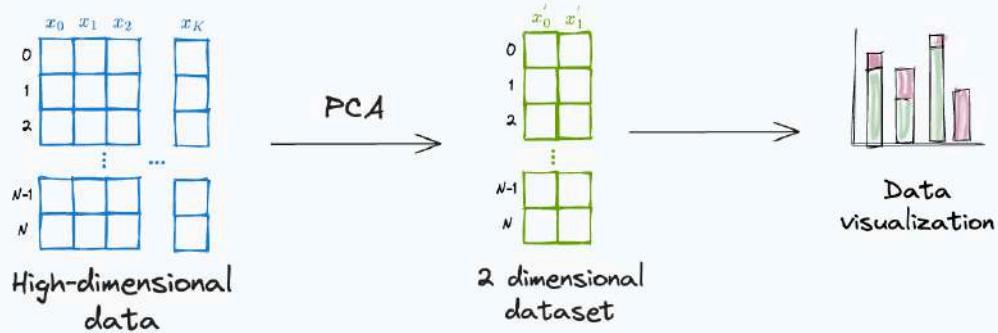
KernelPCA involves the kernel trick, which is quadratically related to the number of data points (n).

Thus, it increases the overall run time.

This is something to be aware of when using KernelPCA.

PCA is not a Visualization Technique

PCA, by its very nature, is a dimensionality reduction technique. Yet, at times, many use PCA for visualizing high-dimensional datasets. This is done by projecting the given data into two dimensions and visualizing it.



While this may appear like a fair thing to do, there's a big problem here that often gets overlooked.

To understand this problem, we first need to understand a bit about how PCA works.

How PCA works?

The core idea in PCA is to linearly project the data to another space using the eigenvectors of the covariance matrix.

$$X' = X \lambda$$

Transform
data

Why eigenvectors?

1. It creates uncorrelated features, which is useful because if features are independent, the features with the least variance can be dropped for dimensionality reduction.
2. It ensures that new features collectively preserve the original data variance.

Coming back to the visualization topic...

As discussed above, after applying PCA, each new feature captures a fraction of the original data variance. Thus, if we intend to use PCA for visualization by projecting the data to 2-dimensions...

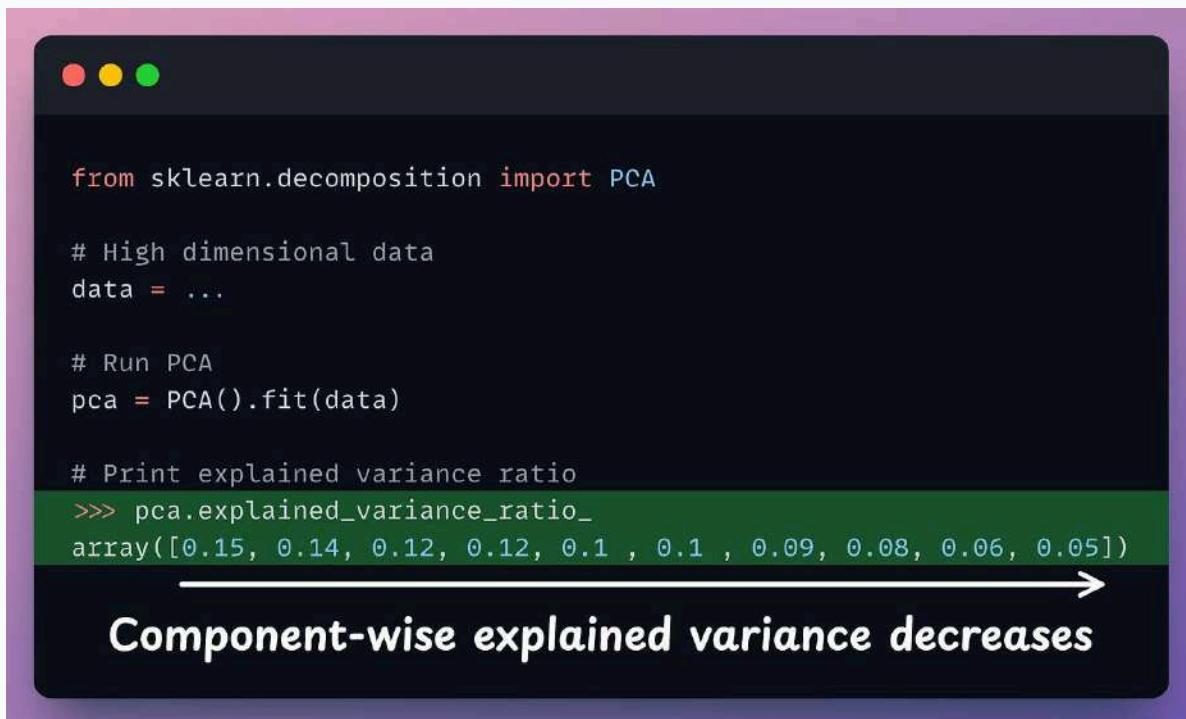
...then this visualization will only be useful if the first two principal components collectively capture most of the original data variance.

If they don't, then the two-dimensional visualization will be highly misleading and incorrect.

How to determine the variance contribution of the first two principal components?

We can avoid this mistake by plotting a cumulative explained variance (CEV) plot. As the name suggests, it plots the cumulative variance explained by principal components.

In sklearn, for instance, the explained variance fraction is available in the `explained_variance_ratio_` attribute:



```

from sklearn.decomposition import PCA

# High dimensional data
data = ...

# Run PCA
pca = PCA().fit(data)

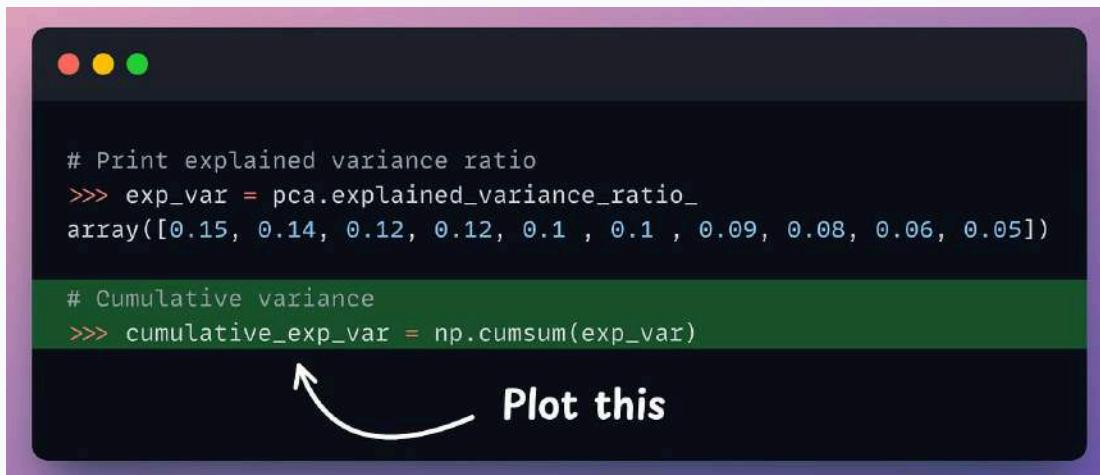
# Print explained variance ratio
>>> pca.explained_variance_ratio_
array([0.15, 0.14, 0.12, 0.12, 0.1 , 0.1 , 0.09, 0.08, 0.06, 0.05])

```

→

Component-wise explained variance decreases

We can create a cumulative plot of explained variance and check whether the first two components explain the majority of variance.

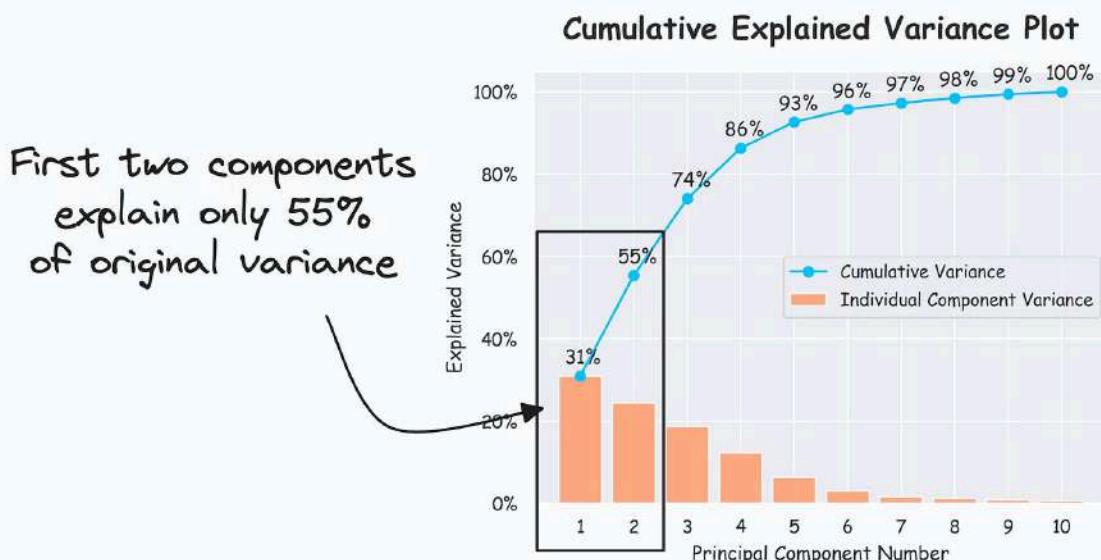


```
# Print explained variance ratio
>>> exp_var = pca.explained_variance_ratio_
array([0.15, 0.14, 0.12, 0.12, 0.1 , 0.1 , 0.09, 0.08, 0.06, 0.05])

# Cumulative variance
>>> cumulative_exp_var = np.cumsum(exp_var)
```

Plot this

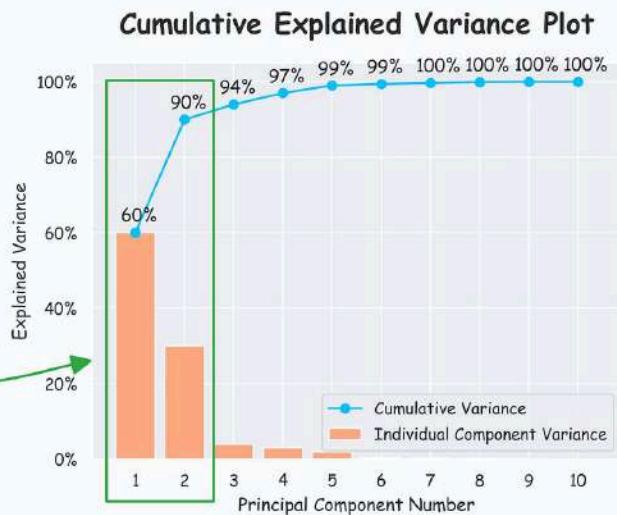
For instance, in the plot below, the first two components only explain 55% of the original data variance.



Thus, visualizing this dataset in 2D using PCA may not be a good choice because plenty of data variance is missing.

However, in the below plot, the first two components explain 90% of the original data variance.

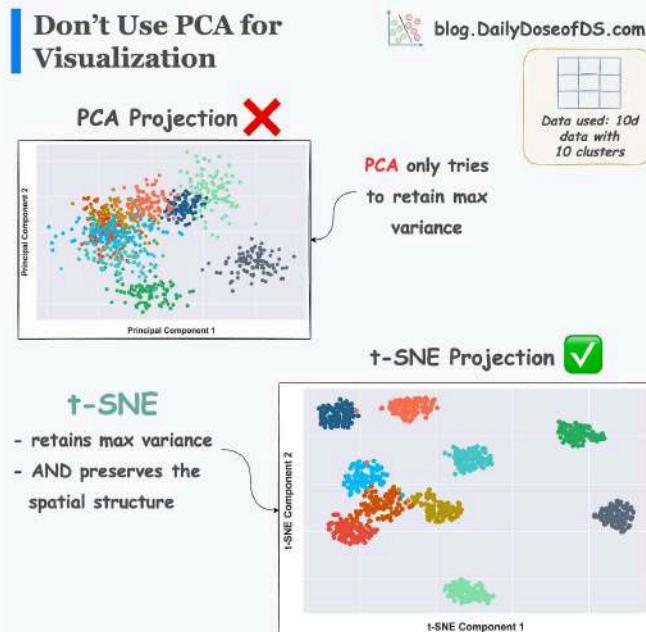
First two components explain 90% of original variance



Thus, using PCA for visualization looks like a fair thing to do. As a takeaway, use PCA for 2D visualization only when the above plot suggests so.

If it does not, then refrain from using PCA for 2D visualization and use other techniques specifically meant for visualization, like t-SNE, UMAP, etc. The efficacy of t-SNE over PCA is depicted below:

2D projections created by PCA do not consider local structure. Instead, its principal components primarily focus on preserving the maximum variance. But t-SNE projections provide much more clarity into data clusters. The clusters produced by t-SNE are well separated. But those in PCA have a significant overlap, and they hardly convey anything.



t-SNE vs. SNE – What's the difference?

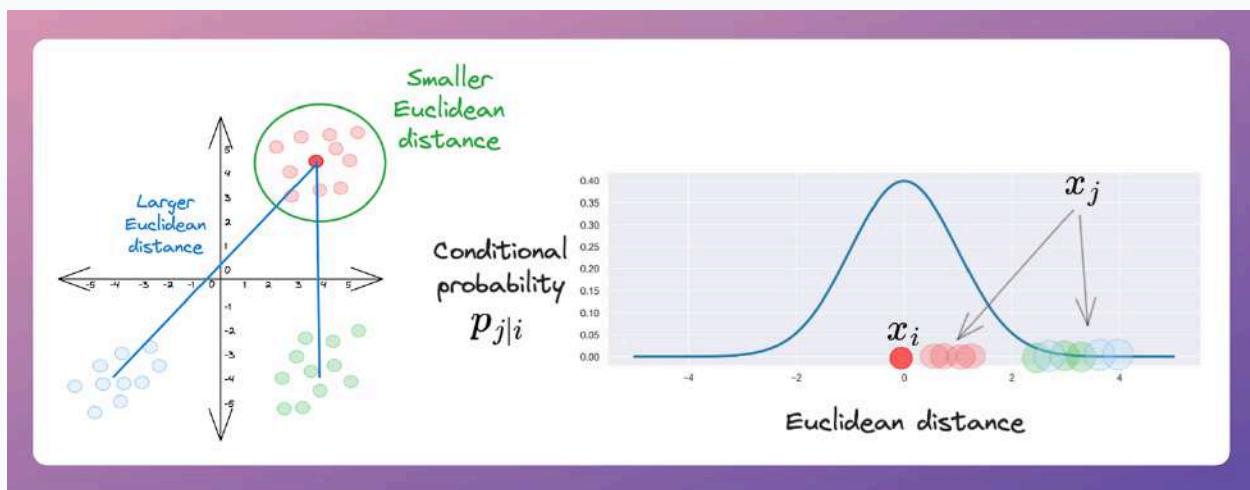
In this chapter, let's continue discussing the t-SNE algorithm.

The t-SNE algorithm is an improved version of the SNE algorithm, both used for dimensionality reduction.

SNE —————→ **t-SNE**

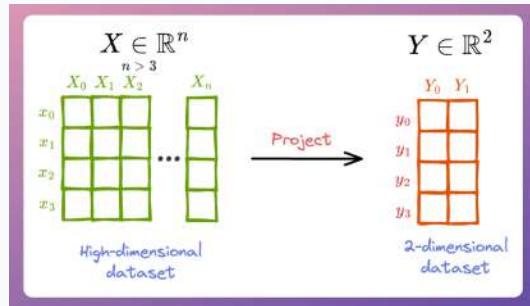
The core idea of SNE (not t-SNE) is the following:

- Step 1) For every point (i) in the given high-dimensional data, convert the high-dimensional Euclidean distances to all other points (j) into conditional Gaussian probabilities.
 - For instance, consider the marked red point in the dataset on the left below.

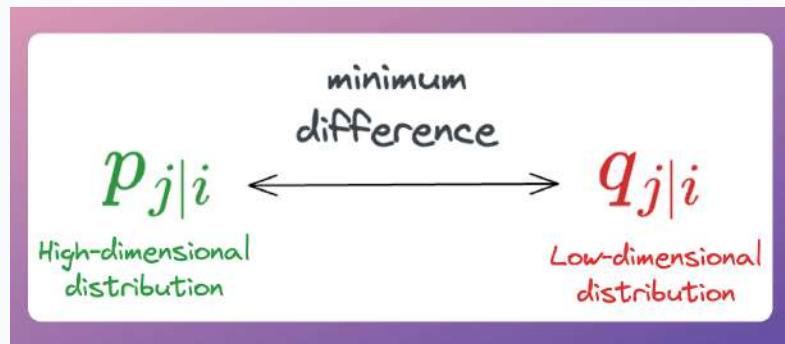


- Converting Euclidean distance to all other points into Gaussian probabilities (the distribution on the right above) shows that other red points have a higher probability of being its neighbor than other points.

- Step 2) For every data point x_i , randomly initialize its counterpart y_i in 2-dimensional space. These will be our projections.



- Step 3) Just like we defined conditional probabilities in the high-dimensional space in Step 1, we define the conditional probabilities in the low-dimensional space, using Gaussian distribution again.
- Step 4) Now, every data point (i) has a high-dimensional probability distribution and a corresponding low-dimensional distribution:



- The objective is to match these two probability distributions. Thus, we can make the positions of counterpart y_i 's learnable such that this difference is minimized.
- Using KL divergence as a loss function helps us achieve this. It measures how much information is lost when we use distribution Q to approximate distribution P.

$$D_{KL}(P \parallel Q) = \sum_x P(x) \cdot \log \left(\frac{P(x)}{Q(x)} \right)$$

- Ideally, we want to have the minimum loss value (which is zero), and this will be achieved when $P=Q$.

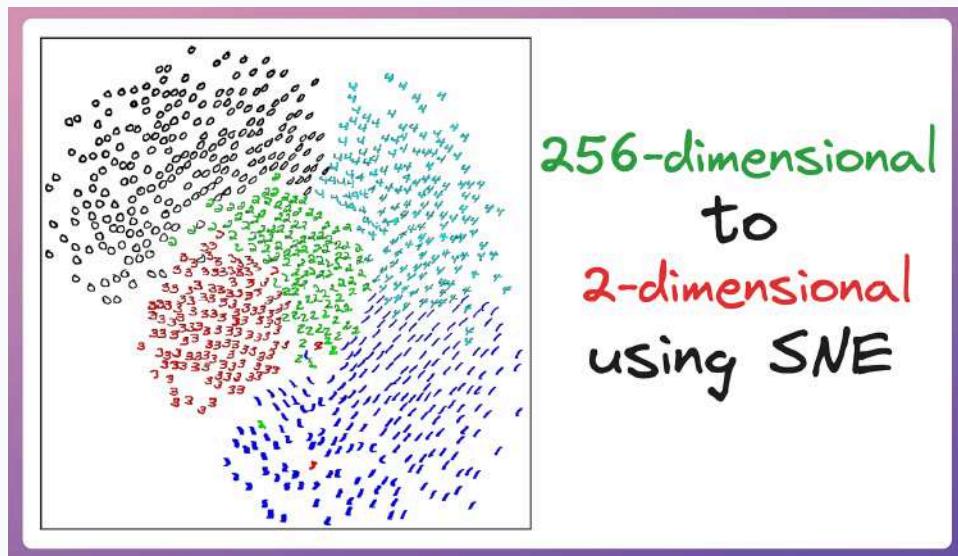
$$D_{KL}(P \parallel Q) = \sum_x P(x) \cdot \log \left(\frac{P(x)}{Q(x)} \right)$$

if $P = Q$

loss = 0

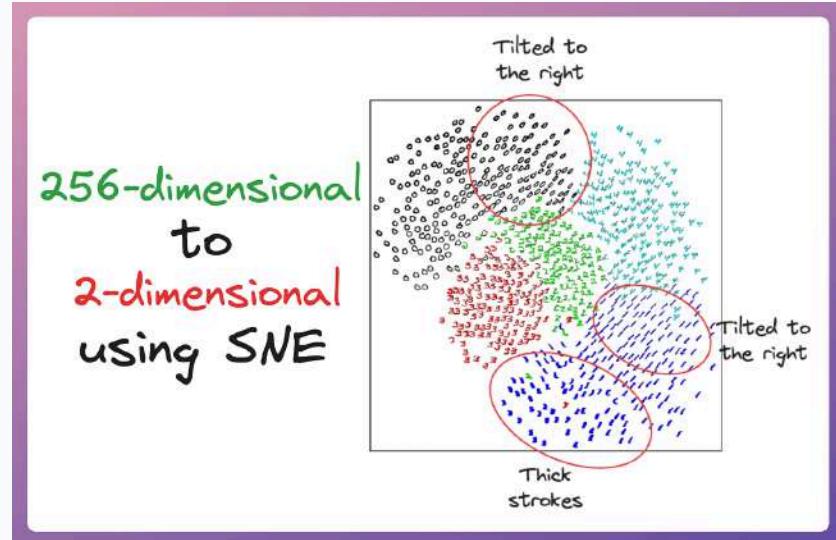
The model can be trained using gradient descent, and it works pretty well.

For instance, the following image depicts a 2-dimensional visualization produced by the SNE algorithm on 256-dimensional handwritten digits:

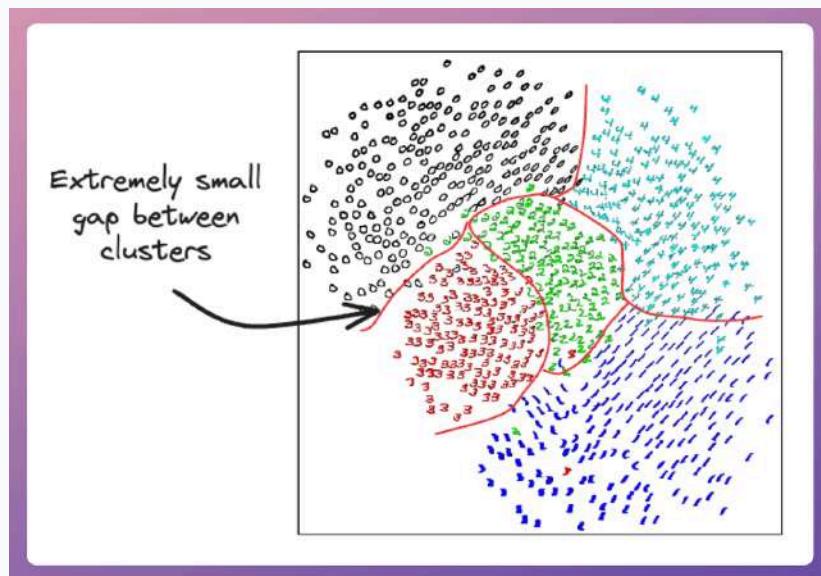


SNE produces good clusters.

What's even more astonishing is that properties like orientation, skew, and strokethickness vary smoothly across the space within each cluster. This is depicted below:



Nonetheless, it has some limitations, which the t-SNE algorithm addresses.
Notice that the clusters produced by SNE are not well separated.



Here, it could be fair to assume that the original data clusters, the one in the 256-dimensional space, most likely would have been well separated. Thus:

- All zeros must have been together but well separated from other digits.
- All ones must have been together but well separated from other digits.
- And so on.

Yet, SNE still produces tightly packed clusters. This is also called the “crowding problem.”

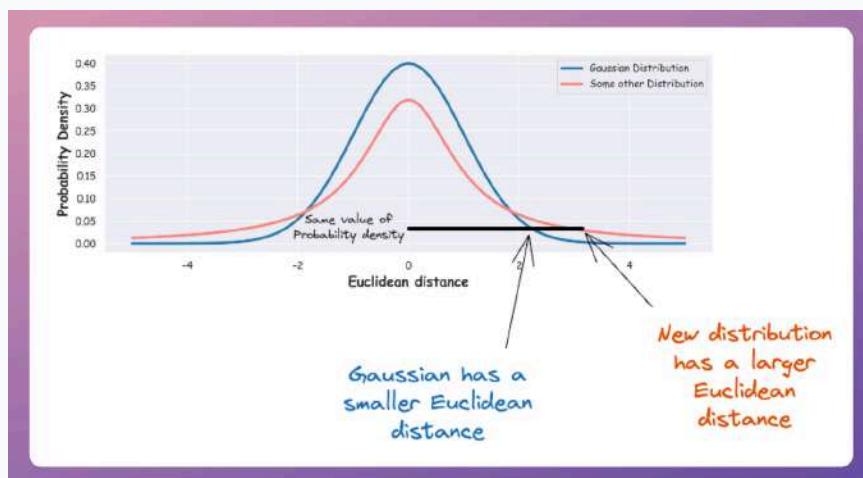
To eliminate this problem, t-SNE was proposed, standing for t-distributed Stochastic Neighbor Embedding (t-SNE).

Here's the difference.

Recall that in SNE, we used a Gaussian distribution to define the low-dimensional conditional probabilities. But it is not producing well-separated clusters.

One solution is to use some other probability distribution, such that for distant points, we get the same value of the conditional probability as we would have obtained from a Gaussian distribution but at a larger Euclidean distance.

Let me simplify that a bit. Compare the following two distributions:

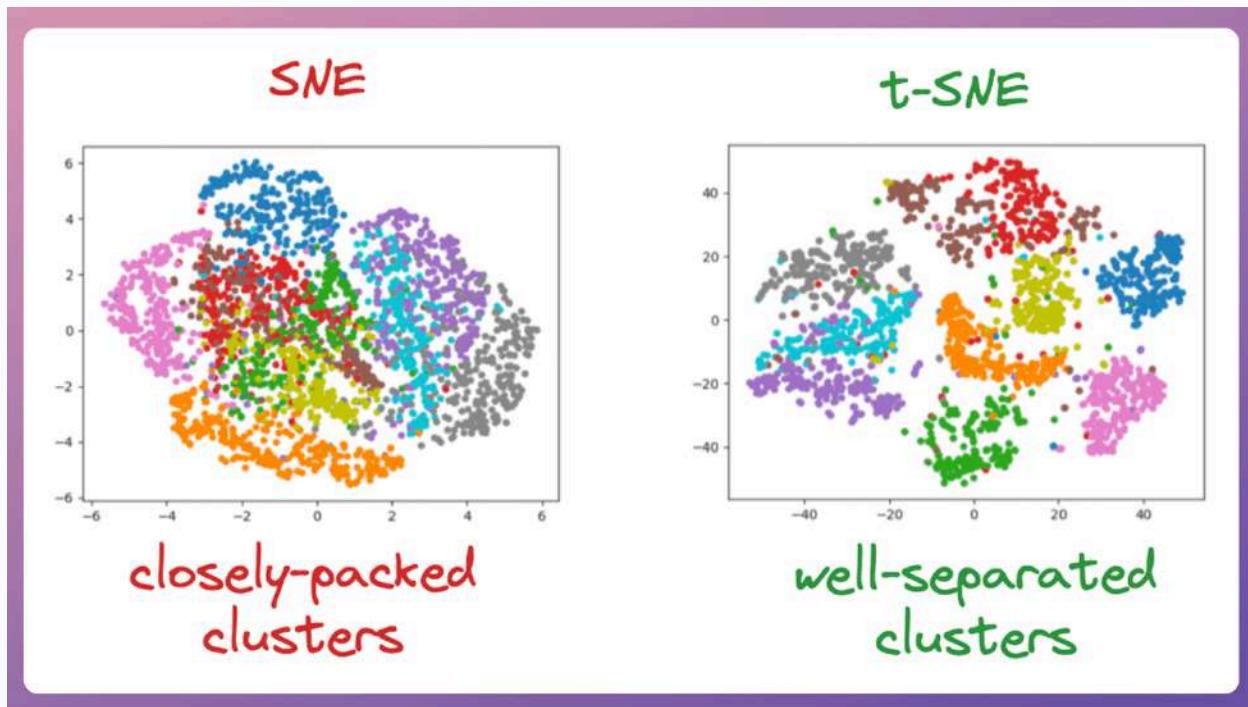


Notice that Gaussian achieves a specific value of low probability density at a smaller distance. But t-distribution achieves it at a larger distance.

This is precisely what we intend to achieve.

We need a heavier tail distribution so that we can still minimize the difference between the two probability distributions but at a larger distance in the low-dimensional space.

The Student t-distribution is a perfect fit for it. The following image depicts the difference this change brings:



As shown above:

- SNE produces closely packed clusters.
- t-SNE produces well-separated clusters.

And that's why t-distribution is used in t-SNE.

That said, besides producing well-separated clusters, using the Student t-distribution has many more advantages.

For instance, it is computationally much faster to evaluate the density of a point under a Student t-distribution than under a Gaussian.

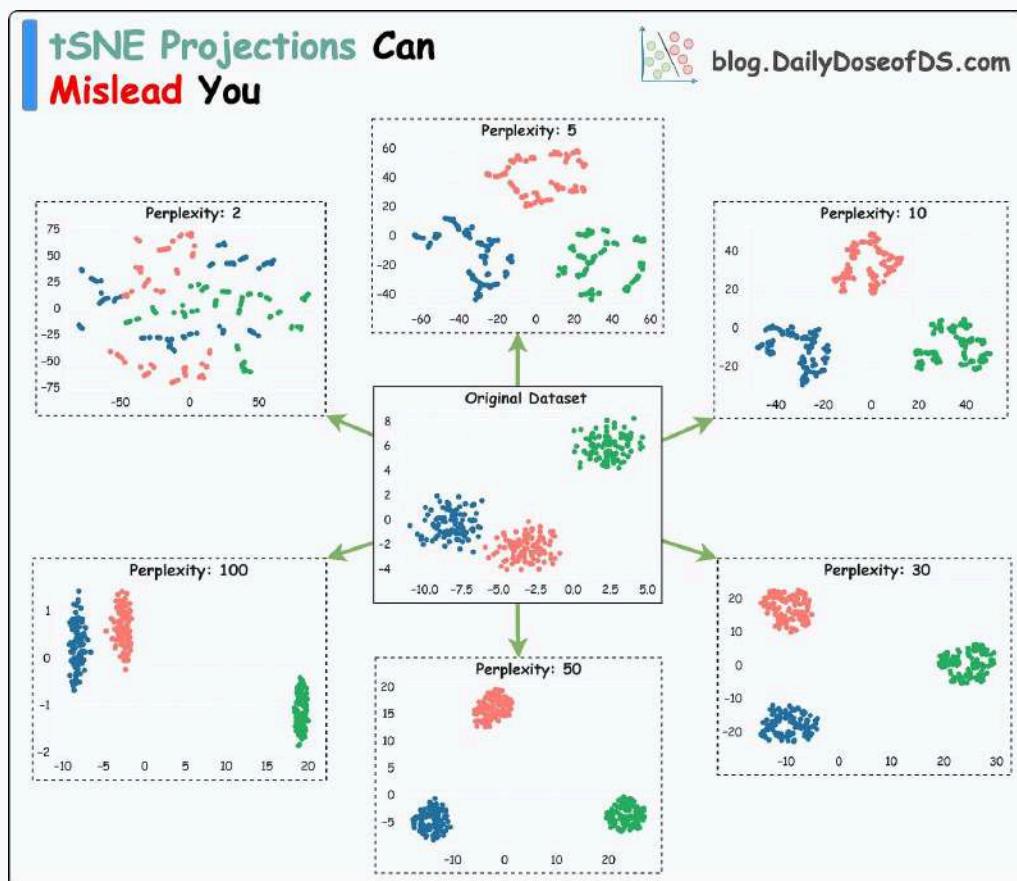
How To Avoid Getting Misled by t-SNE Projections?

From the above discussion, t-SNE sounds promising, doesn't it? While the algorithm is quite powerful, many consistently make misleading conclusions from the t-SNE projections of their high-dimensional data.

In this chapter, I want to point out a few of these mistakes so that you don't make those mistakes ever.

To begin, the performance of the t-SNE algorithm is primarily reliant on **perplexity** — a hyperparameter of t-SNE. That is why it is considered the most important hyperparameter in the t-SNE algorithm.

Simply put, the perplexity value provides a rough estimate for the number of neighbors a point may have in a cluster. And different values of perplexity create very different low-dimensional cluster spaces, as depicted below:



As shown above, most projections do depict the original clusters. However, they vary significantly in shape.

There are five takeaways from the above image:

- NEVER make any conclusions about the original cluster shape by looking at these projections.
 - Different projections have different low-dimensional cluster shapes, and they do not resemble the original cluster shape.
 - For low perplexity values (5 and 10), cluster shapes significantly differ from the original ones.
 - Although, in this case, the clusters were color-coded, which provided more clarity. But it may not always be the case, as tSNE is an unsupervised algorithm.
- Cluster sizes in a t-SNE plot do not convey anything either.
- The dimensions (or coordinates of data points) created by t-SNE in low dimensions have no inherent meaning.
 - The axes tick labels of the low-dimensional plots are different and somewhat random.
 - Similar to PCA's principal components, they offer no interpretability.
- The distances between clusters in a projection do not mean anything.
 - In the original dataset, the blue and red clusters are close.
 - Yet, most projections do not preserve the global structure of the original dataset.
- Strange things happen at **perplexity=2** and **perplexity=100**.
 - At **perplexity=2**, the low-dimensional mapping conveys nothing.
 - As discussed earlier, the perplexity value provides a rough estimate of the number of neighbors a point may have in a cluster.
 - t-SNE tries to maintain approx. 2 points per cluster. That is why the distortion.

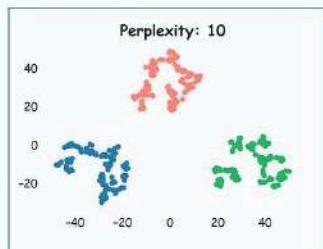
- At $\text{perplexity}=100$, the global structure is preserved, but the local structure gets distorted.
- Thus, tweaking the perplexity hyperparameter is extremely critical here.
- That is why I mentioned above that it is the most important hyperparameter of this algorithm.

As a concluding note, it is found that the ideal perplexity values typically lie in the range [5,50]. So try experimenting in that range and see what looks promising.

Next time you use t-SNE, consider the above points, as these plots can get tricky to interpret. This is especially true if you don't understand the internal workings of this algorithm.

Nonetheless, understanding the algorithm will massively help you develop an intuition on its interpretability.

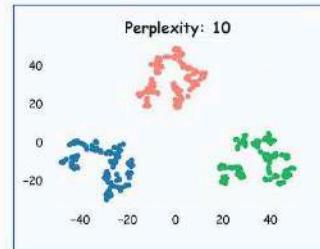
Before understanding



What does
this mean?



After understanding tSNE

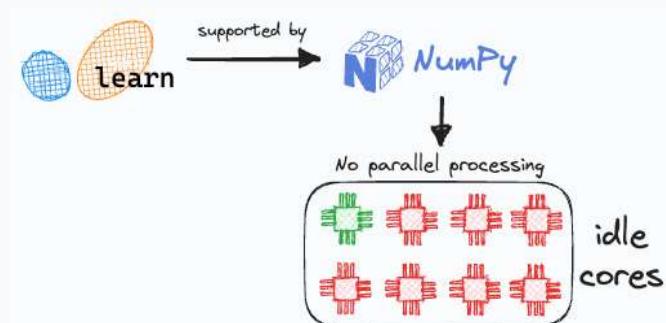


Aha! I know
what this means.

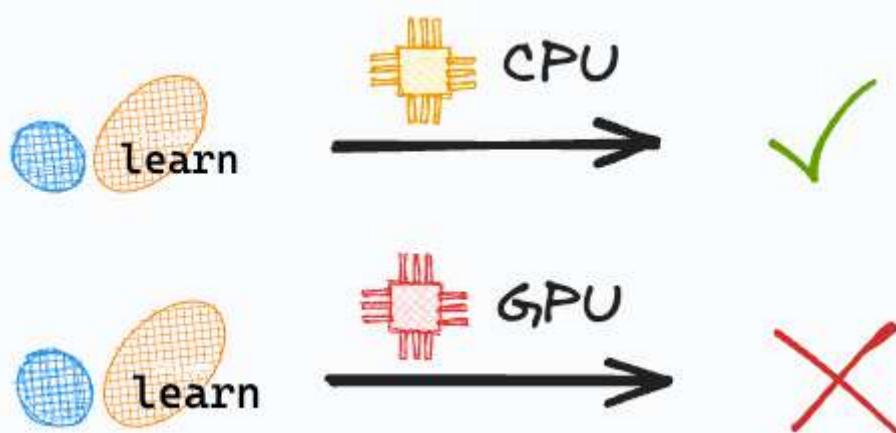


Accelerate tSNE with GPU

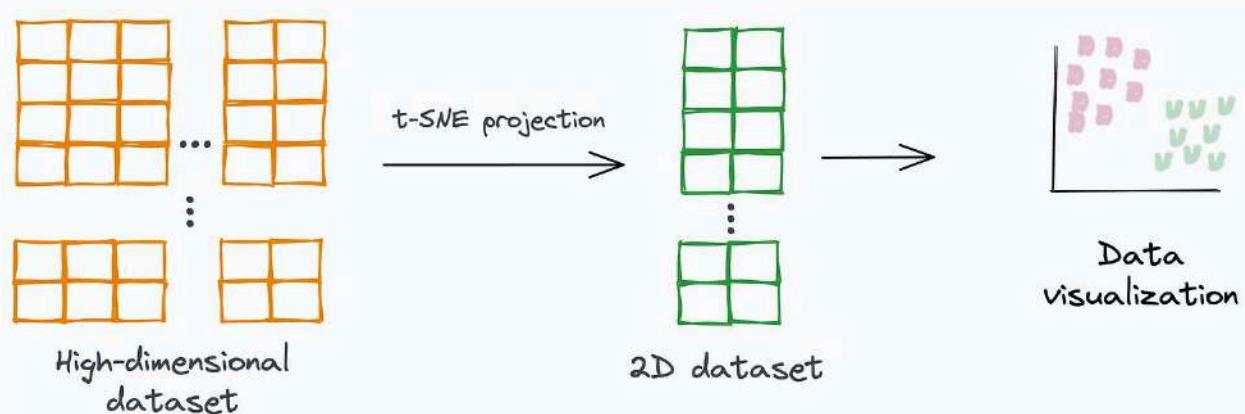
Sklearn implementations are driven by NumPy, which typically run on a single core of a CPU. Thus, the ability to run parallelized operations is quite limited at times.



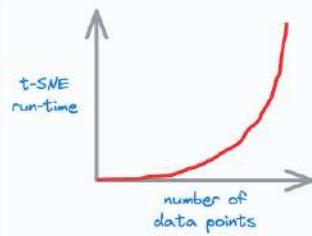
Another major limitation is that scikit-learn models cannot run on GPUs.



This bottleneck provides massive room for run-time improvement. The same applies to the tSNE algorithm, which is among the most powerful dimensionality reduction techniques to visualize high-dimensional datasets.



Because the biggest issue with tSNE is that its run-time is quadratically related to the number of data points.



Thus, beyond, say, 20k-25k data points, it becomes pretty difficult to use tSNE from Sklearn implementations. There are two ways to handle this:

- Either keep waiting.
- Or use optimized implementations that could be possibly accelerated with GPU.

I was stuck due to the same issue in one of my projects, and I found a pretty handy solution.

[tSNE-CUDA](#) is an optimized CUDA version of the tSNE algorithm, which, as the name suggests, can leverage hardware accelerators. As a result, it provides immense speedups over the standard Sklearn implementation:

```
from sklearn.manifold import TSNE
X_projected = TSNE().fit_transform(X)
```

tSNE Sklearn

Run-time: 100 seconds X

```
from tsnecuda import TSNE
X_projected = TSNE().fit_transform(X)
```

tSNE GPU

Run-time: 3 seconds ✓ (33x Faster)

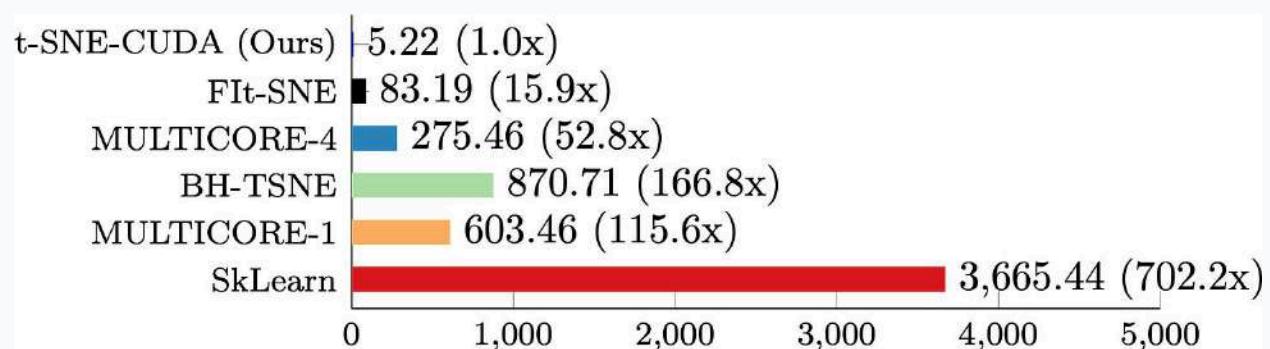
As depicted above, the GPU-accelerated implementation:

- Is 33 times faster than the Sklearn implementation.
- Produces similar quality clustering as the Sklearn implementation.

Do note that this implementation only supports `n_components=2`, i.e., you can only project to two dimensions. As per the docs, the authors have no plans to support more dimensions, as this will require significant changes to the code.

But in my opinion, that doesn't matter because, for more than 99% of cases, tSNE is used to obtain 2D projections. So we are good here.

Before I conclude this chapter, I also found the following benchmarking results by the authors:



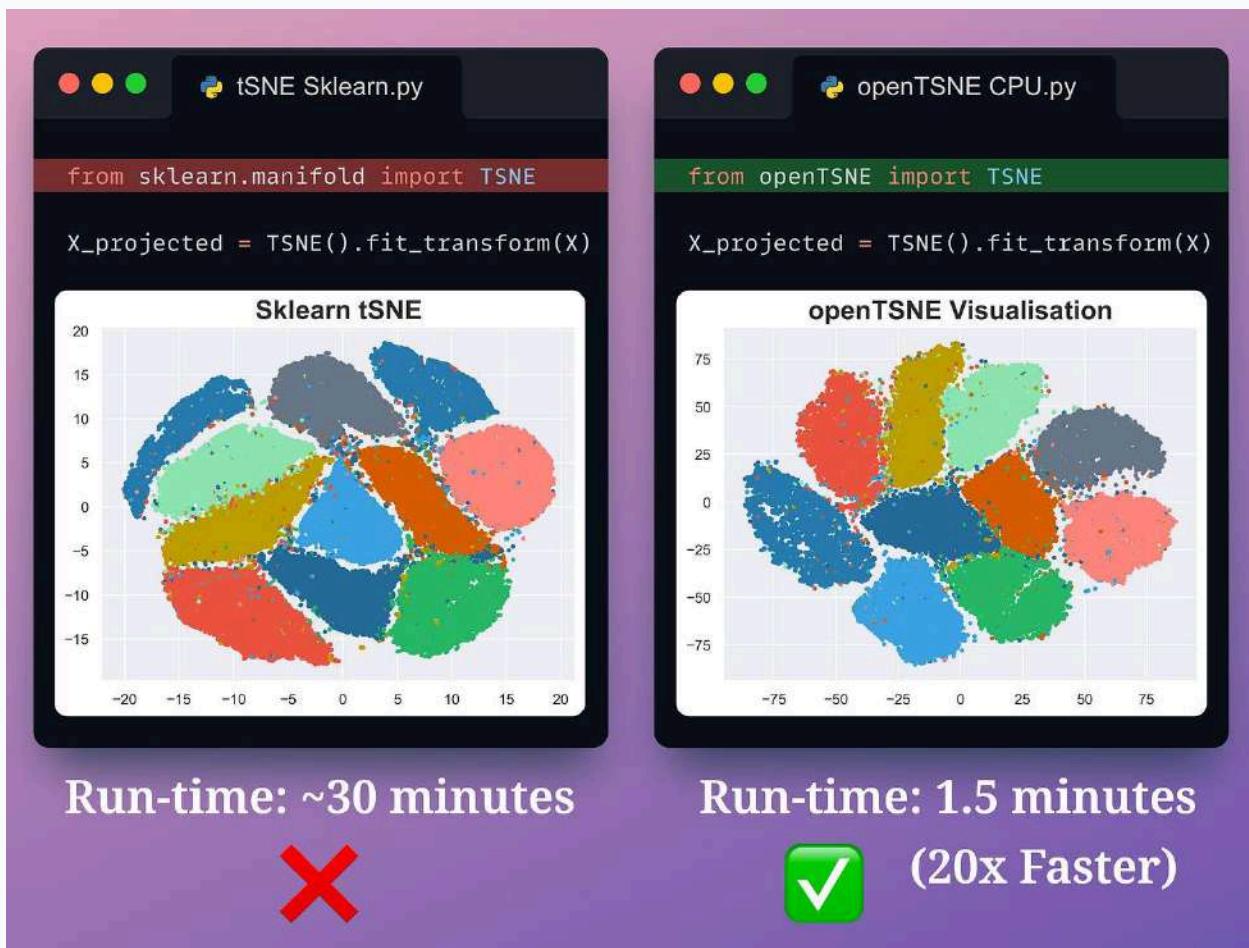
It depicts that on the CIFAR-10 training set (50k images), tSNE-CUDA is 700x Faster than Sklearn, which is an insane speedup.

Here's my Colab notebook for you to get started: <https://bit.ly/3WdEone>.

Scale tSNE to Millions of Data Points With openTSNE

tSNE-CUDA discussed above provides immense speedups over the standard Sklearn implementation using a GPU. But what if you don't have access to a GPU?

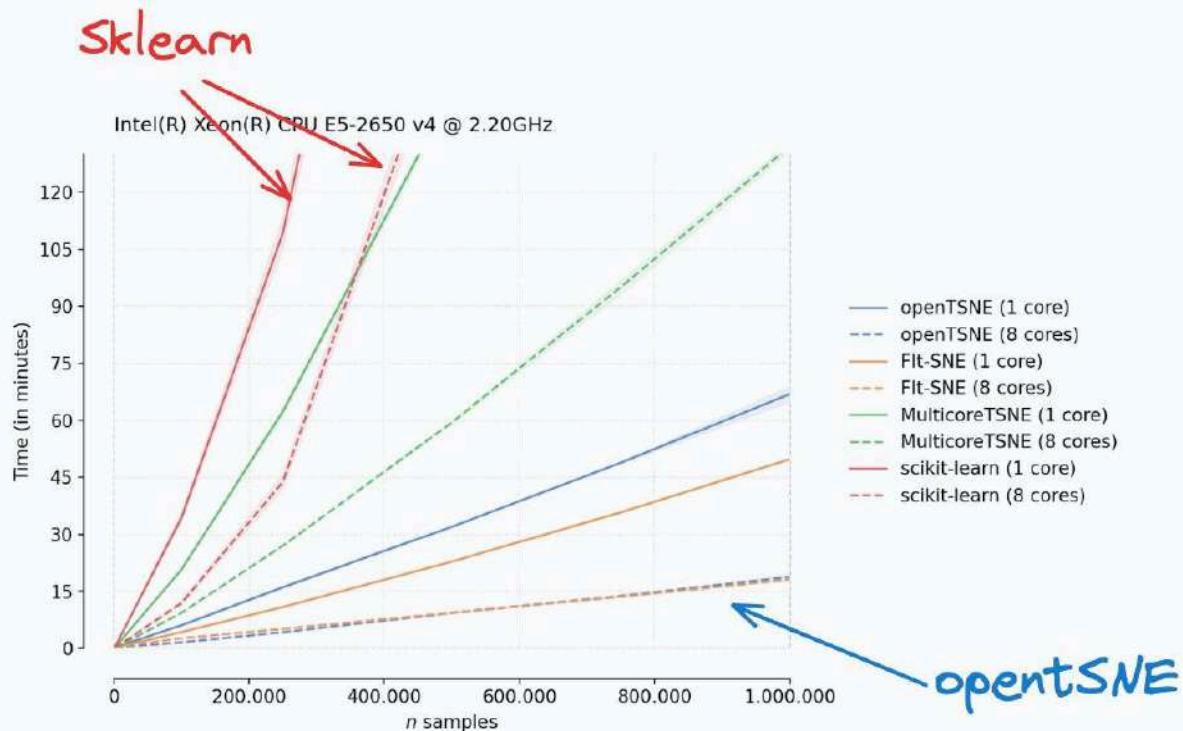
openTSNE is another optimized Python implementation of t-SNE, which provides massive speed improvements and enables us to scale t-SNE to millions of data points — a place where Sklearn implementation may never reach.



As depicted above, the openTSNE implementation:

- is 20 times faster than the Sklearn implementation.
- produces similar quality clustering as the Sklearn implementation.

The authors have also provided the following benchmarking results:



As depicted above, openTSNE can produce low dimensional visualization of a million data points in just ~15 minutes.

However, it is clear from their benchmarks that the run-time of the Sklearn implementation has already reached a couple of hours with just ~250k data points.

Download the notebook here to try out openTSNE: <https://bit.ly/3zAHCZ8>.

PCA vs. t-SNE

Finally, it would be good to look at the core differences between PCA and t-SNE algorithms. The following table for you which neatly summarizes the major differences between the two algorithms:

PCA vs. t-SNE

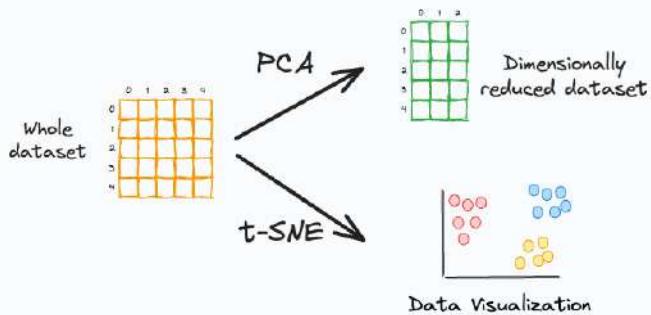
 blog.DailyDoseofDS.com

		PCA	t-SNE
Purpose		Dimensionality reduction	Data visualisation in low dimensions (typically 2D)
Type of algorithm		Deterministic	Stochastic
Unique solution?		Yes	No
Projection type		Linear	Non-linear
Type of approach		Global	Local/Global*
Interpretation		PCA is just a rotation of axes	Subjective interpretation

*We do not explicitly specify global structure preservation in tSNE.

#1) Purpose

- While many interpret PCA as a data visualization algorithm, it is primarily a dimensionality reduction algorithm.
- t-SNE, however, is a data visualization algorithm. We use it to project high-dimensional data to low dimensions (primarily 2D).

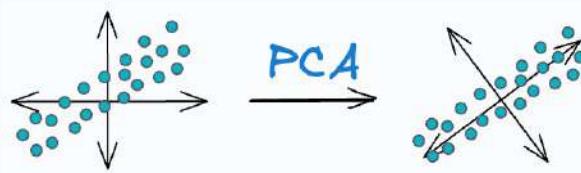


#2) Type of algorithm

- PCA is a deterministic algorithm. Thus, if we run the algorithm twice on the same dataset, we will ALWAYS get the same result.
- t-SNE, however, is a stochastic algorithm. Thus, rerunning the algorithm can lead to entirely different results.

#3) Uniqueness of solution

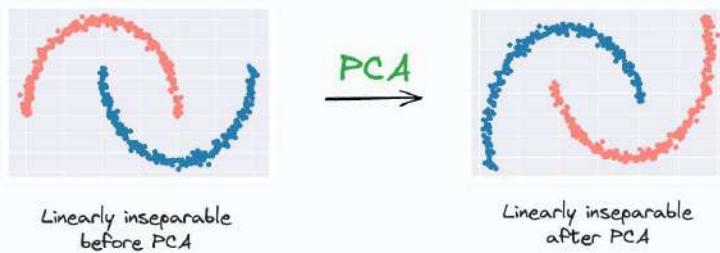
As far as uniqueness and interpretation of results is concerned...



- PCA always has a unique solution for the projection of data points. Simply put, PCA is just a rotation of axes such that the new features we get are uncorrelated.
- t-SNE, as discussed above, can provide entirely different results, and its interpretation is subjective in nature.

#4) Projection type

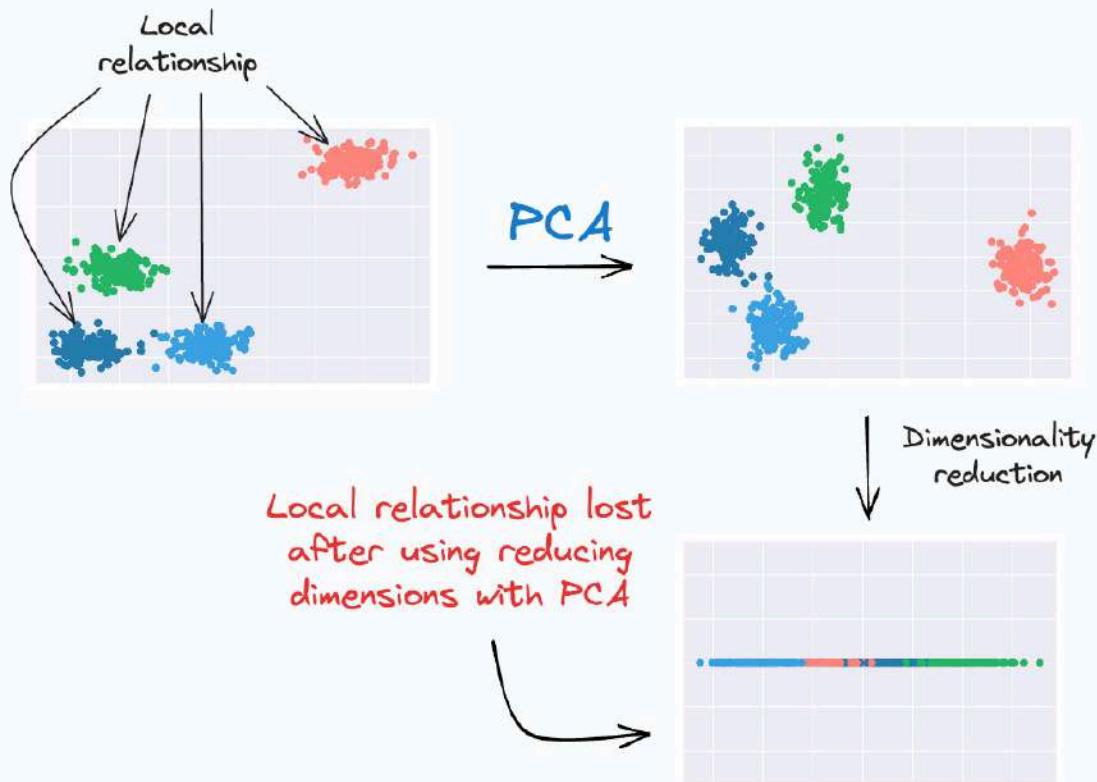
- PCA is a linear dimensionality reduction approach. It can only find a linear subspace to project the given dataset. KernelPCA addresses this, which we covered earlier in this section.



- t-SNE is a non-linear approach. It can handle non-linear datasets.

#5) Projection approach

- PCA only aims to retain the global variance of the data. Thus, local relationships (such as clusters) are often lost after projection, as shown below:



- t-SNE preserves local relationships. Thus, data points in a cluster in the high-dimensional space are much more likely to lie together in the low-dimensional space.
 - In t-SNE, we do not explicitly specify global structure preservation. But it typically does create well-separated clusters.
 - Nonetheless, as discussed in the previous chapter, the distance between two clusters in low-dimensional space is NEVER an indicator of cluster separation in high-dimensional space.

Clustering

Beyond KMeans: 6 Must-Know Types of Clustering Algorithms

Clustering is one of the core branches of unsupervised learning in ML. The first (and sometimes the only) clustering algorithm folks learn is KMeans. Yet, it is important to note that KMeans is not a universal solution to all clustering problems. In fact, there's a whole world of clustering algorithms beyond KMeans, which we must be familiar with. The visual below summarizes 6 different types of clustering algorithms:

Centroid-based: Cluster data based on proximity to centroids.

Connectivity-based: Cluster points based on proximity between clusters.

Density-based: Cluster points based on their density. It is more robust to clusters with varying densities and shapes than centroid-based clustering.

Graph-based: Cluster points based on graph distance.

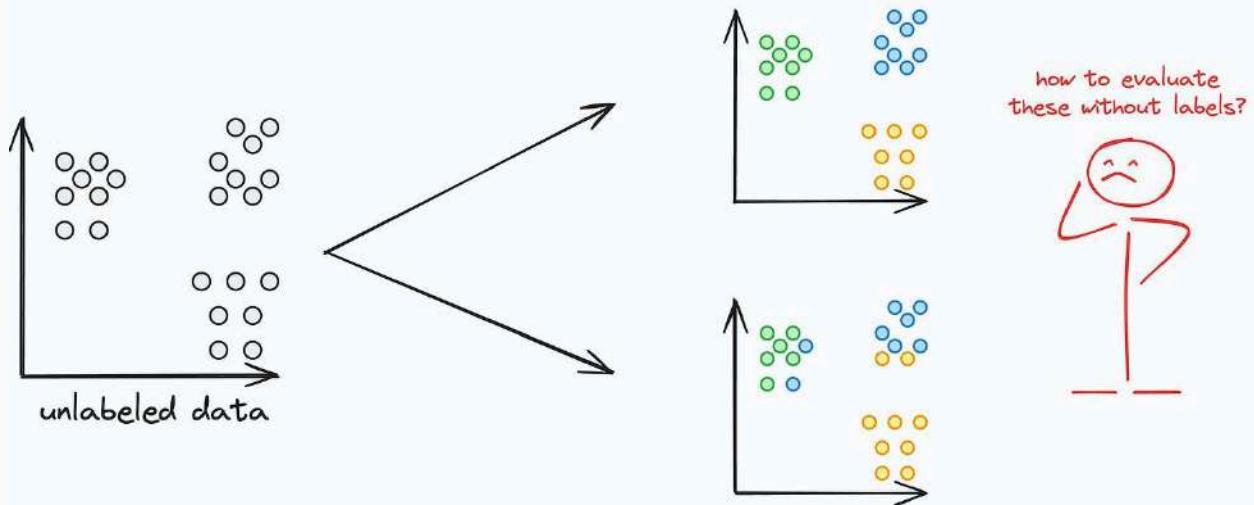
Distribution-based: Cluster points based on their likelihood of belonging to the same distribution. Gaussian Mixture Model is one example. We will discuss Gaussian mixture models (GMMs) in an upcoming chapter.

Compression-based: Transform data to a lower dimensional space and then perform clustering

6 Types of Clustering Algorithms in Machine Learning			
Clustering Algorithm Type	Clustering Methodology	Algorithm(s)	
	Centroid-based	Cluster points based on proximity to centroid	KMeans KMeans++ KMedoids
	Connectivity-based	Cluster points based on proximity between clusters	Hierarchical Clustering (Agglomerative and Divisive)
	Density-based	Cluster points based on their density instead of proximity	DBSCAN OPTICS HDBSCAN
	Graph-based	Cluster points based on graph distance	Affinity Propagation Spectral Clustering
	Distribution-based	Cluster points based on their likelihood of belonging to the same distribution.	Gaussian Mixture Models (GMMs)
	Compression-based	Transform data to a lower dimensional space and then perform clustering	BIRCH

Intrinsic Measures for Clustering Evaluation

Without labeled data, it is not possible to objectively measure how well the clustering algorithm has grouped similar data points together and separated dissimilar ones.



Also, most clustering datasets are multi-dimensional, so directly visualizing the results is not possible.

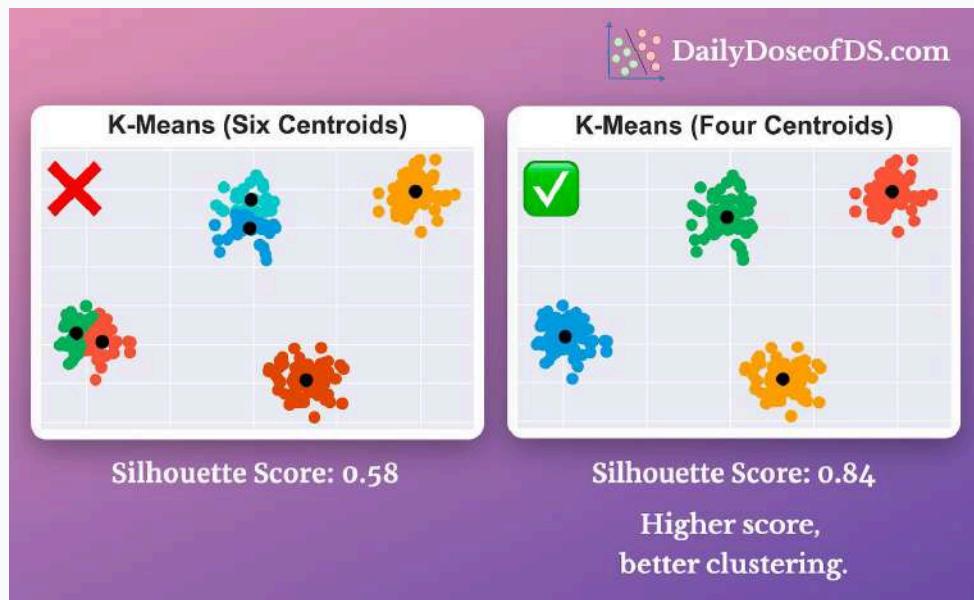
Thus, we must rely on intrinsic measures to determine clustering quality in such cases.

The way I like to use them is as follows:

- Say I am using KMeans:
 - Run KMeans with a range of k values.
 - Evaluate the performance.
 - Select the value of k based on the most promising cluster quality metric.

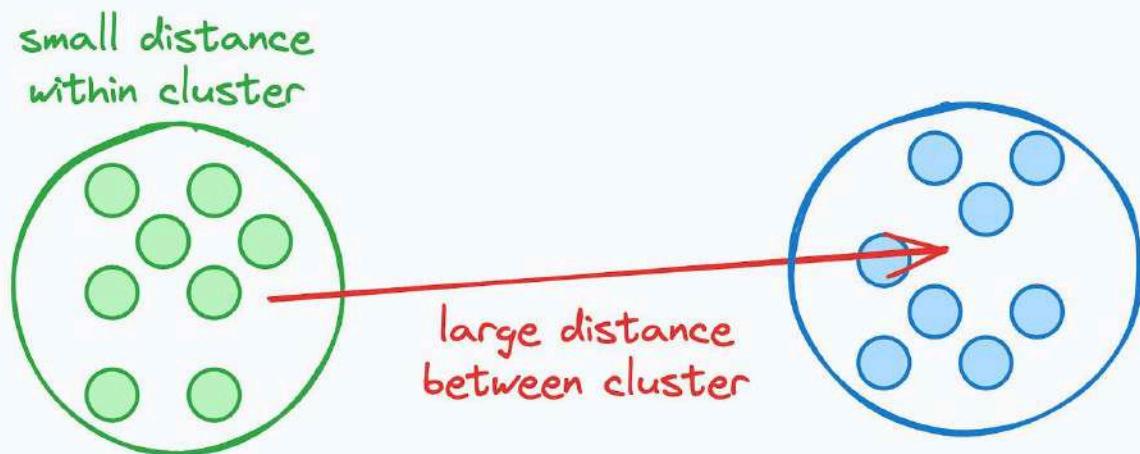
Let's understand a few of these metrics.

#1) Silhouette Coefficient:



The Silhouette Coefficient indicates how well each data point fits into its assigned cluster.

The idea is that if the average distance to all data points in the same cluster is small but that to another cluster is large, this intuitively indicates that the clusters are well separated and somewhat reliable.



It is measured as follows:

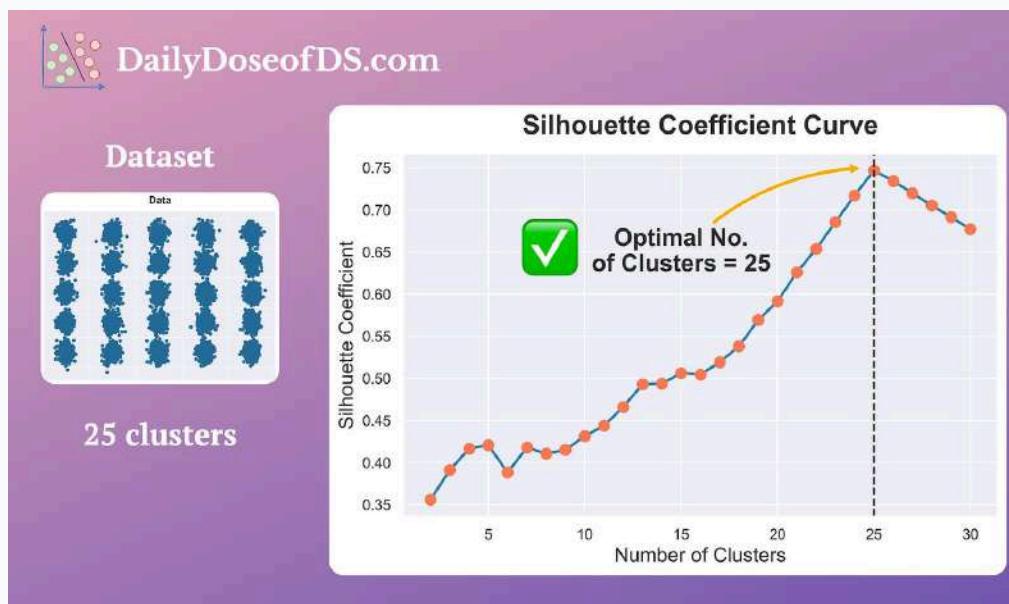
- For every data point:
 - find the average distance to all other points within its cluster → A
 - find the average distance to all points in the nearest cluster → B
 - point's score = $(B-A)/\max(B, A)$
- Compute the average of all individual scores to get the overall clustering score.

Now consider this:

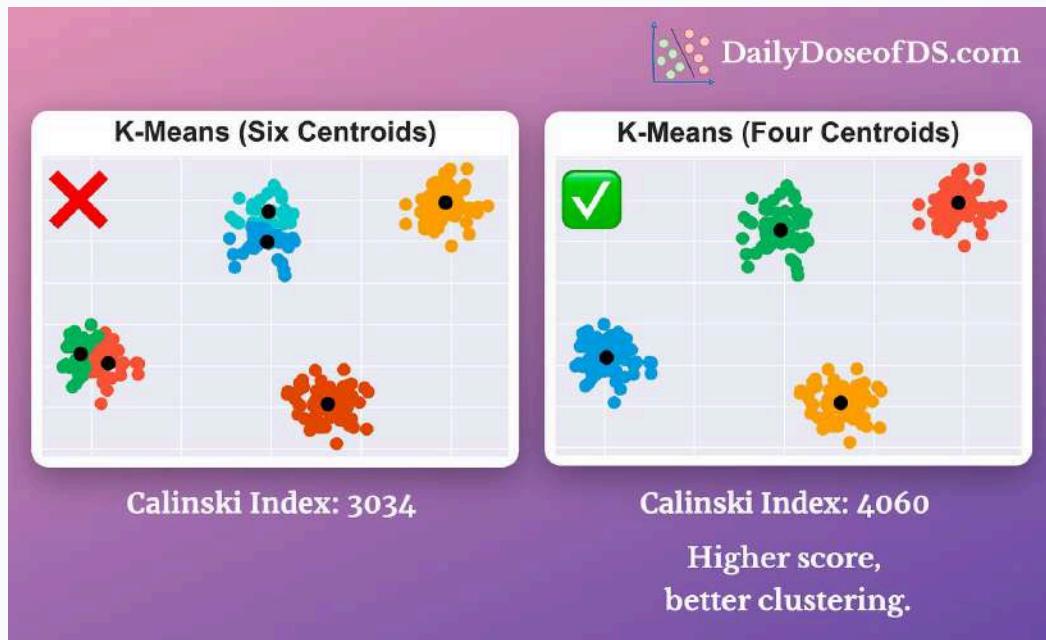
- If, for all data points, B is much greater than A, this would mean that:
 - The average distance to points in the closest cluster is large.
 - The average distance to points in the same cluster is small.
 - Thus, the overall score will be close to 1, indicating that the clusters are well separated.

So, a higher Silhouette Coefficient would mean that, generally speaking, all data points fit into their ideal clusters.

Measuring it across a range of centroids (k) can reveal which clustering results are most promising:



#2) Calinski-Harabasz Index



The issue with the Silhouette score is that its run-time grows quadratically with the number of data points.

Calinski-Harabasz Index is another metric that is measured quite similarly to the Silhouette Coefficient.

Here's what it computes:

- A → sum of squared distance between all centroids and overall dataset center.
- B → sum of squared distance between all points and their specific centroid.
- metric is computed as A/B (with an additional scaling factor).

Yet again, if A is much greater than B:

- The distance of centroids to the dataset center is large.
- The distance of data points to their specific centroid is small
- Thus, it will result in a higher score, indicating that the clusters are well separated.

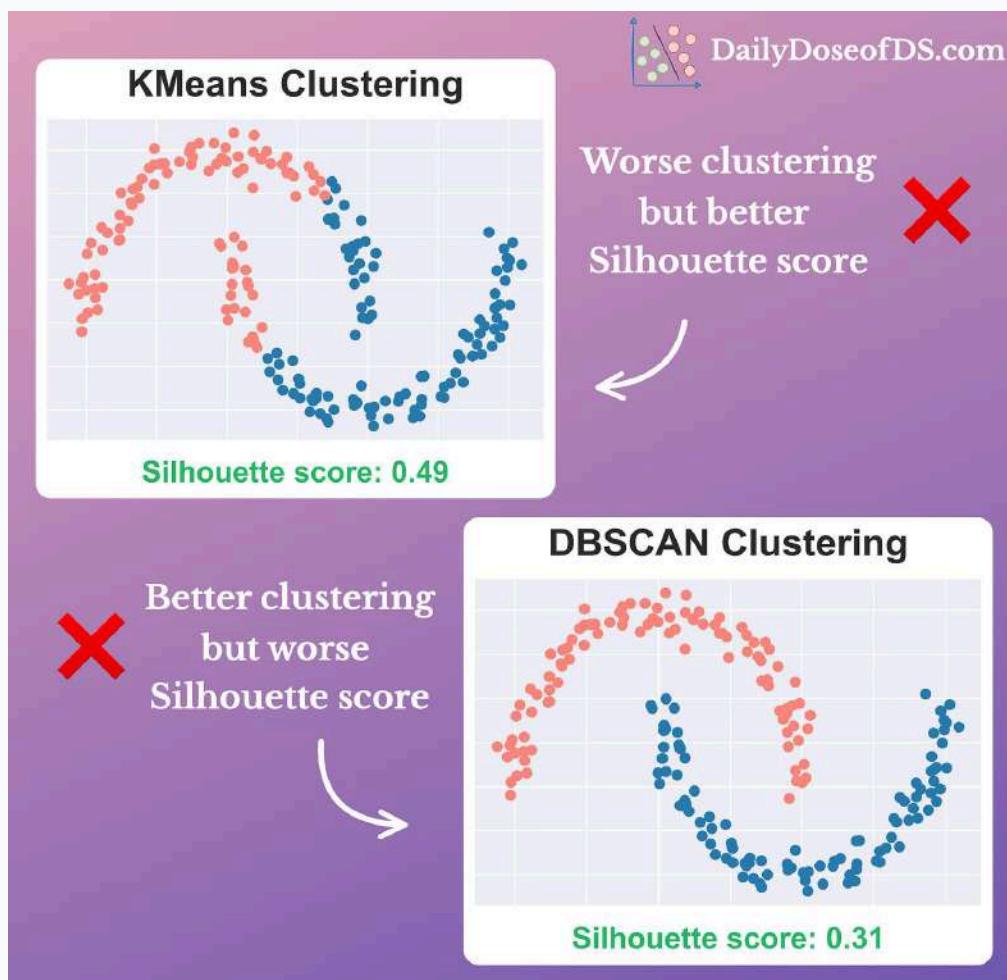
The reason I prefer the Calinski-Harabasz Index over the Silhouette Coefficient is that:

- It is relatively much faster to compute than the Silhouette Coefficient.
- It makes the same intuitive sense for interpretation as the Silhouette Coefficient.

#3) DBCV

The issue with the Silhouette score and Calinski-Harabasz index is that they are typically higher for convex (or somewhat spherical) clusters.

However, using them to evaluate arbitrary-shaped clustering, like those obtained with density-based clustering, can produce misleading results. This is evident from the image below:



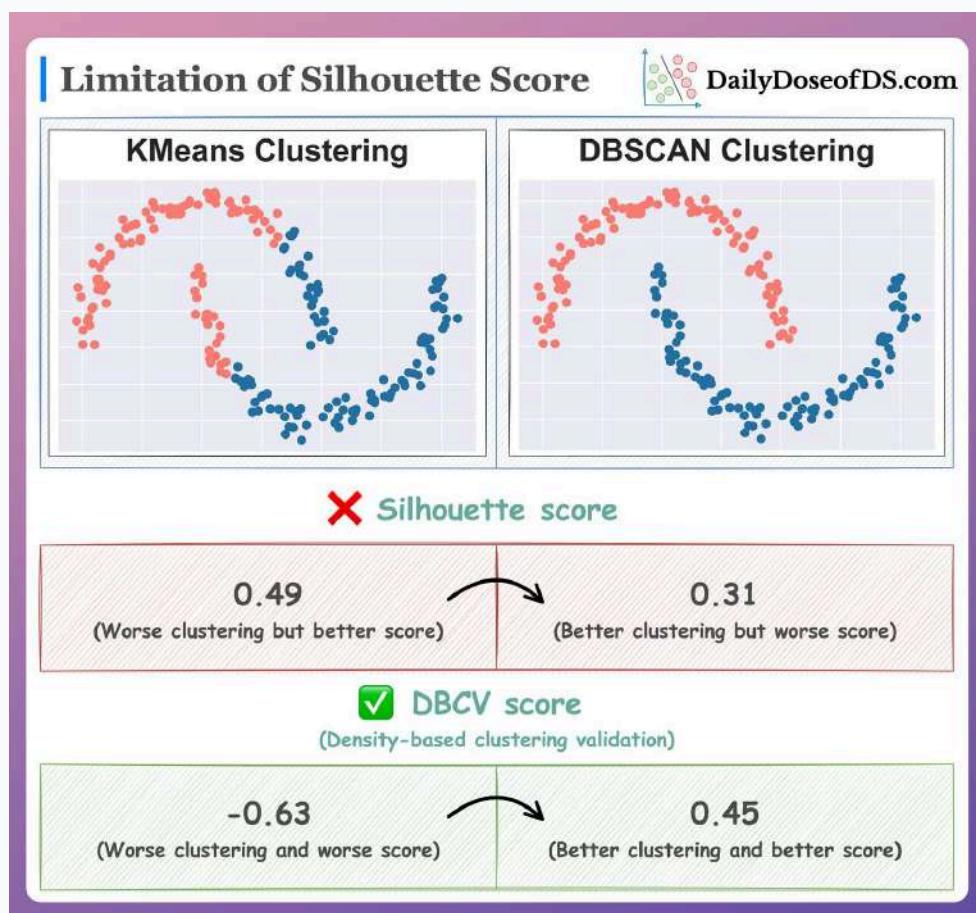
As depicted above, while the clustering output of KMeans is worse, the Silhouette score is still higher than Density-based clustering.

DBCV — density-based clustering validation is a better metric in such cases.

As the name suggests, it is specifically meant to evaluate density-based clustering. Simply put, DBCV computes two values:

- The density within a cluster
- The density overlap between clusters

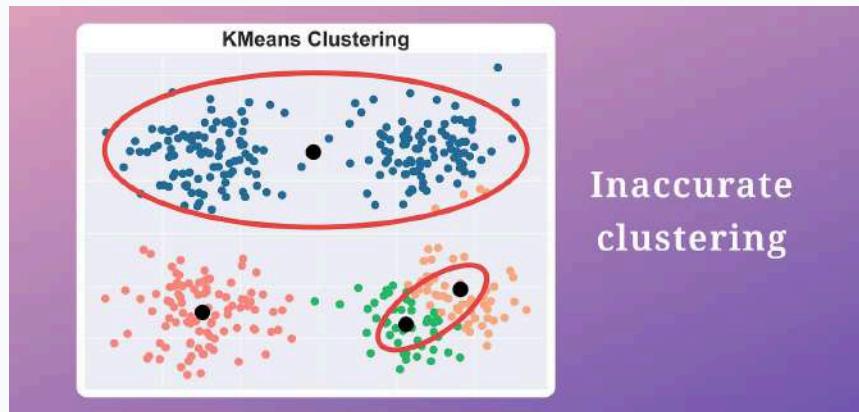
A high density within a cluster and a low density overlap between clusters indicate good clustering results. The effectiveness of DBCV is also evident from the image below:



With DBCV, the score for the clustering output of KMeans is worse, and that of density-based clustering is higher.

Breathing KMeans vs KMeans

KMeans is widely used for its simplicity and effectiveness as a clustering algorithm. Yet, we all know that its performance is entirely dependent on the centroid initialization step. Thus, it is likely that we may obtain inaccurate clusters, as depicted below:

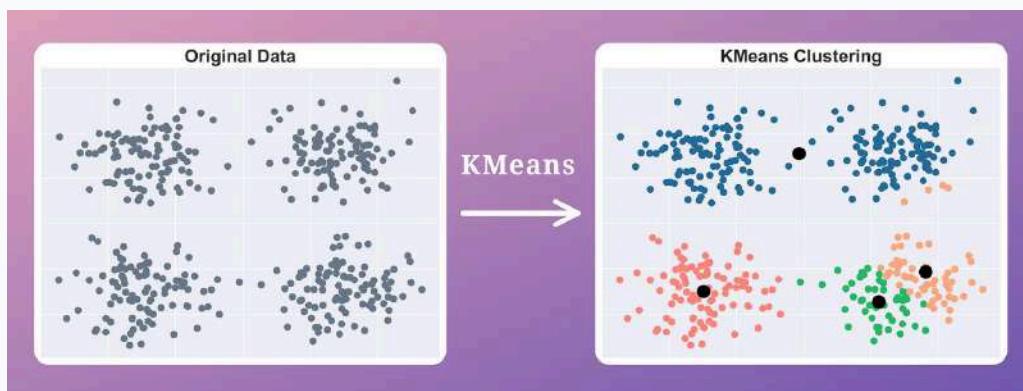


Of course, rerunning with different initialization does help at times. But I have never liked the unnecessary run-time overhead it introduces.

In this chapter, let's learn a supercharged upgrade to KMeans, which addresses this issue while also producing better clustering results. It's called the **Breathing KMeans** algorithm.

Step 1: Run Kmeans

First, we run the usual KMeans clustering only once, i.e., without rerunning the algorithm with a different initialization.



Clustering results of KMeans after running the algorithm without repetition and until convergence

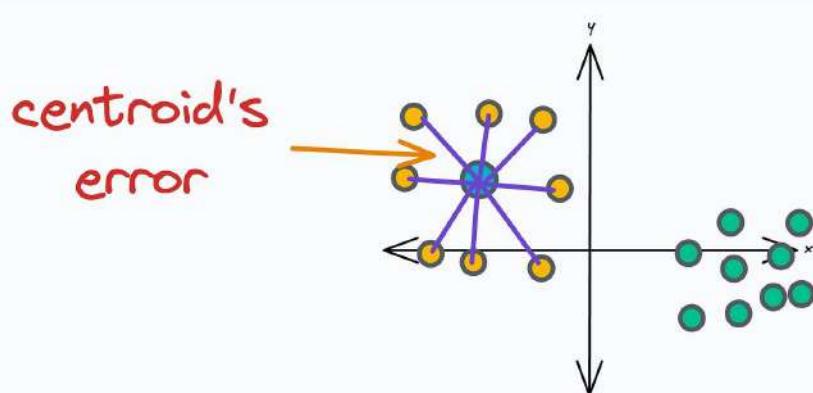
This gives us the location of “k” centroids, which may or may not be accurate.

Step 2: Breathe in step

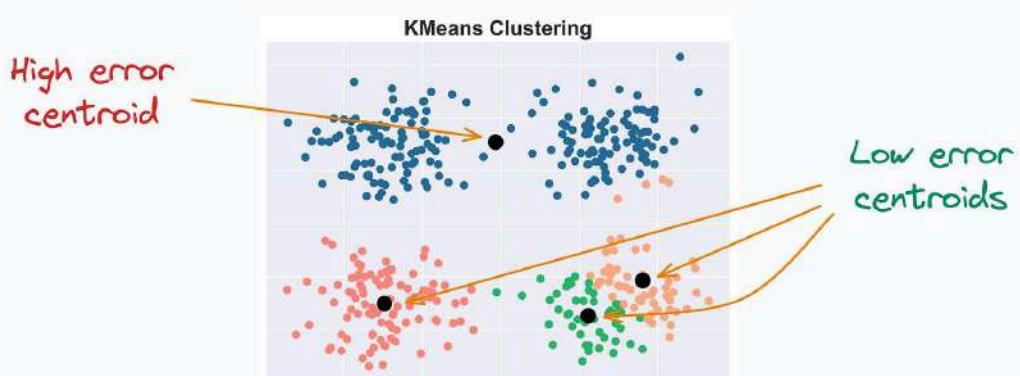
To the “k” centroids obtained from Step 1, we add “m” new centroids.

As per the research paper of Breathing Kmeans, m=5 was found to be good value. Now, you might be thinking, where do we add these “m” centroids?

The addition of new centroids is decided based on the error associated with a centroid. Simply put, a centroid’s error is the sum of the squared distance to the points associated with that centroid.



Thus, we add “m” centroids in the vicinity of centroids with high error. Let’s understand more intuitively why this makes sense.



In the above KMeans clustering results:

- The centroid at the top has a high error.
- All other centroids have relatively low error.

Intuitively speaking, if a centroid has a very high error, it is possible that multiple clusters are associated with it. Thus, we would want to split this cluster. Adding new centroids near clusters with high error will precisely fulfill this objective.

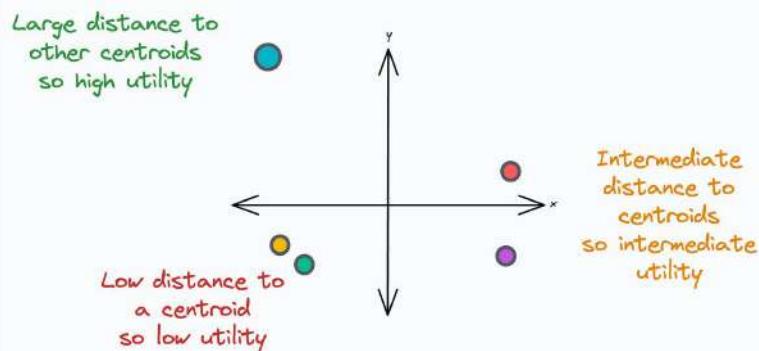
After adding “m” new centroids, we get a total of “k+m” centroids.

Finally, we run KMeans again with “k+m” centroids only once. This gives us the location of “k+m” centroids.

Step 3: Breathe out step

Next, we want to remove “m” centroids from the “k+m” centroids obtained above. Here, you might be thinking, which “m” centroids should we remove? The removal of centroids is decided based on the “utility” of a centroid.

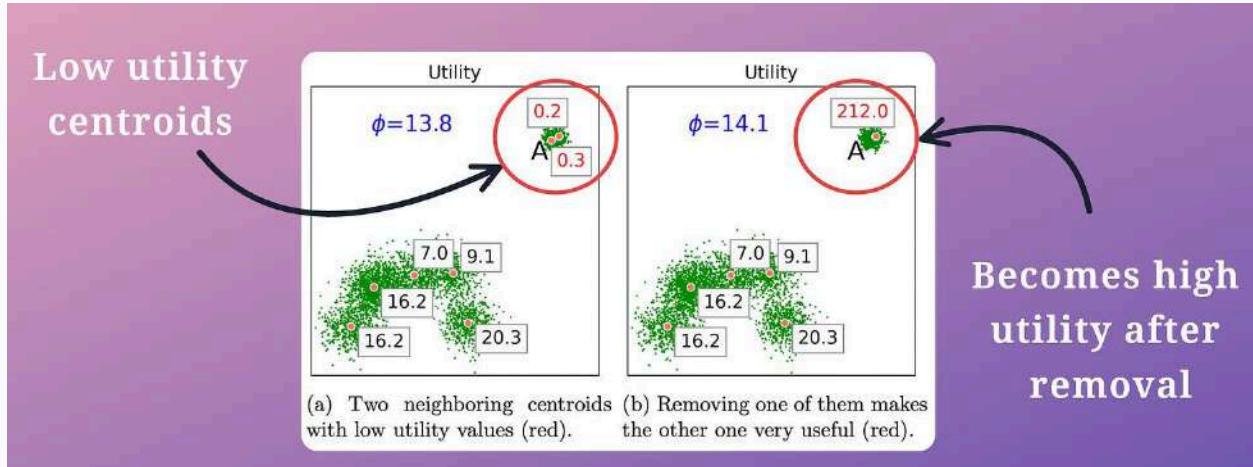
Simply put, a centroid’s utility is proportional to its distance from all other centroids. The greater the distance, the more isolated it will be. Hence, the more the utility.



This makes intuitive sense as well. If two centroids are pretty close, they will have low utility.

Thus, they are likely in the same cluster, and we must remove one of them.

This is demonstrated below:



After removing one of the low-utility centroids, the other centroid becomes very useful. So, in practice, after removing one centroid, we update the utility values of all other centroids.

We repeat the process until all “m” low-utility centroids have been removed. This gives back “k” centroids.

Finally, we run KMeans with these “k” centroids only once.

Step 4: Decrease m by 1.

Step 5: Repeat Steps 2 to 4 until m=0.

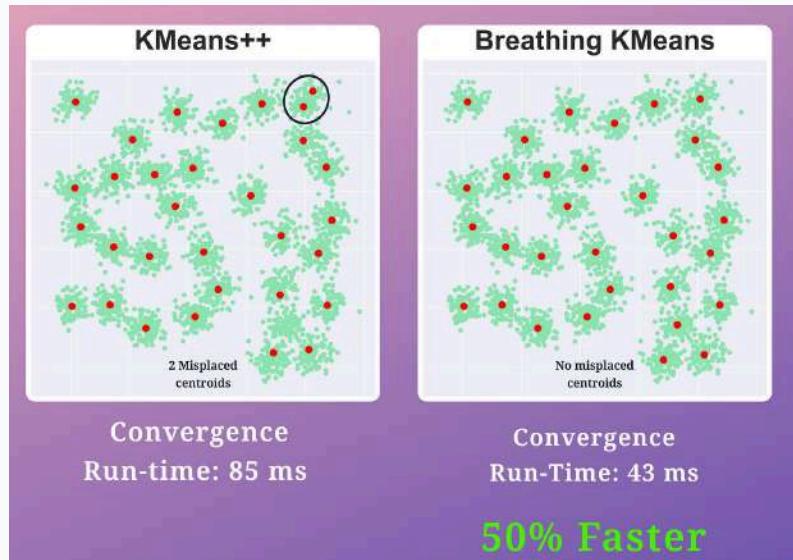
Done!

These repeated breathing cycles (breathe-in and breathe-out steps) almost always provide a faster and better solution than standard KMeans with repetitions.

In each cycle:

- New centroids are added at “good” locations. This helps in splitting clusters occupied by a single centroid.
- Low-utility centroids are removed. This helps in eliminating centroids that are likely in the same cluster.

As a result, it is expected to converge to the optimal solution faster. The effectiveness of Breathing KMeans over KMeans is evident from the image below:



- KMeans produced two misplaced centroids
- Breathing KMeans accurately clustered the data with a 50% run-time improvement.

There is also an open-source implementation of Breathing KMeans, with a sklearn-like API. To get started, install the `bkmeans` library:

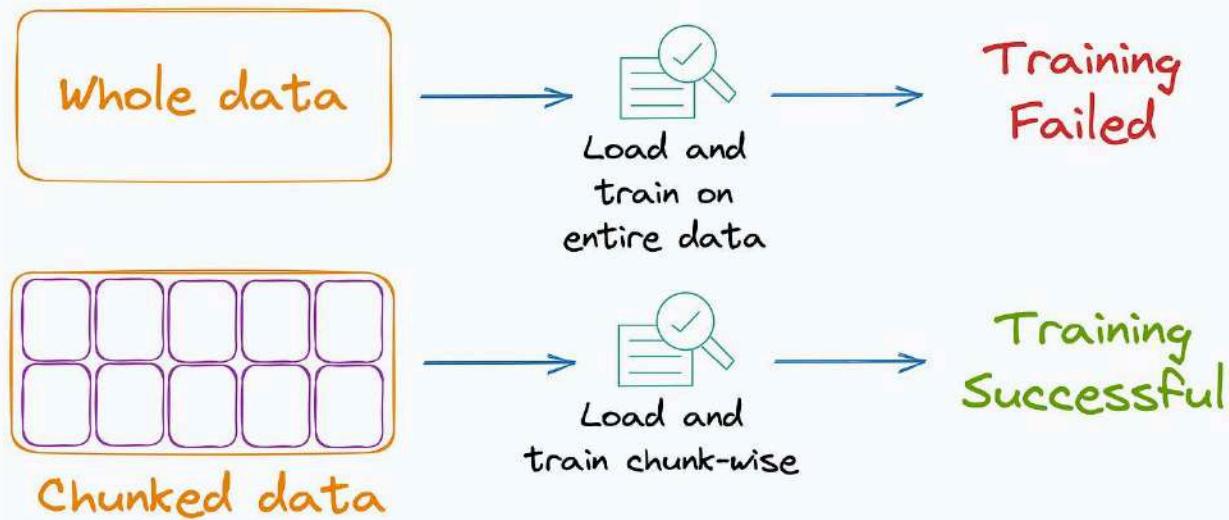
```
Python
pip install bkmeans
```

Next, run the algorithm as follows:

```
Python
from bkmeans import BKMeans
bkm = BKMeans(n_clusters = 100)
bkm.fit(X)
```

How Does MiniBatchKMeans Works?

Many ML algorithms implemented by Sklearn support incremental learning, i.e., mini-batch-based training. This makes them quite useful when the entire data does not fit into memory.



These are the supported algorithms for your reference:

<ul style="list-style-type: none"> • Classification <ul style="list-style-type: none"> ◦ <code>sklearn.naive_bayes.MultinomialNB</code> ◦ <code>sklearn.naive_bayes.BernoulliNB</code> ◦ <code>sklearn.linear_model.Perceptron</code> ◦ <code>sklearn.linear_model.SGDClassifier</code> ◦ <code>sklearn.linear_model.PassiveAggressiveClassifier</code> 	<ul style="list-style-type: none"> • Clustering <ul style="list-style-type: none"> ◦ <code>sklearn.cluster.MiniBatchKMeans</code>
<ul style="list-style-type: none"> • Regression <ul style="list-style-type: none"> ◦ <code>sklearn.linear_model.SGDRegressor</code> ◦ <code>sklearn.linear_model.PassiveAggressiveRegressor</code> 	<ul style="list-style-type: none"> • Decomposition / feature Extraction <ul style="list-style-type: none"> ◦ <code>sklearn.decomposition.MiniBatchDictionaryLearning</code> ◦ <code>sklearn.cluster.MiniBatchKMeans</code>

If you look closely, there's KMeans in that list too. Being a popular clustering, it's good that we can run it on large datasets if needed.

But when I first learned that KMeans can be implemented in a mini-batch fashion, I wondered how exactly a batch implementation of KMeans might work?

Giving it a thought for a minute or so gave me an idea, which I wish to share with you in this chapter.

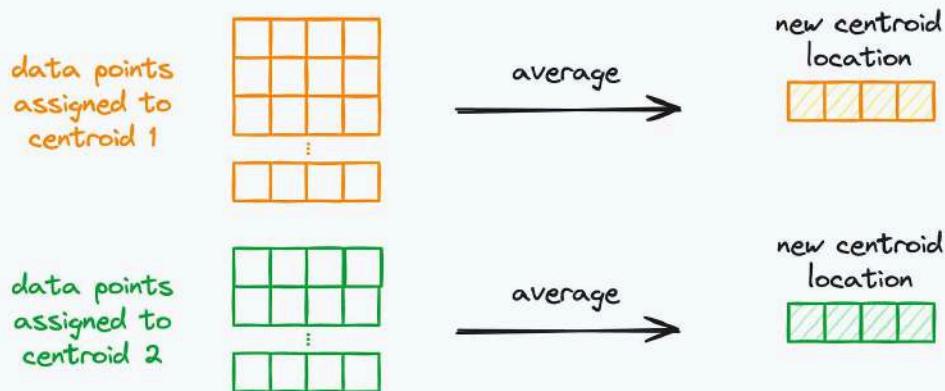
The source of overhead in KMeans

Before understanding the mini-batch implementation, we must know the issue with the usual KMeans implementation.

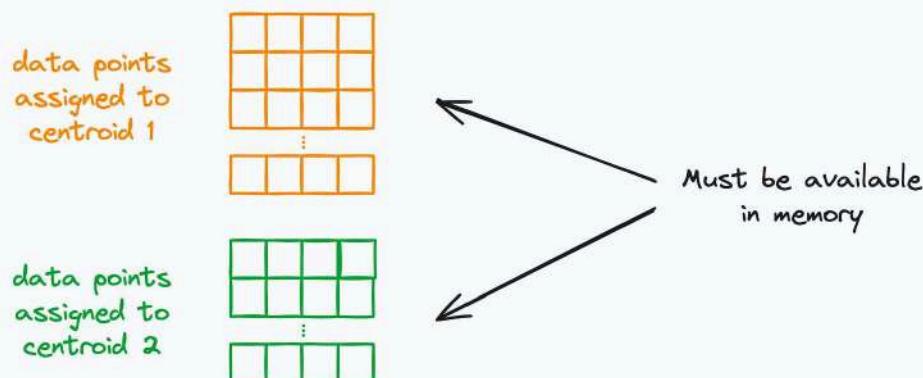
To recall, KMeans works as follows:

- Step 1) Initialize centroids.
- Step 2) Find the nearest centroid for each point.
- Step 3) Reassign centroids as the average of points assigned to them.
- Step 4) Repeat until convergence.

In this implementation, the primary memory bottleneck originates from Step 3. This is because, in this step, we compute the average of all points assigned to a centroid to shift the centroid.



Thus, traditionally, all data points assigned to a centroid must be available in memory to take an average over.



But in our case, all data points DO NOT fit into memory at once. Thus, direct averaging isn't feasible. Nonetheless, here's how we can smartly address this issue.

Before I proceed further:

Honestly speaking, I haven't checked the precise implementation of MiniBatchKMeans in Sklearn. The docs don't tell much either.

The solution I am about to share is the one I came up with when I was trying to figure out how a mini-batch KMeans can be possibly implemented.

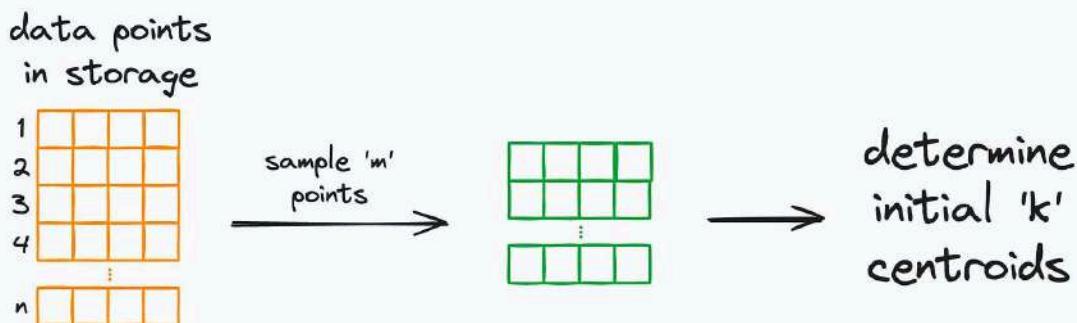
So if there's a perfect overlap with Sklearn, great!

If not, it's fine because the solution I am about to share is perfectly viable too.

MiniBatchKMeans

We begin the algorithm by selecting a value of k — the number of centroids.

Next, we need an initial location for the centroids so we can sample m ($m \ll n$) random data points to initialize the centroids.



Now, recall the bottleneck we discussed above — All data points must be available in memory to average over.

Thus, to address this, we need a way to keep track of the data points assigned to each centroid as we load and process each mini-batch.

We cannot maintain a growing list of data points assigned to each centroid because that would eventually take up the memory of the entire dataset.

So here's how we can smartly fulfill this objective.

We can utilize the following mathematical property that relates sum and average:

$$\text{average} = \frac{\text{sum}}{n}$$

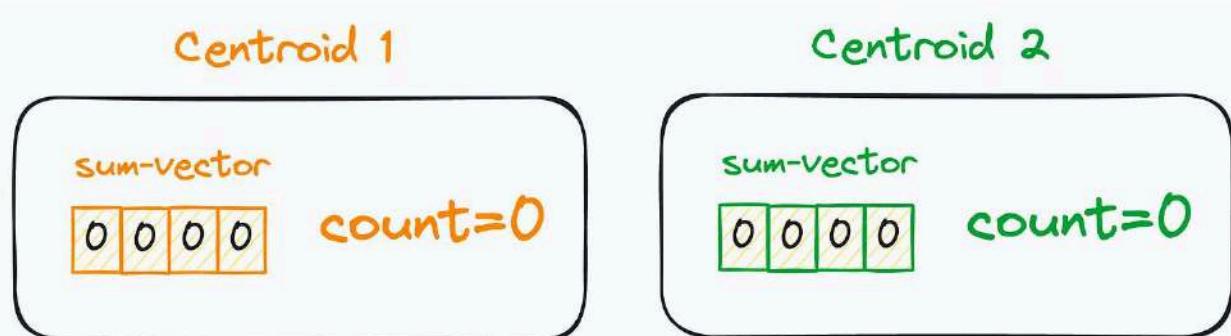
How?

See, in a typical KMeans setting, we average over all data points assigned to a centroid to determine its new location. But as discussed above, we can't keep all data points in memory at once.

So how about we do this in the mini-batch setting of KMeans:

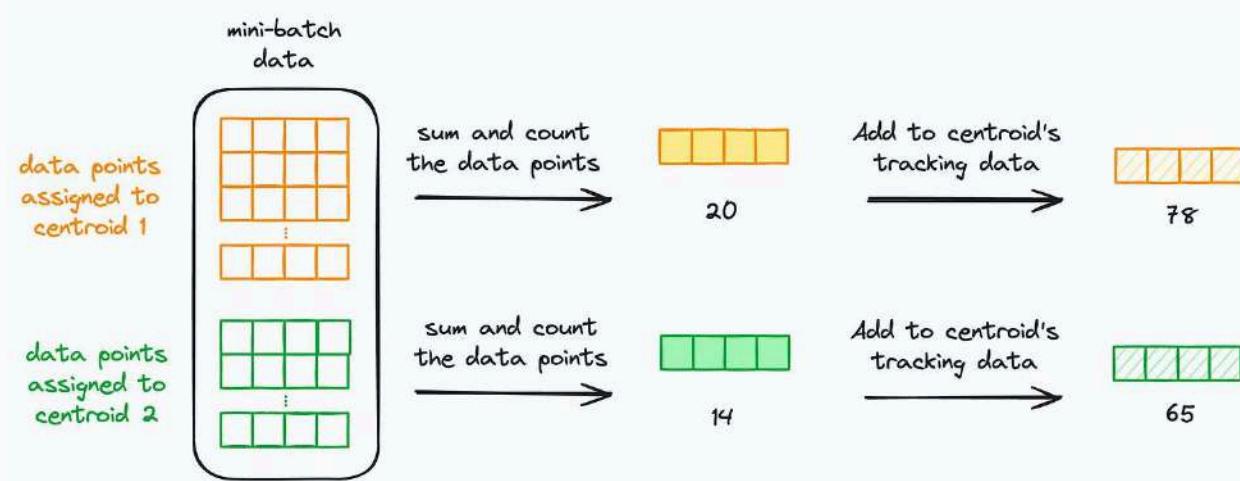
For every centroid:

- We maintain a “sum-vector”, which stores the sum of the vectors of all data points assigned to that centroid. To begin, this will be a zero vector.
- A “count” variable, which tracks the number of data points assigned to a centroid. To begin, the count will be zero.

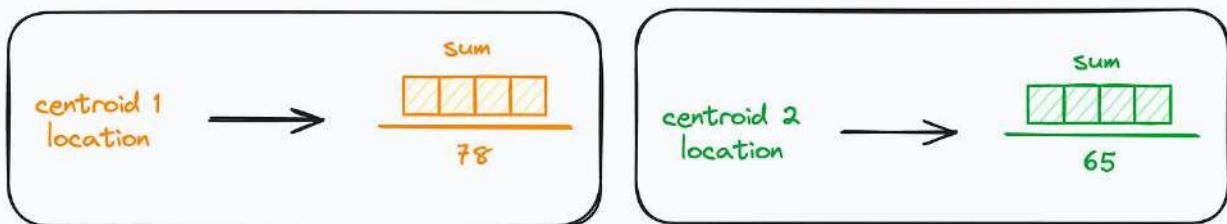


While iterating through mini-batches:

- Every “sum-vector” of a centroid will continue to accumulate the sum of all data points that have been assigned to it so far.
- Every “count” variable of a centroid will continue to track the number of data points that have been assigned to it so far.



After iterating through all mini-batches, we can perform a scalar division between sum-vector and count to get the average vector:



This average will precisely tell us the new location of the centroid as if all data points were present in memory. We can repeat this until convergence.

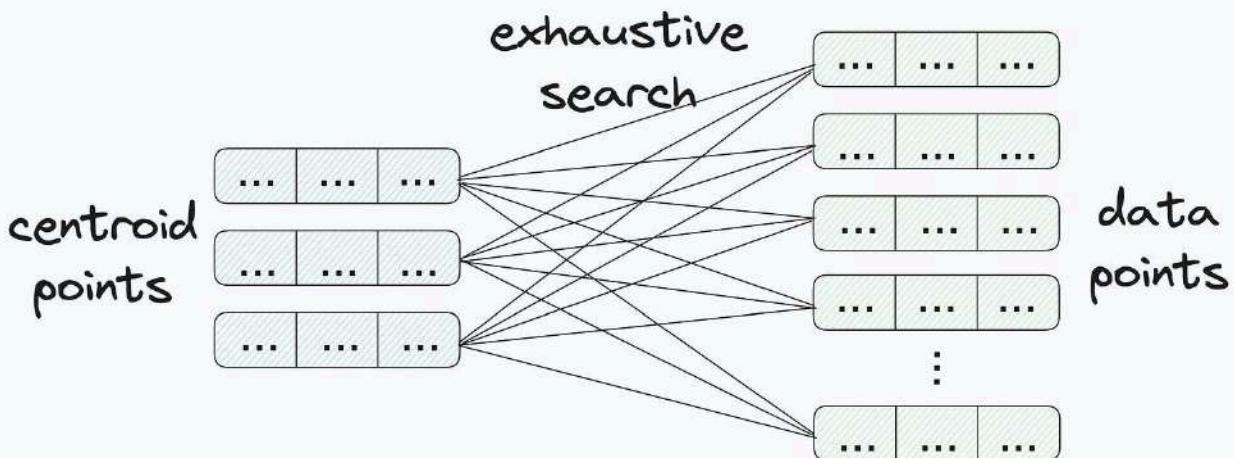
Also, before starting the next iteration of MiniBatchKMeans, we must reset the sum-vector to a zero vector and count to zero. And that's how a mini-batch version of KMeans can be implemented.

ANN-driven KMeans with Faiss

As discussed in the previous chapter, KMeans is trained as follows:

- Step 1) Initialize centroids
- Step 2) Find the nearest centroid for each point
- Step 3) Reassign centroids
- Step 4) Repeat until convergence

But in this implementation, “Step 2” has a run-time bottleneck, as this step involves a brute-force and exhaustive search. In other words, this finds the distance of every data point from every centroid.



As a result, this step isn't optimized, and it takes plenty of time to train and predict. This is especially challenging with large datasets.

To speed up KMeans, one of the implementations I usually prefer, especially on large datasets, is Faiss by Facebook AI Research.

To elaborate further, Faiss provides a much faster nearest-neighbor search using approximate nearest-neighbor search algorithms.

It uses an “Inverted Index,” which is an optimized data structure to store and index the data point.

This makes performing clustering extremely efficient, especially on large datasets, which is also evident from the image below:

DailyDoseofDS.com

sklearn.py

```
from sklearn.cluster import KMeans
kmeans = KMeans(8).fit(x_train)
# Training Time: 162s
```

x_train shape: (500000, 1024)

faiss.py

```
import faiss
kmeans = faiss.Kmeans(d=1024, k=8)
kmeans.train(x_train)
# Training Time: 7.8s
```

~20x Faster

As shown above, on a dataset of 500k data points (1024 dimensions), Faiss is roughly 20x faster than KMeans from Sklearn, which is an insane speedup.

What's more, Faiss can also run on a GPU, which can further speed up your clustering run-time performance.

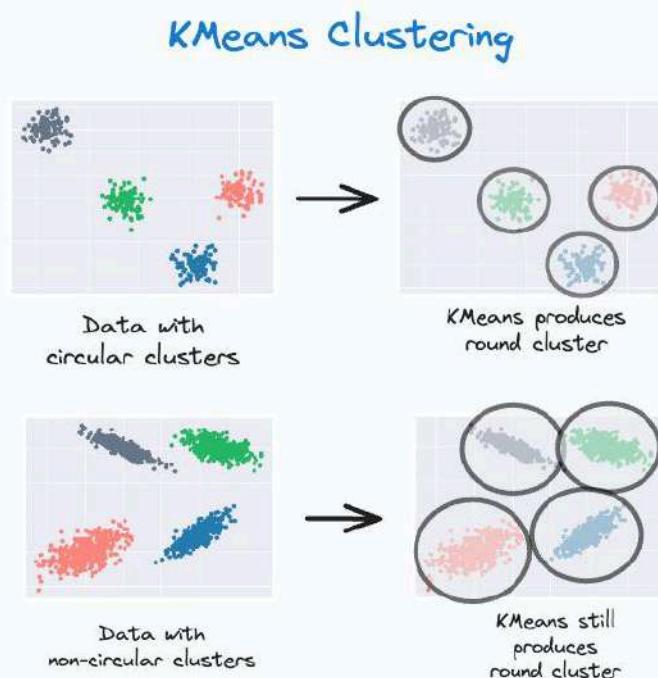
Get started with Faiss here: <https://github.com/facebookresearch/faiss>.

KMeans vs. Gaussian Mixture Models

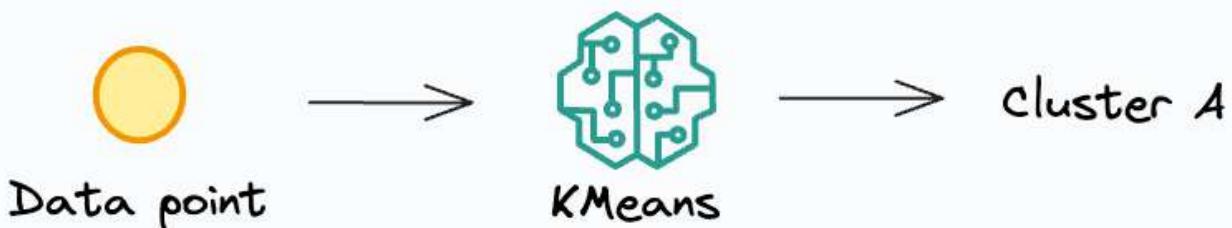
I like to think of Gaussian Mixture Models as a more generalized version of KMeans, which addresses some of the widely known limitations of KMeans.

Limitations of KMeans

To begin, KMeans can only produce globular clusters. For instance, as shown below, even if the data has non-circular clusters, it still produces round clusters.

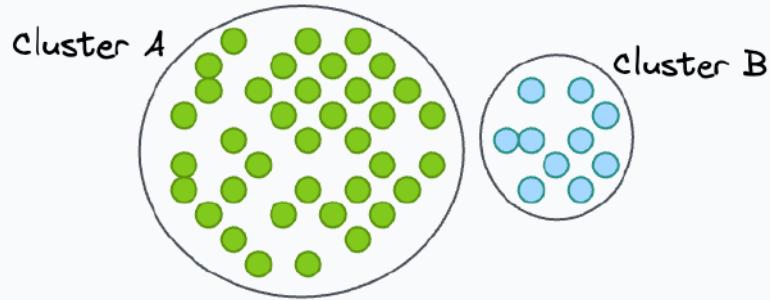


Moreover, it performs a hard assignment. There are no probabilistic estimates of each data point belonging to each cluster.

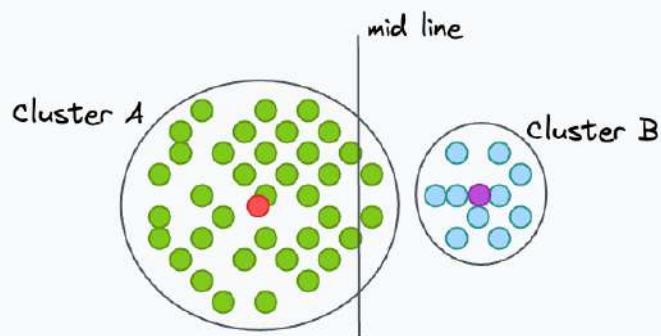


Lastly, it only relies on distance-based measures to assign data points to clusters.

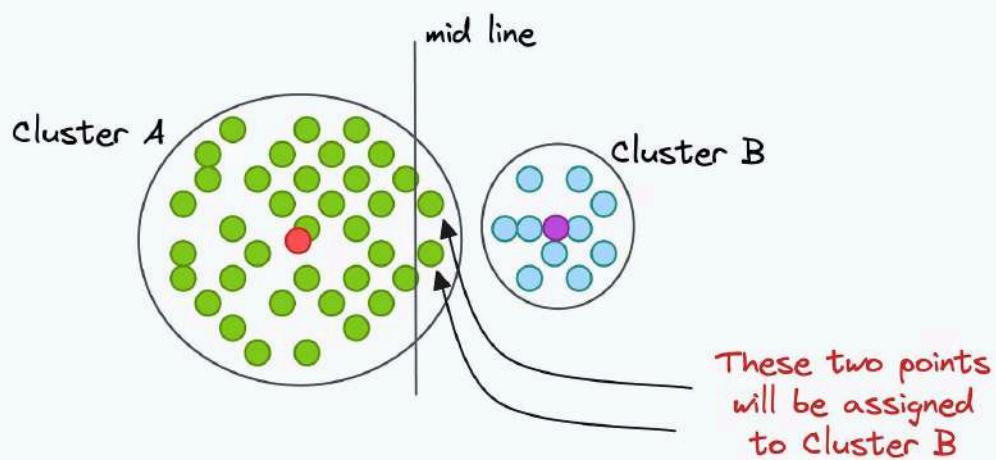
To understand this point better, consider two clusters in 2D — A and B. Cluster A has a higher spread than B.



Now consider a line that is mid-way between centroids of A and B.



Although A has a higher spread, even if a point is slightly right to the midline, it will get assigned to cluster B.

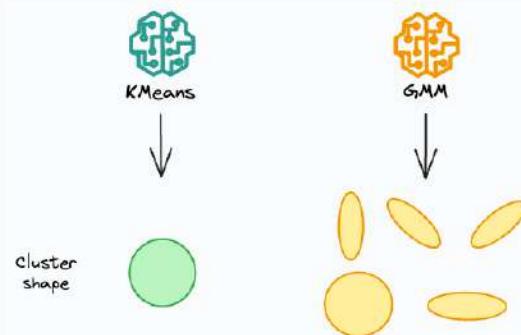


Ideally, however, cluster A should have had a larger area of influence.

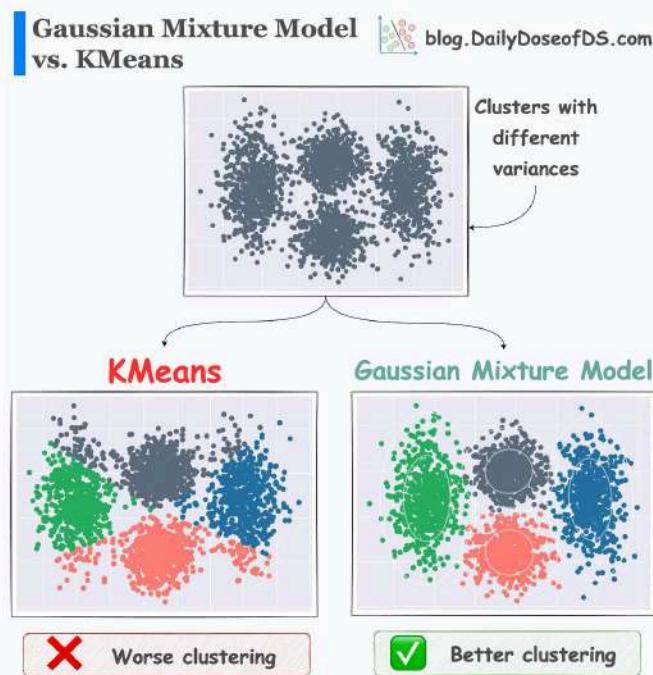
Gaussian Mixture Models

These limitations often make KMeans a non-ideal choice for clustering. Gaussian Mixture Models are often a superior algorithm in this respect. As the name suggests, they can cluster a dataset that has a mixture of many Gaussian distributions. They can be thought of as a more flexible twin of KMeans.

The primary difference is that, KMeans learns centroids but Gaussian mixture models learn a distribution. For instance, in 2 dimensions, KMeans can only create circular clusters but GMM can create oval-shaped clusters.



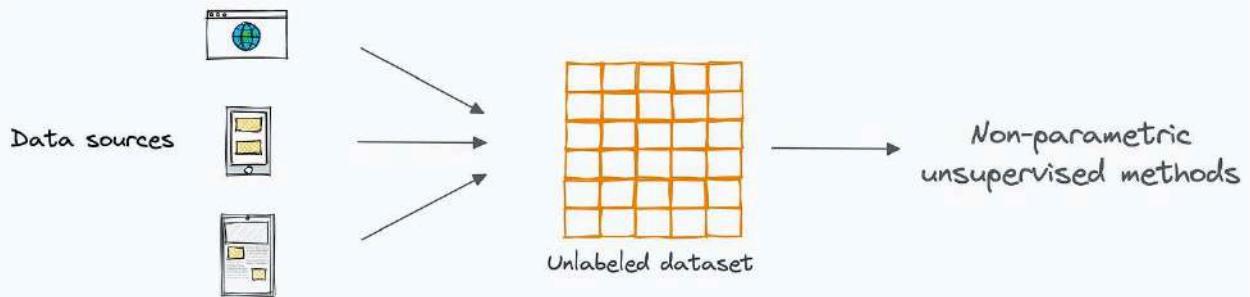
The effectiveness of GMMs over KMeans is evident from the image below.



- KMeans just relies on distance and ignores the distribution of each cluster.
- GMM learns the distribution and produces better clustering.

DBSCAN++: A Faster and Scalable Alternative to DBSCAN

Non-parametric unsupervised algorithms are widely used in industry to analyze large datasets. This is because in many practical applications, gathering true labels is pretty difficult or infeasible.



In such cases:

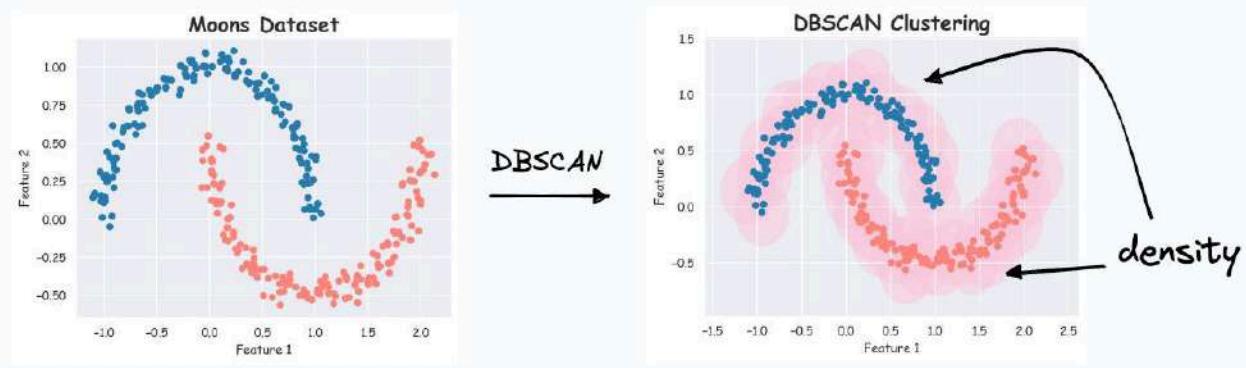
- Either data teams annotate the data, which can be practically impossible at times,
- Or use unsupervised methods to identify patterns.

While KMeans is widely used here due to its simplicity and effectiveness as a clustering algorithm, it has many limitations as we also discussed in the previous chapter:

- It does not account for cluster covariance.
- It can only produce spherical clusters. As shown below, even if the data has non-circular clusters, it still produces round clusters.

Density-based algorithms, like DBSCAN, quite effectively address these limitations.

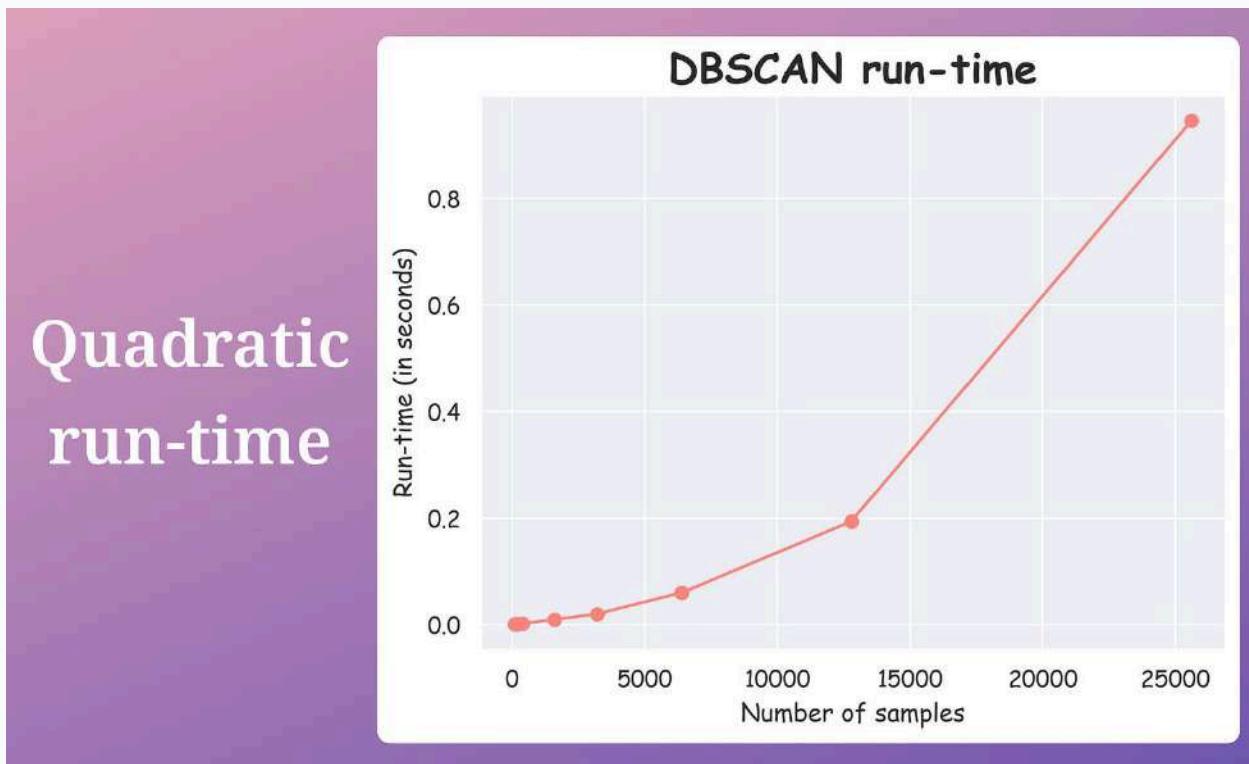
The core idea behind DBSCAN is to group together data points based on “density”, i.e., points that are close to each other in a high-density region and are separated by lower-density regions.



But unfortunately, with the use of DBSCAN, we again run into an overlooked problem. One of the things that makes DBSCAN infeasible to use at times is its run-time.

Until recently, it was believed that DBSCAN had a run-time of $O(n \log n)$, but it was proven to be $O(n^2)$ in the worst case.

In fact, this can also be verified from the figure below. It depicts the run-time of DBSCAN with the number of samples:



It is clear that DBSCAN has a quadratic relation with the dataset's size.

Thus, there is an increasing need to establish more efficient versions of DBSCAN.

DBSCAN++ is a major step towards a fast and scalable DBSCAN.

Simply put, DBSCAN++ is based on the observation that we only need to compute the density estimates for a subset "m" of the "n" data points in the given dataset, where "m" can be much smaller than "n" to cluster properly. The effectiveness of DBSCAN++ is evident from this image.

The image shows two side-by-side terminal windows comparing the performance of DBSCAN and DBSCAN++ on a dataset of 60k data points.

Top Window (DBSCAN):

- Code:

```
from sklearn.cluster import DBSCAN
```
- Execution:

```
dbscan = DBSCAN(eps=0.3, min_samples=5)
```
- Output:

```
>>> clusters = dbscan.fit(X)
# Run-time: ~150 seconds
```
- Output:

```
>>> score(clusters, y)
# Accuracy: 0.56
```

Bottom Window (DBSCAN++):

- Code:

```
from DBSCANpp import DBSCANpp
```
- Execution:

```
dbscanpp = DBSCANpp(eps=0.3, min_samples=5, sampling="kcenter")
```
- Output:

```
>>> clusters = dbscanpp.fit(X)
# Run-time: ~7 seconds
```
- Output:

```
>>> score(clusters, y)
# Accuracy: 0.56
```

Annotations:

- A red checkmark icon is placed next to the DBSCAN++ code.
- A green checkmark icon is placed next to the DBSCAN++ run-time output.
- A red X icon is placed next to the DBSCAN run-time output.
- A callout box on the right says "Dataset: 60k data points" with a grid icon.
- A large red X is overlaid on the DBSCAN run-time output.
- The text "DBSCAN is slow" is overlaid on the DBSCAN run-time output.
- The text "DBSCAN++ is 20x Faster and has same score" is overlaid on the DBSCAN++ run-time output.

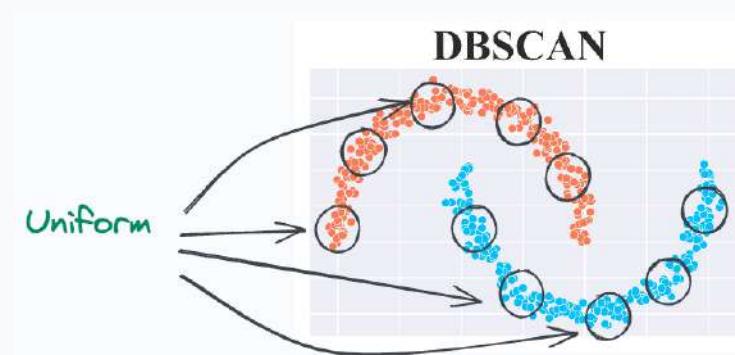
As depicted above, on a dataset of 60k data points:

- DBSCAN++ is 20x faster than DBSCAN.
- DBSCAN++ produces better clustering scores than DBSCAN.

HDBSCAN vs. DBSCAN

In this chapter, let's dive into HDBSCAN and understand how it differs from DBSCAN.

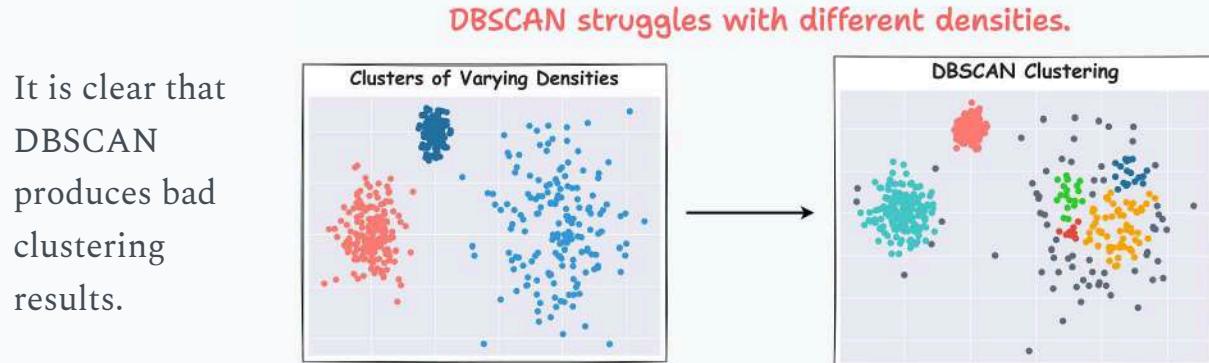
Like any other algorithm, DBSCAN also has some limitations. To begin, DBSCAN assumes that the local density of data points is (somewhat) globally uniform. This is governed by its `eps` parameter.



Thus, it may struggle to identify clusters with varying densities. This may need several hyperparameter tuning attempts to get promising results.

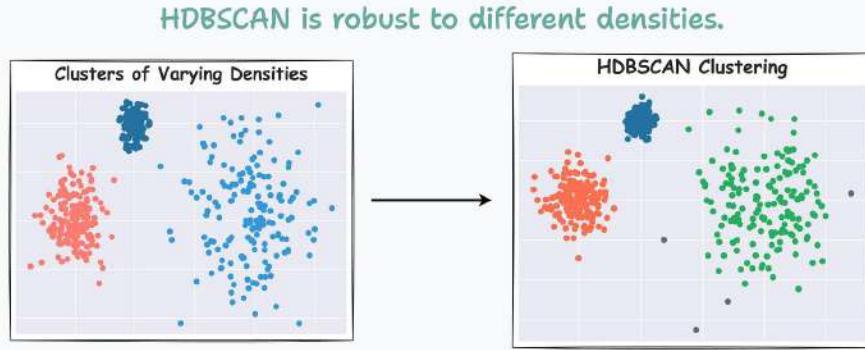
HDBSCAN can be a better choice for density-based clustering. It relaxes the assumption of local uniform density, which makes it more robust to clusters of varying densities by exploring many different density scales.

For instance, consider the clustering results obtained with DBSCAN on the dummy dataset below, where each cluster has different densities:



Now compare it with HDBSCAN results depicted below:

On a dataset with three clusters, each with varying densities, HDBSCAN is found to be more robust.

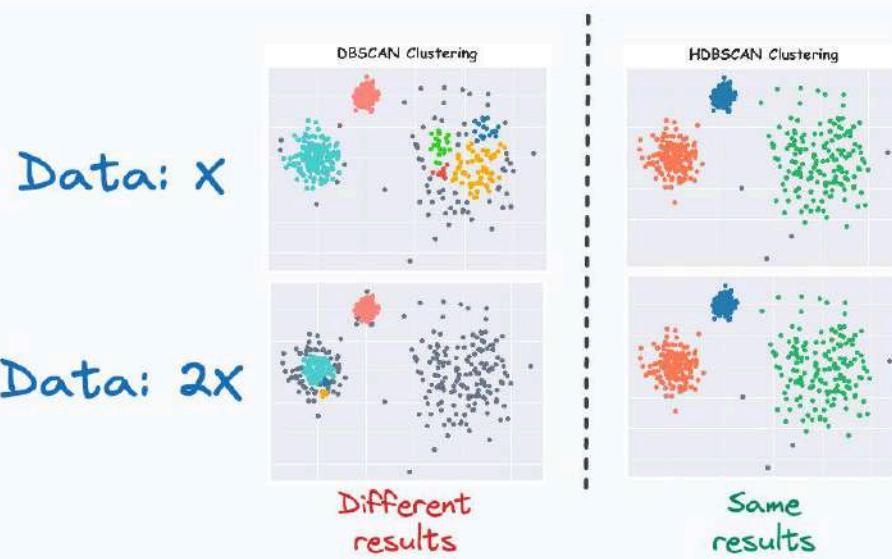


There's one more thing I love about HDBSCAN:

- DBSCAN is a scale variant algorithm. Thus, clustering results for data X, 2X, 3X, etc., can be entirely different.
- On the other hand, HDBSCAN is scale-invariant. So, clustering results remain the same across different scales of data.

We have DBSCAN on the left, and we can see that the results vary with the scale of the data.

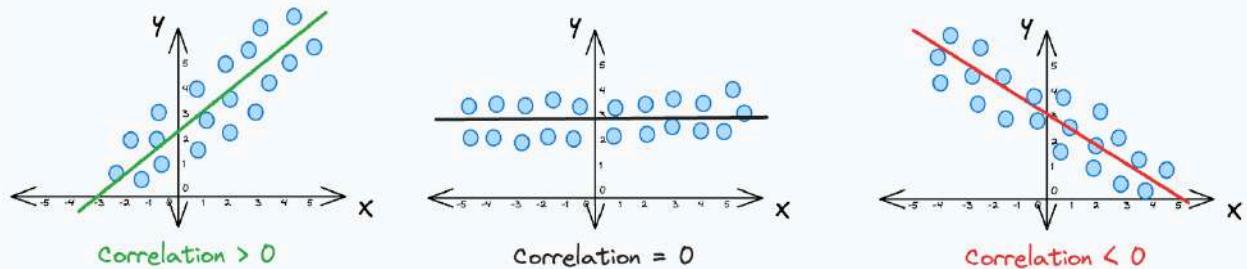
However, clustering from HDBSCAN (on the right) remains unaltered with the scale for HDBSCAN.



Correlation analysis

Correlation != Predictiveness

Correlation measures how two features vary with one another linearly (or monotonically).



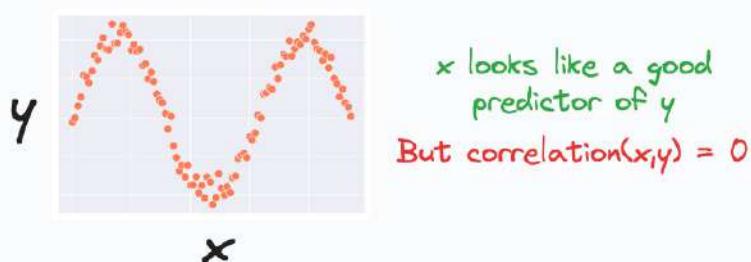
This makes correlation symmetric: $\text{corr}(A, B) = \text{corr}(B, A)$.

Yet, associations are often asymmetric. For instance, given a date, it is easy to tell the month. But given a month, you can never tell the date.

28/08/2023 → August

August → ?

Correlation, being symmetric, entirely ignores this notion. What's more, it is not meant to quantify how well a feature can predict the outcome, as demonstrated below:

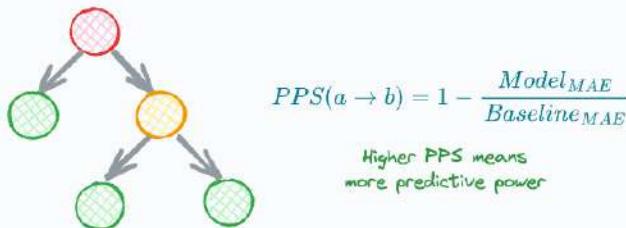


Yet, at times, it is misinterpreted as a measure of “predictiveness”. Lastly, correlation is mostly limited to numerical data. But categorical data is equally important for predictive models.

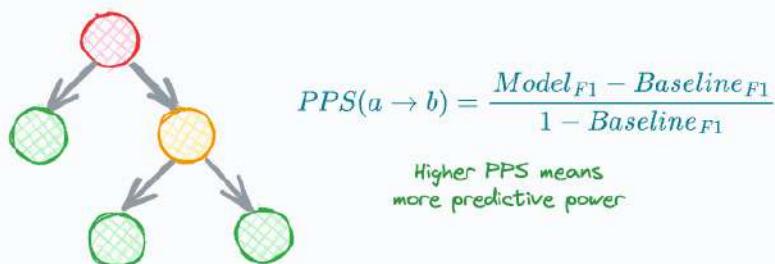
The Predictive Power Score (PPS) addresses each of these limitations. As the name suggests, it measures the predictive power of a feature.

$\text{PPS}(a \rightarrow b)$ is calculated as follows:

- If the target (b) is numeric:
 - Train a Decision Tree Regressor that predicts b using a .



- Find PPS by comparing its MAE to the MAE of a baseline model (median prediction).
- If the target (b) is categorical:
 - Train a Decision Tree Classifier that predicts b using a .



- Find PPS by comparing its F1 to the F1 of a baseline model (random or most frequent prediction).

Thus, PPS:

- is asymmetric, meaning $\text{PPS}(a, b) \neq \text{PPS}(b, a)$.
- can be used on categorical targets (b).

- can be used to measure the predictive power of categorical features (a).
- works well for linear and non-linear relationships.
- works well for monotonic and non-monotonic relationships.

Its effectiveness is evident from the image below. For all three datasets:

Dataset	Pearson Correlation	Predictive Power Score (PPS)	
		$x \rightarrow y$	$y \rightarrow x$
	0.01 ✗	0.81 ✓	0.00 ✓
	0.03 ✗	0.87 ✓	0.00 ✓
	0.01 ✗	0.83 ✓	0.00 ✓

Low Correlation High PPS for $x \rightarrow y$ Low PPS for $y \rightarrow x$

- Correlation is low.
- PPS ($x \rightarrow y$) is high.
- PPS ($y \rightarrow x$) is zero.

That being said, it is important to note that correlation has its place.

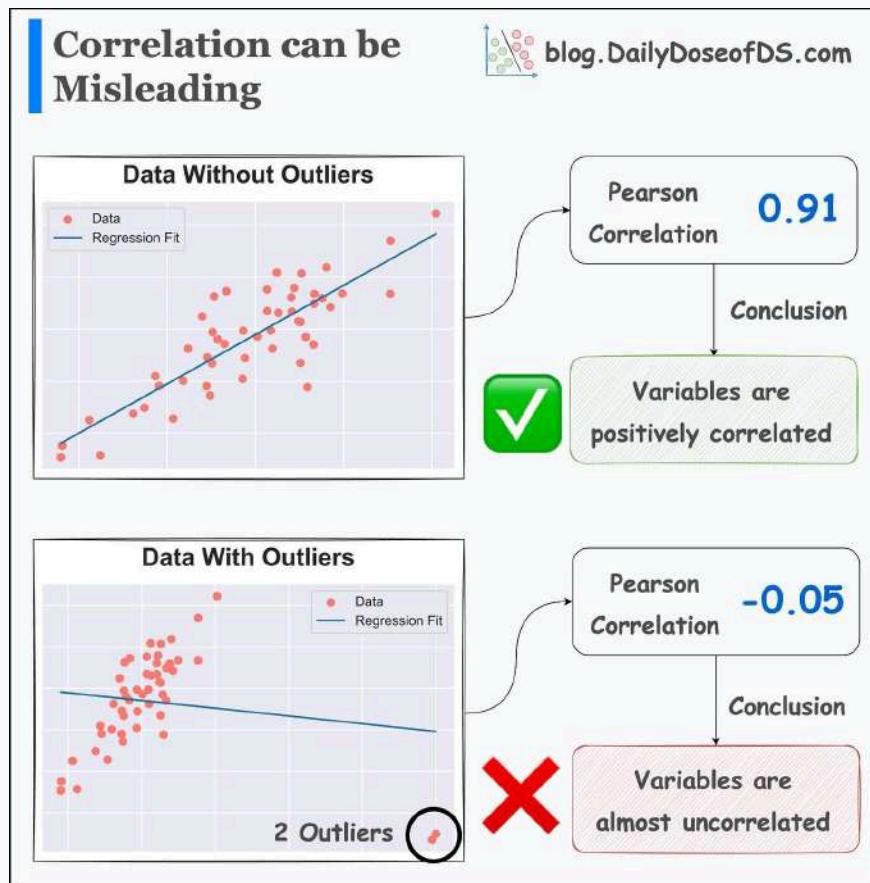
When selecting between PPS and correlation, first set a clear objective about what you wish to learn about the data:

- Do you want to know the general monotonic trend between two variables? Correlation will help.
- Do you want to know the predictiveness of a feature? PPS will help.

Get started with PPS: <https://github.com/8080labs/ppscode>.

Beware of Summary Statistics

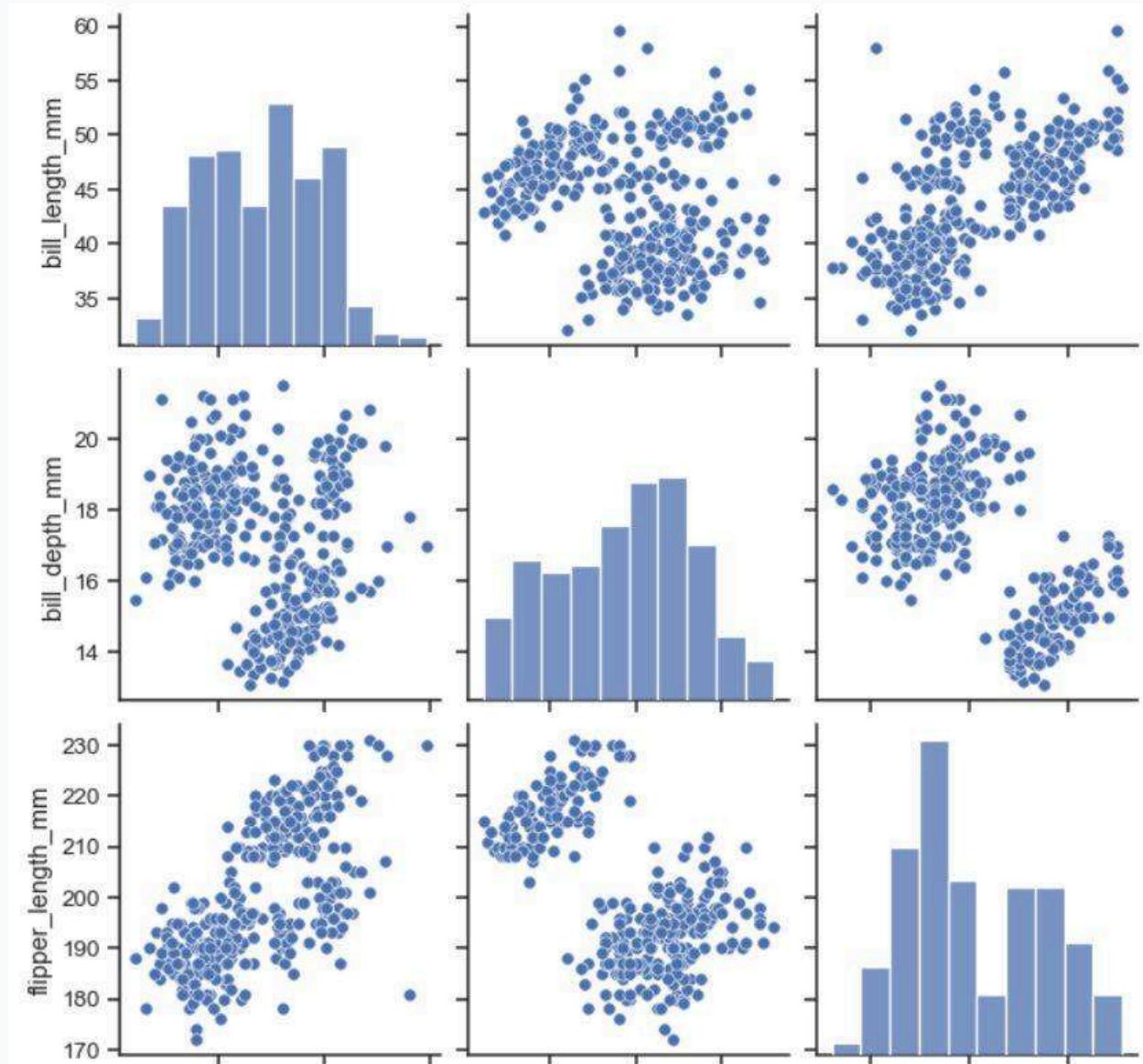
Many data scientists solely rely on the correlation matrix to study the association between variables. But unknown to them, the obtained statistic can be heavily driven by outliers. This is evident from the image below.



Adding just two outliers drastically changed the correlation coefficient and the regression fit. Thus, plotting the data is highly important.

This can save you from drawing wrong conclusions, which you may have drawn otherwise by solely looking at the summary statistics.

One thing that I often do when using a correlation matrix is creating a PairPlot as well (shown below).



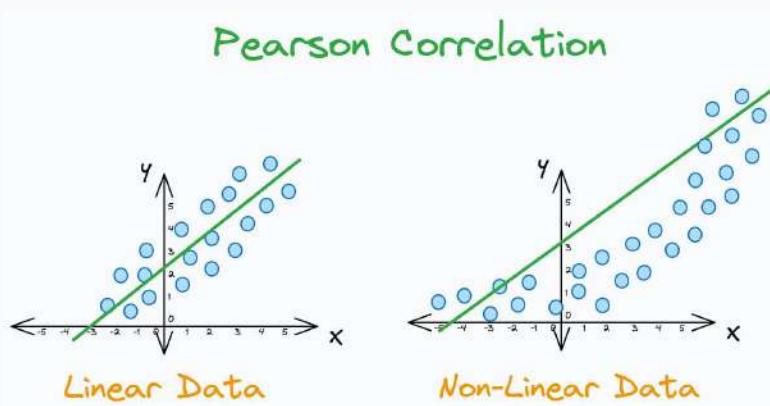
This lets me infer if the scatter plot of two variables and their corresponding correlation measure resonate with each other or not. We shall continue this discussion in the next chapter.

Pearson correlation can only measure linear association

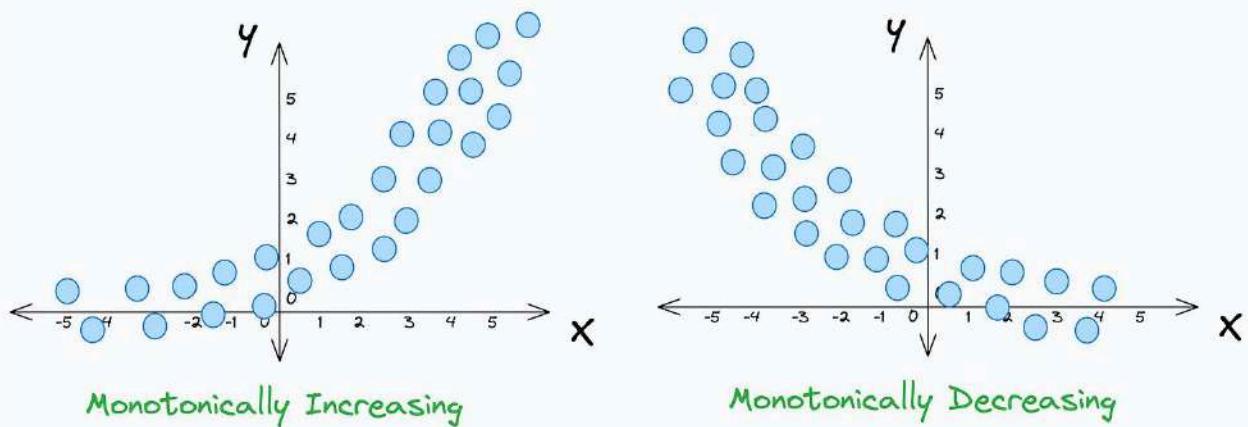
Pearson correlation is commonly used to determine the association between two continuous variables. Many frameworks (like Pandas, for instance) have it as their default correlation metric.

Yet, unknown to many, Pearson correlation:

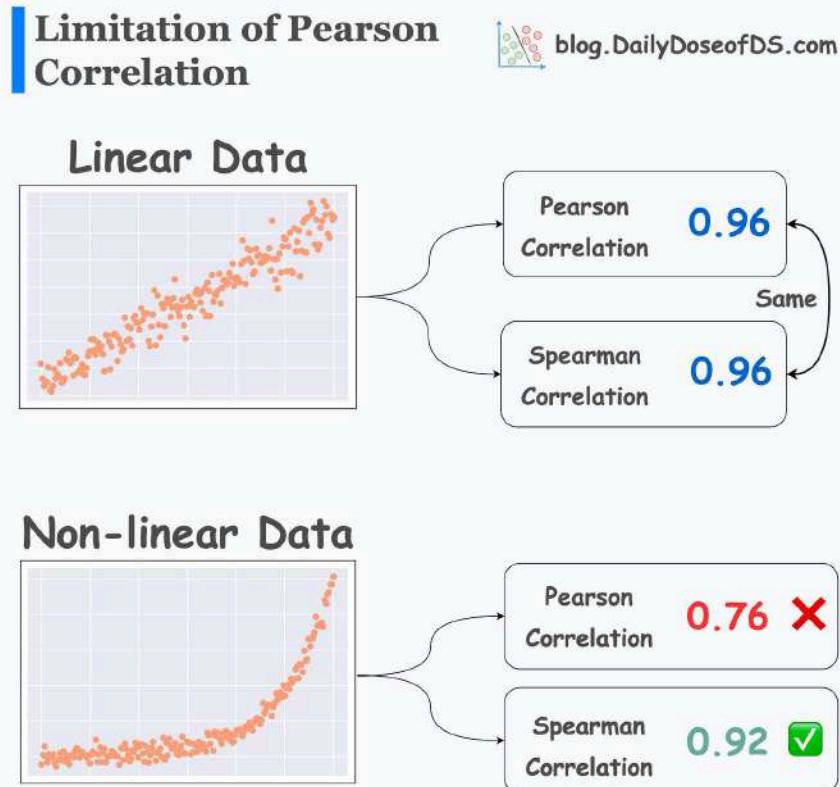
- Only measures the linear relationship.
- Penalizes a non-linear yet monotonic association.



Instead, Spearman correlation is a better alternative. It assesses monotonicity, which can be linear as well as non-linear.



This is evident from the illustration below:



Pearson vs. Spearman on linear and non-linear data

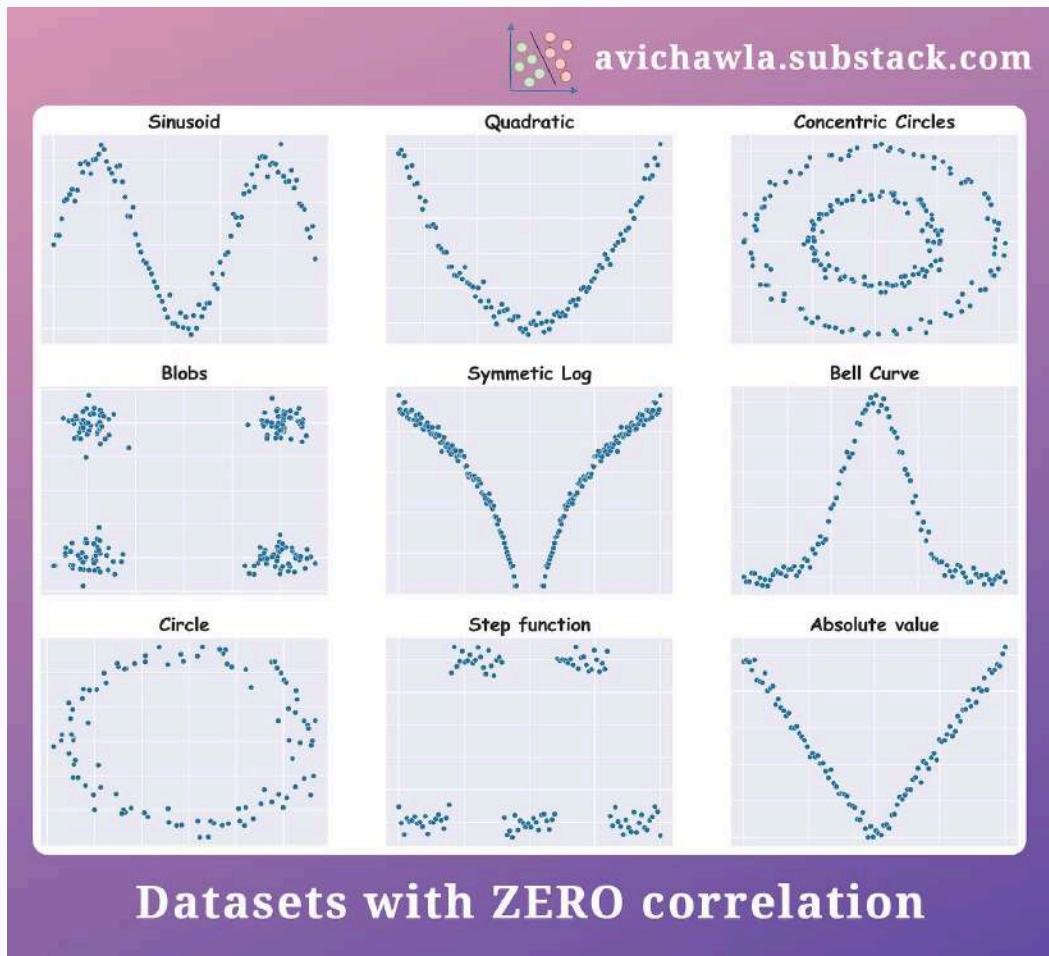
- Pearson and Spearman correlation is the same on linear data.
- But Pearson correlation underestimates a non-linear association.

Spearman correlation is also useful when data is ranked or ordinal, which will shall discuss in the next chapter.

Also, before we end this chapter, remember to always be cautious before drawing any conclusions using summary statistics. We also saw this in the last chapter.

While analyzing data, so many people get tempted to draw conclusions solely based on its statistics. Yet, the actual data might be conveying a totally different story.

This is also evident from the image below:



All nine datasets have approx. zero correlation between the two variables. However, the summary statistic, Pearson correlation in this case, gives no clue about what's inside the data because it is always zero.

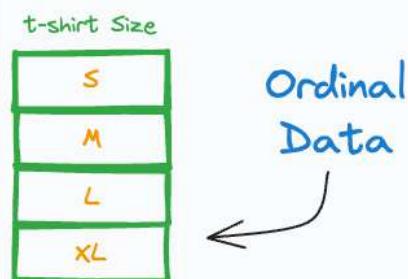
In fact, this is not just about Pearson correlation but applies to all summary statistics. The idea is that whenever you generate any summary statistic, you lose essential information.

Thus, the importance of looking at the data cannot be stressed enough. It saves us from drawing wrong conclusions, which we could have made otherwise by looking at the statistics alone.

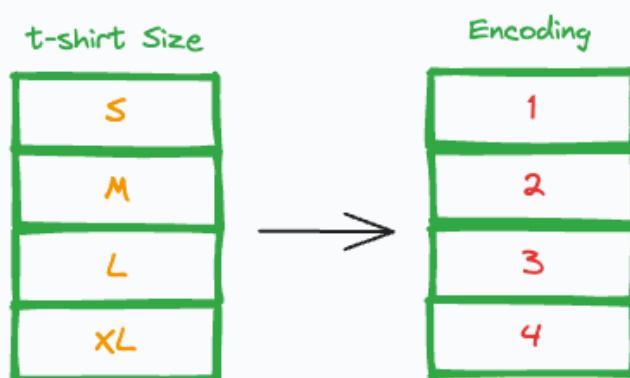
Correlation with Ordinal Categorical Data

Imagine you have an ordinal categorical feature. You want to measure its correlation with other continuous features.

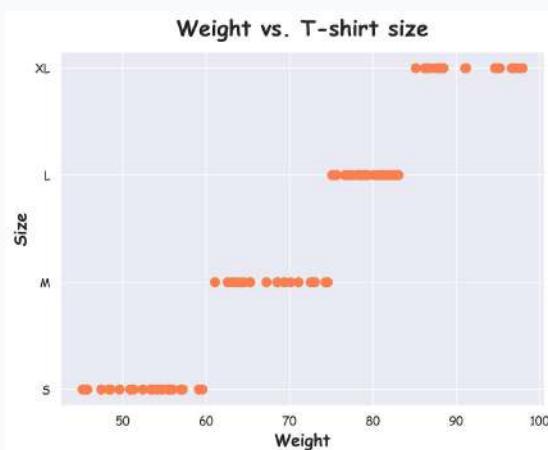
Ordinal feature is categorical data with a natural ordering in categories.



Before proceeding with the correlation analysis, you will encode the feature, which is a fair thing to do.



Yet, unknown to many, the choice of encoding can largely affect the correlation results. For instance, consider the dataset below:



Here, we have:

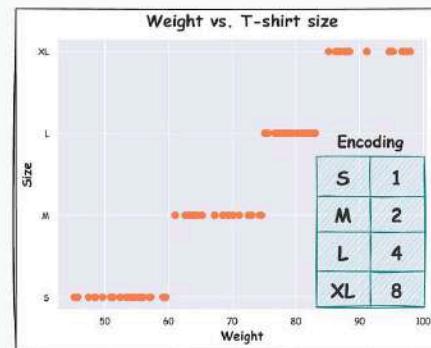
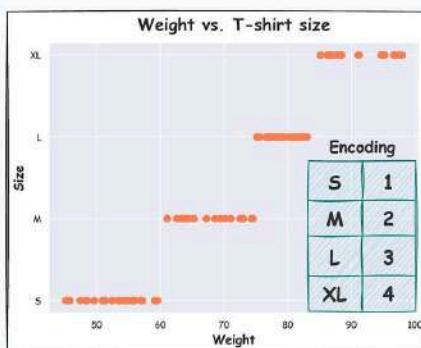
- An ordinal categorical feature: t-shirt size (S, M, L, XL).
- A continuous feature: weight.

If we look at it graphically, there does exist a monotonic relationship between the two features. However, as depicted below, altering the categorical encoding affects the Pearson correlation.

Never Use Pearson Correlation on Ordinal Categorical Data



blog.DailyDoseofDS.com



Pearson Correlation **0.96**

Dropped

Pearson Correlation **0.90**

Spearman Correlation **0.96**

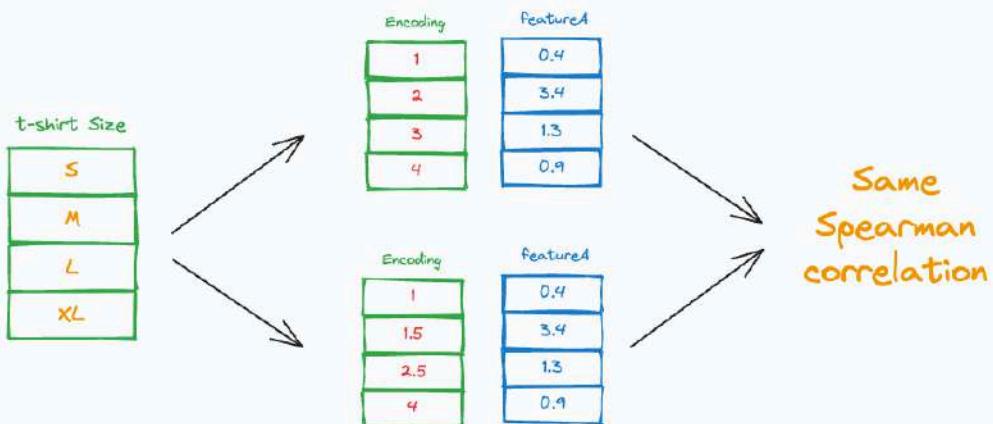
Unchanged

Spearman Correlation **0.96**

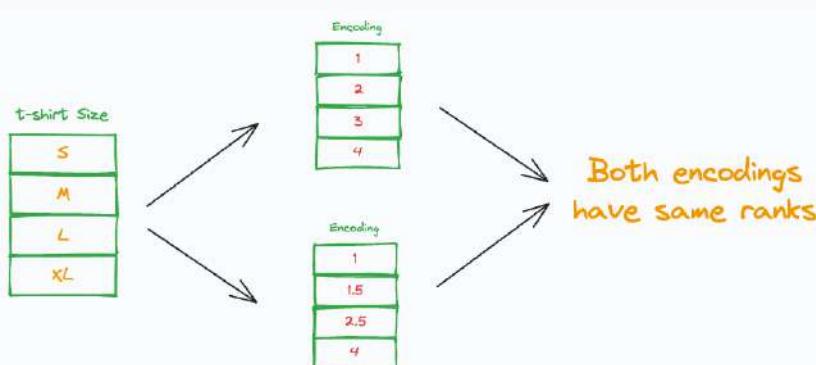
- In the left plot, we have the following encoding → S(1), M(2), L(3) and XL(4).
- In the right plot, we have the following encoding → S(1), M(2), L(4) and XL(8).

It is clear that the Spearman correlation is a better alternative to assess the monotonicity between ordinal and continuous features.

It always remains the same, irrespective of the choice of categorical encoding. This is because the Spearman correlation is rank-based.



It operates on the ranks of the data, which makes it more suitable for such cases of correlation analysis.



Using Spearman correlation is pretty simple as well. For instance, if you are using Pandas, just specify the desired correlation measure as follows:

```

Correlation.py

df = ...

df.corr(method = "pearson")
df.corr(method = "spearman")

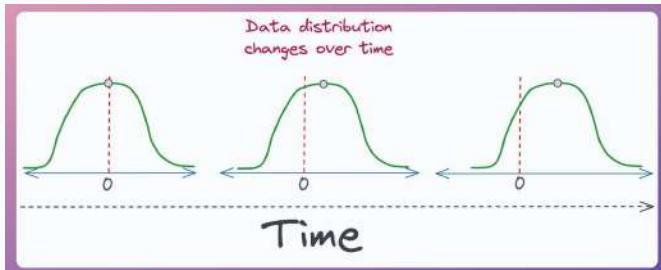
```

Pearson
Spearman

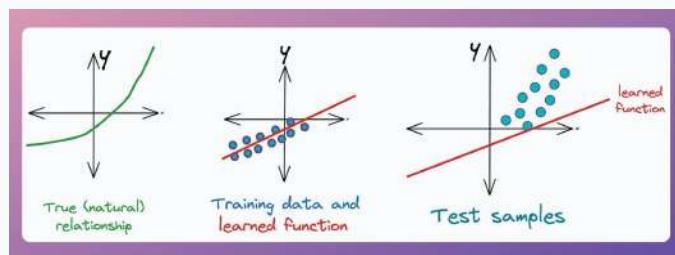
Drift

Multivariate Covariate Shift

Almost all real-world ML models gradually degrade in performance due to covariate shift.



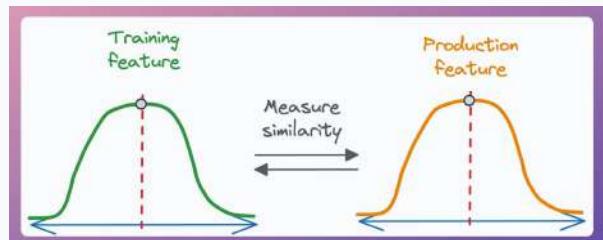
For starters, covariate shift happens when the distribution of features changes over time, but the true (natural) relationship between input and output remains the same.



It is a serious problem because we trained the model on one distribution, but it is being used to predict on another distribution in production. Thus, it is critical to detect covariate shift early so that models continue to work well.

Issues with detection methods

One of the most common and intuitive ways to detect covariate shift is by simply comparing the feature distribution in training data to that in production.

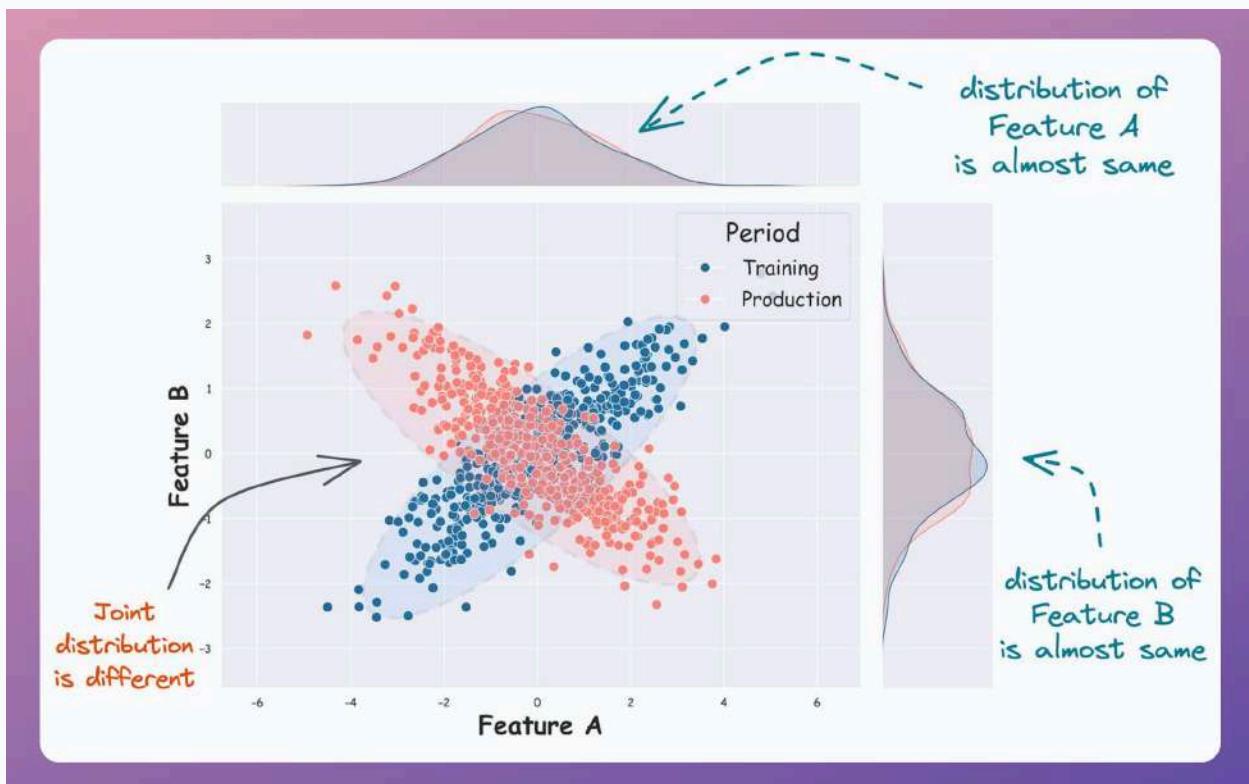


This could be done in many ways, such as:

- Compare their summary statistics — mean, median, etc.
- Inspect visually using distribution plots.
- Perform hypothesis testing.
- Measure distances between training/production distributions using Bhattacharyya distance, KS test, etc.

While these approaches are often effective, the biggest problem is that they work on a single feature at a time. But, in real life, we may observe multivariate covariate shift (MCS) as well. It happens when:

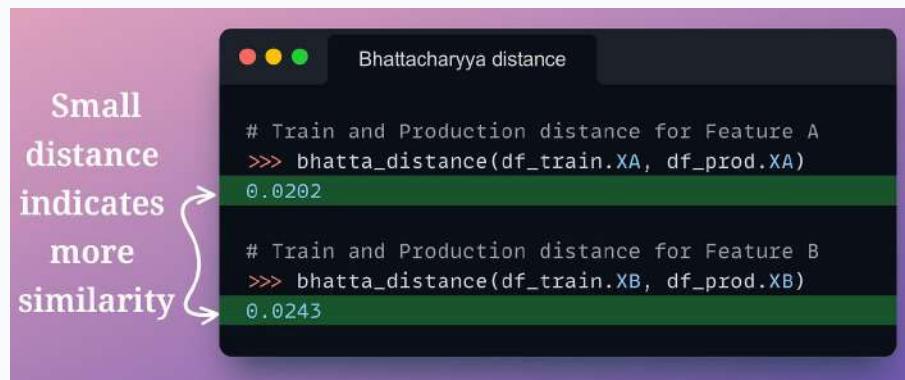
- The distribution of individual distributions remains the same:
 - $P(X_1)$ and $P(X_2)$ individually remain the same.
- But their joint distribution changes:
 - $P(X_1, X_2)$ changes.



- From the KDE plots on the top and the right, it is clear that the distribution of both features (covariates) is almost the same.

- But, the scatter plot reveals that their joint distribution in training (Blue) differs from that in production (Red).

And it is easy to guess that the univariate covariate shift detection methods discussed above will produce misleading results. For instance, as demonstrated below, measuring the Bhattacharyya distance between a training and production feature gives a very low distance value, indicating high similarity:

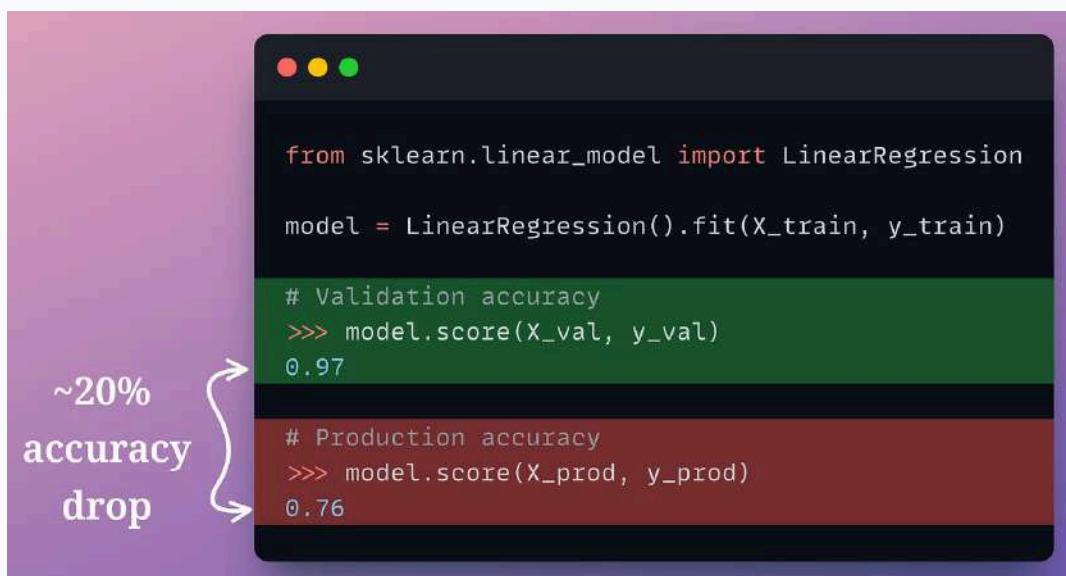


```
# Train and Production distance for Feature A
>>> bhatta_distance(df_train.XA, df_prod.XA)
0.0202

# Train and Production distance for Feature B
>>> bhatta_distance(df_train.XB, df_prod.XB)
0.0243
```

In fact, even though the individual feature distribution is the same, we can confirm experimentally that this will result in a drop in model performance:

Let's say the true output $y = 2*\text{Feature_A} + 3*\text{Feature_B} + \text{noise}$.



```
from sklearn.linear_model import LinearRegression

model = LinearRegression().fit(X_train, y_train)

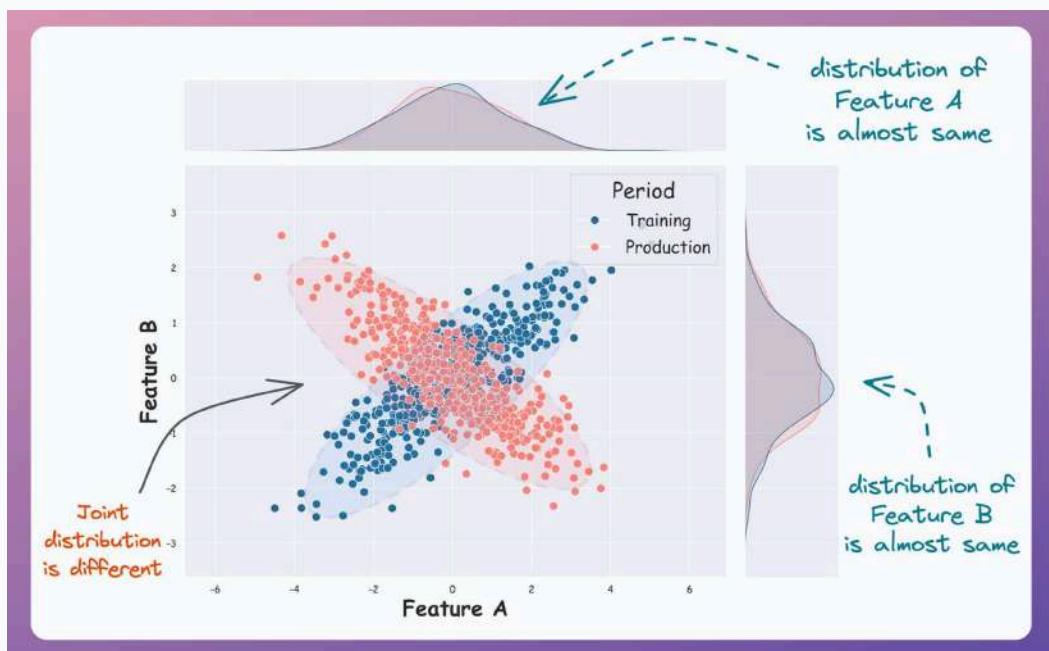
# Validation accuracy
>>> model.score(X_val, y_val)
0.97

# Production accuracy
>>> model.score(X_prod, y_prod)
0.76
```

In most cases, the true output predictions on production data are not immediately available. In the above demonstration, just assume that we have already gathered the true predictions somehow.

As depicted above, the model performance drops by 20% in production, which is huge. Here, we may completely rule out the possibility of covariate shift if we think that covariate shift can never be multivariate in nature.

Of course, in the figure below, it was easy to identify multivariate covariate shift because we are only looking at two features.



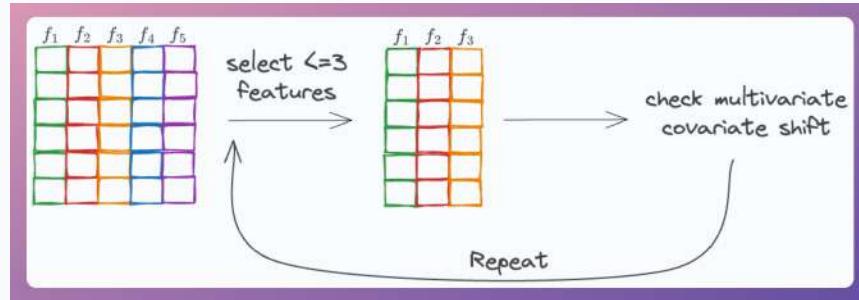
But multivariate covariate shift can happen with more than two features as well. Unlike the bivariate case above, visual inspection will not be possible for higher dimensions.

Detecting multivariate covariate shift

To begin, it is important to understand that multivariate covariate shift is a big problem, and there is no direct (or single) approach to handle this. Below, I will share a couple of ideas that I often use myself and have seen others use as well to handle multivariate covariate shift.

Idea #1: MCS is rare

As long as we are checking two (or even three) features at a time, visual inspection can be used to detect covariate shift. So here, at times, many simply ignore the possibility of any covariate shift beyond three features.



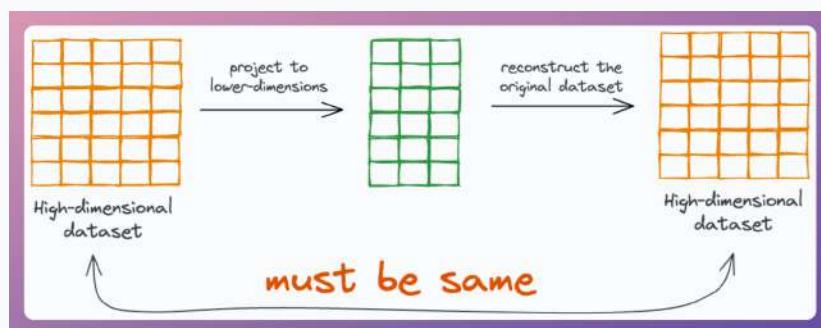
The rationale is that beyond three features, it is pretty unlikely that:

- $P(X_1), P(X_2), P(X_3), P(X_4), \dots \rightarrow$ all of them individually will almost remain the same.
- But their joint distribution $P(X_1, X_2, X_3, X_4, \dots)$ will change.

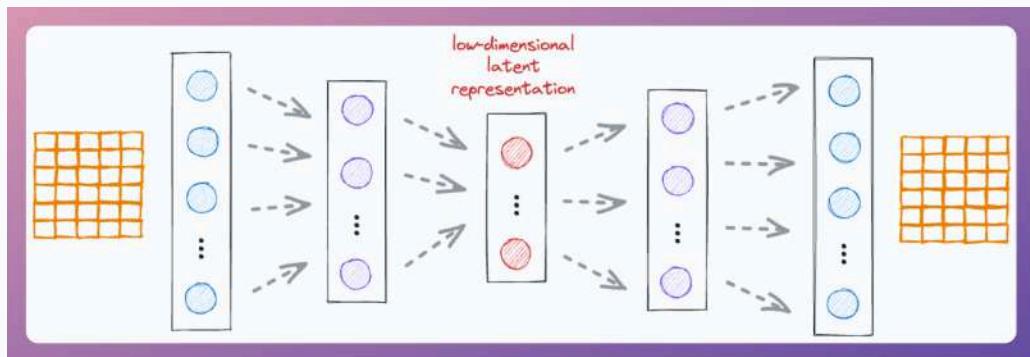
Thus, it might be fair to limit multivariate covariate shift analysis to just one, two, and three features at a time. But of course, this may not be always true, which brings us to another idea:

Idea #2: Data reconstruction

This is another cool and practical idea that I used in one of my projects. Data reconstruction, as the name suggests, revolves around learning a mapping that projects the data to low dimensions and then reconstructs the original data from the low dimensions.



This is precisely what Autoencoders do. They are a class of neural networks that learn to encode data in a lower-dimensional space, and then decode it back to the original data space.

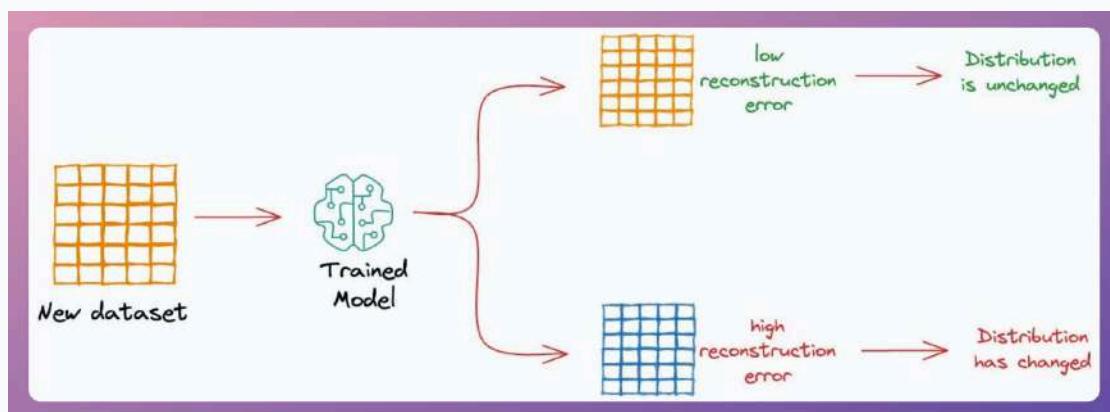


The objective is to minimize the data reconstruction error. It's like asking the model to learn a mapping that:

- Takes some data.
- Projects the data to low dimensions.
- And then gives us the exact input data back.

So here's what we can do:

- Train an Autoencoder on the original training dataset. This will give us the weights for the neural network model that reconstructs the dataset.
- Use this model on new data to check multivariate covariate shift:
 - If the reconstruction loss is high, the distribution has changed.
 - If the reconstruction loss is low, the distribution is almost the same.



One of the best things about this approach is that it does not need a labeled dataset. Essentially, Autoencoders aim to reconstruct the dataset without labels.

This is quite useful in real-world models because, as we discussed yesterday, in most cases, the true output predictions on production data are never immediately available. Instead, they always take some time.

For instance, I remember when I was working on a transactional fraud detection model at Mastercard, a card holder's issuer bank may take upto 45-50 days to send the fraud label for the transactions that went through Mastercard's network.

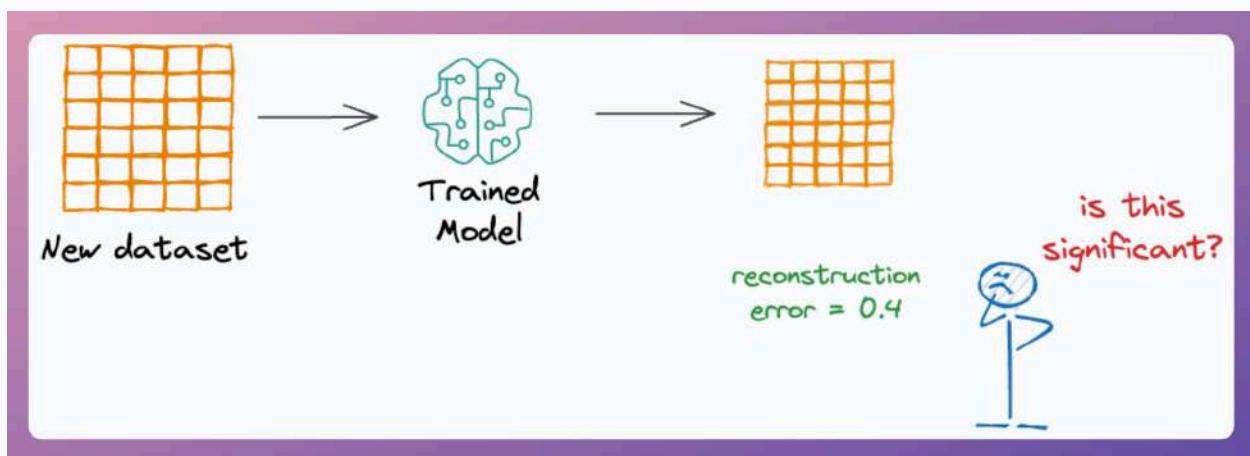
This is a lot of time, isn't it?

But using Autoencoders, we can still check data reconstruction errors on the unlabeled data. In fact, we don't necessarily have to use Autoencoders. It was just an example here.

Other data reconstruction techniques can also be used, such as PCA is one of them. Nonetheless, when using any data reconstruction approach, it is important to take care of one thing. Earlier, we discussed that:

- If the reconstruction loss is high, the distribution has changed.
- If the reconstruction loss is low, the distribution is almost the same.

But interpreting reconstruction loss can be pretty subjective, and it also needs some context.

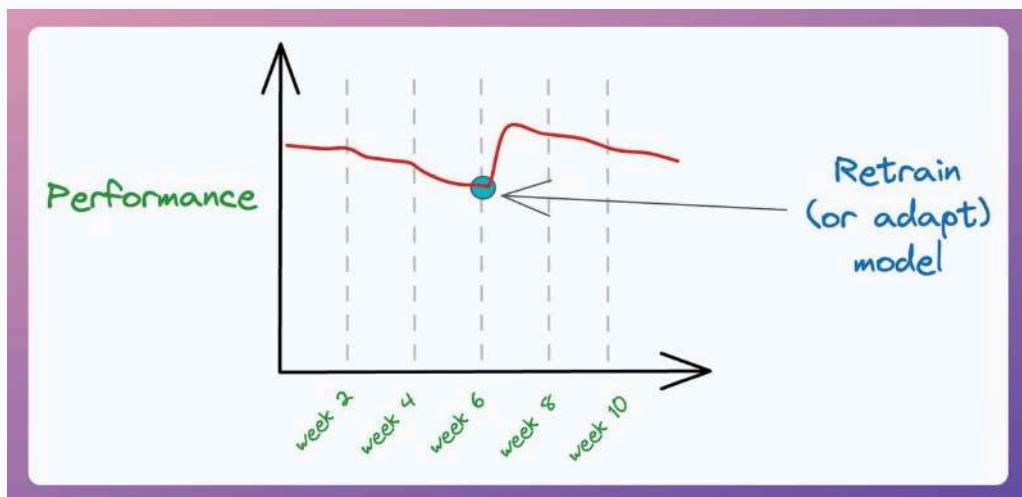


For instance, if the reconstruction loss is 0.4 (say), how do we determine whether this is significant? Contrary to what many think, covariate shift happens gradually in ML models. In other words, it is unlikely that our model is performing pretty well one day, and the next day, it starts underperforming DUE TO COVARIATE SHIFT.



Of course, abrupt degradation is possible but it is rarely due to covariate shift. Most of the time, it happens due to other issues like - May be due to some compliance issue, the team has stopped collecting a feature. As a result, your model no longer has access to a feature you trained it with. Or may be there were some updates to the inference logic which was not thoroughly tested and more.

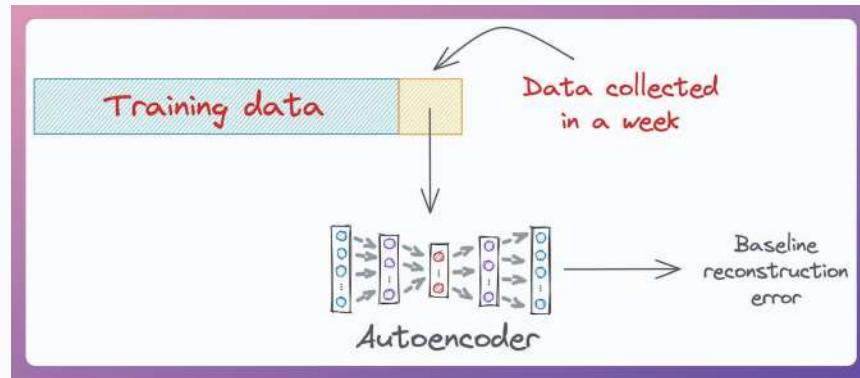
But performance degradation happens gradually due to covariate shift. And in that phase of performance degradation, we must decide to update the model at some point to adjust for distributional changes.



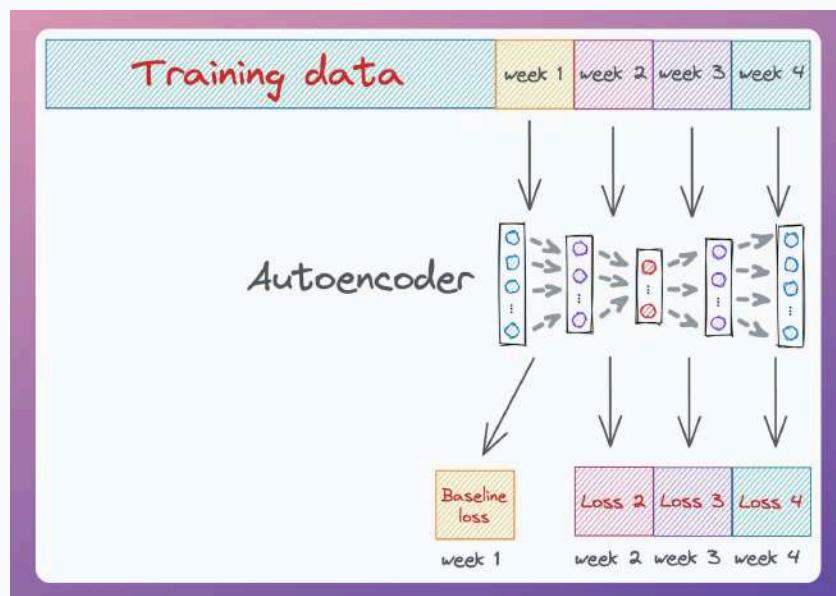
While data reconstruction loss is helpful in deciding model updates, it conveys very little meaning alone. Simply put, we need a baseline reconstruction loss value to compare our future reconstruction losses.

So here's what we can do.

Let's say we decided to review our model every week. After training the model on the gathered data:



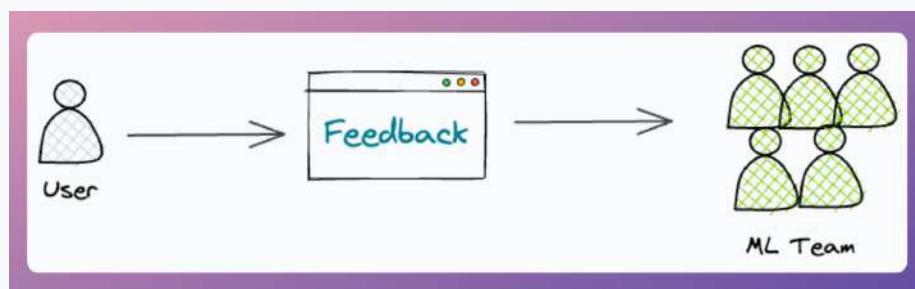
- We determine our baseline reconstruction score using the new data gathered over the first week in the post-training phase.
- We use this baseline score to compare future losses of every subsequent week.



A high difference between reconstruction loss and the baseline score will indicate a covariate shift.

Simple and intuitive, right? The above reconstruction metric, coupled with the model performance, becomes a reliable way to determine if covariate shift has stepped in. Of course, as discussed above, true labels might not be immediately available at times.

Thus, determining performance in production can be difficult. To counter this, all ML teams consistently try to gather feedback from the end user to learn if the model did well or not.



This is especially seen in recommendation systems:

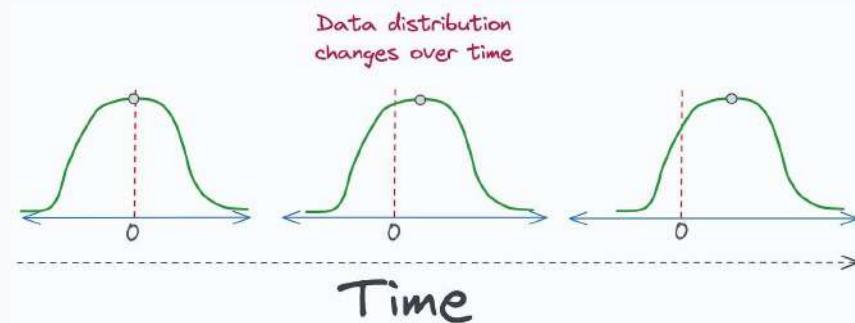
- Was this advertisement relevant to you?
- Was this new video relevant to you?
- And more.

These things are a part of the model logging phase, where we track model performance in production. This helps us decide whether we should update our model or not.



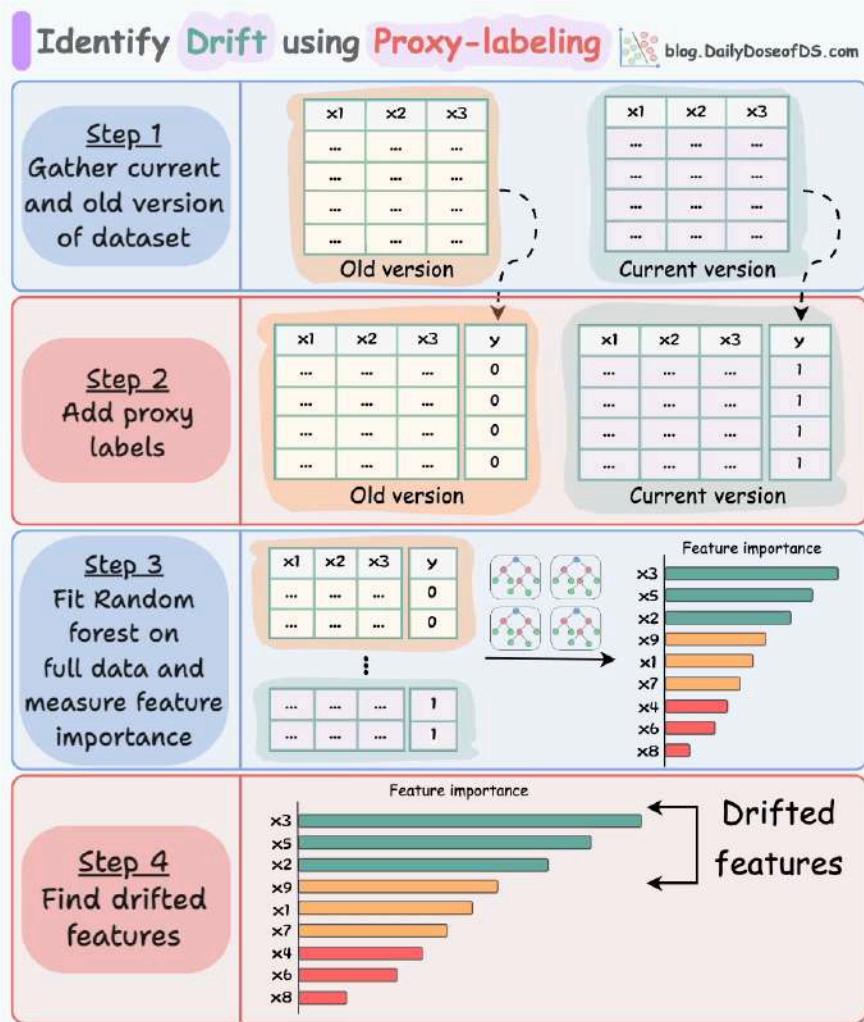
Using Proxy-Labelling to Identify Drift

Almost all real-world ML models gradually degrade in performance due to a drift in feature distribution:



It is a serious problem because we trained the model on one distribution, but it is being used to generate predictions on another distribution in production.

Thus, it is critical to detect drift early so that models continue to work well. The image depicts an intuitive technique I often use to determine which features are drifting in my dataset.



Consider we have two versions of the dataset — the old version (one on which the model was trained) and the current version (one on which the model is generating predictions):

Old dataset					
	0	1	2	3	4
0					
1					
2					
3					

Current dataset					
	0	1	2	3	4
0					
1					
2					
3					

The core idea is to assess if there are any distributional dissimilarities between these two versions.

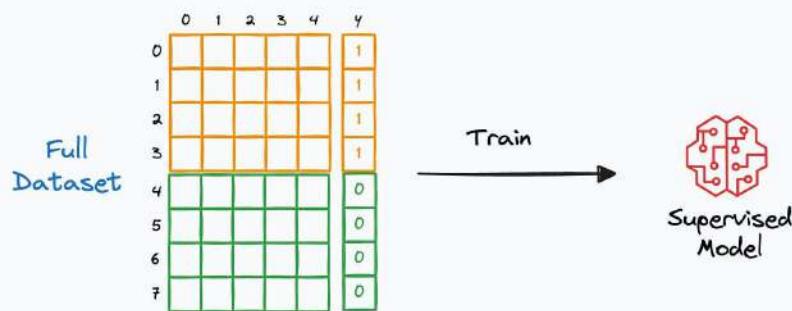
So here's what we do:

- Append a `label=1` column to the old dataset.
- Append a `label=0` column to the current dataset.

Old dataset						5
	0	1	2	3	4	5
0						1
1						1
2						1
3						1

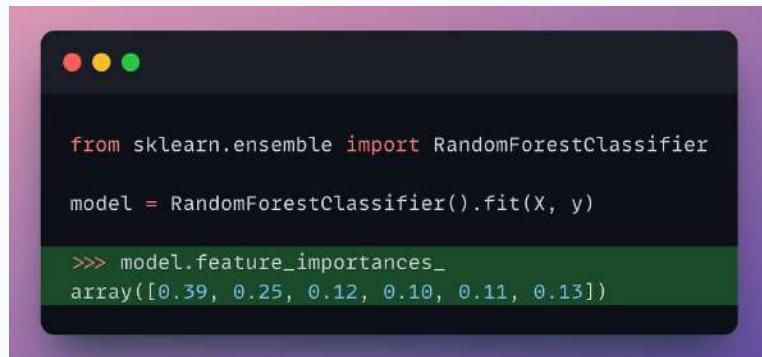
Current dataset						5
	0	1	2	3	4	5
0						0
1						0
2						0
3						0

Now, merge these two datasets and train a supervised learning classification model on the combined dataset:



The choice of the classification model could be a bit arbitrary, but it should be ensured that it is possible to reliably determine feature importance.

Thus, I personally prefer a random forest classifier because it has an inherent mechanism to determine feature importance:

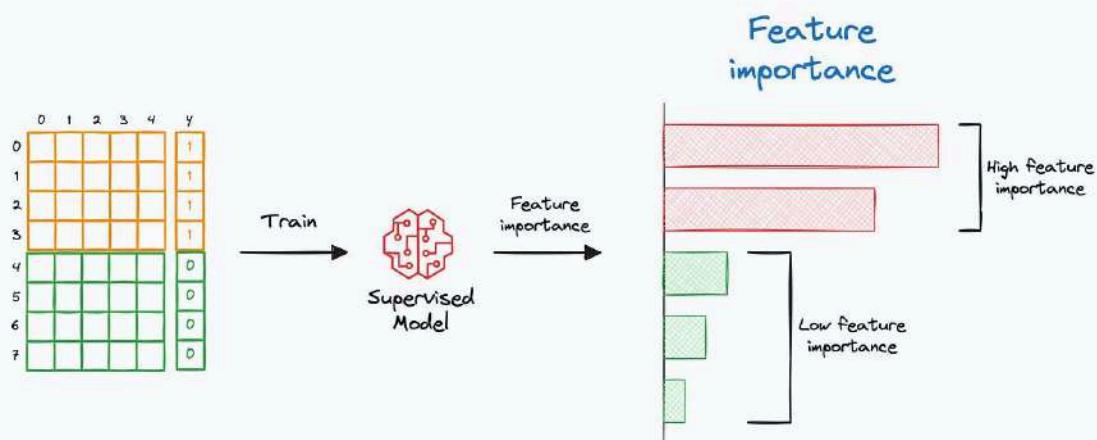


```
from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier().fit(X, y)
>>> model.feature_importances_
array([0.39, 0.25, 0.12, 0.1, 0.11, 0.13])
```

That said, it is not necessary to use a random forest.

Techniques like shuffle feature importance (which we discussed above) on top of a classification model can be used as well.

If the feature importance values suggest that there are features with high feature importance, this means that those features are drifting.



Why?

This is because if some features can reliably distinguish between the two versions of the dataset, then it is pretty likely that their distribution corresponding to `label=1` and `label=0` (conditional distribution) are varying.

- If there are distributional differences, the model will capture them.
- If there are no distributional differences, the model will struggle to distinguish between the classes.

This idea makes intuitive sense as well.

At this point, one question that many have is that...

Why can't we just monitor the model accuracy to determine drift?

Of course, we can do that as long as we have the true labels for the current version or have a way to compare the model performance on the old version and the new version.

But in many cases, the true output predictions on production data are never immediately available.

Instead, they always take some time.

For instance, when I was working on a transactional fraud detection model at Mastercard, a cardholder's issuer bank may take up to 45-50 days to send the fraud label for the transactions that went through Mastercard's network.

This is a lot of time, isn't it?

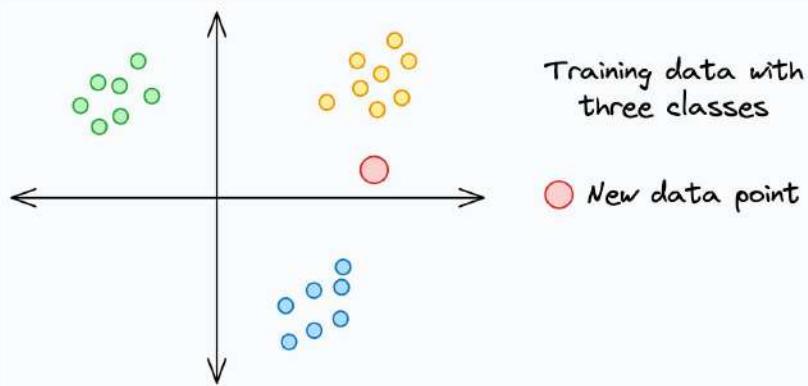
Thus, one must rely on some relevant feedback from the system to determine whether the model's performance is dropping or not.

Using this “proxy-labeling technique” discussed above is something I have found to be immensely useful.

kNN

Using kNNs on Imbalanced Datasets

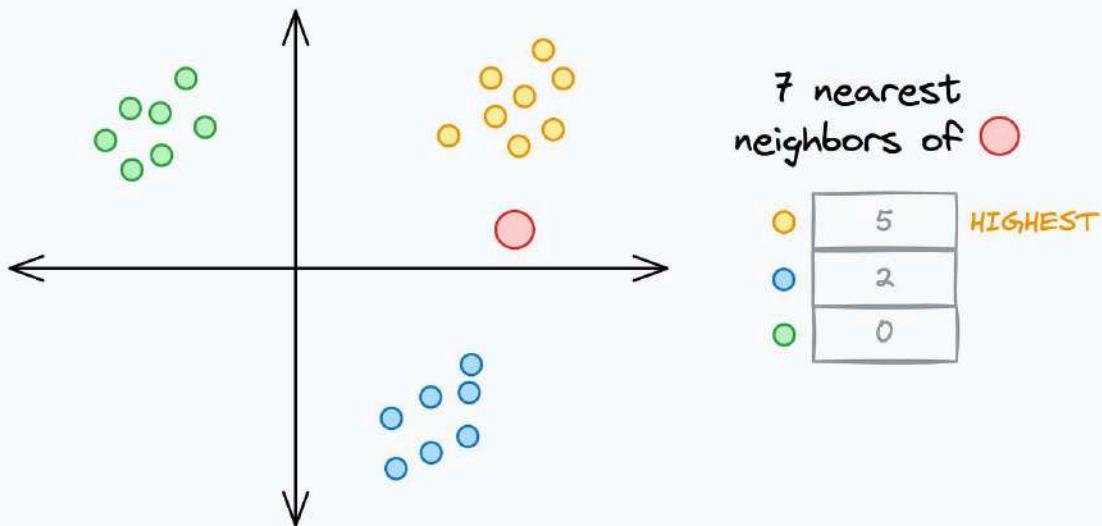
One of the things that always makes me a bit cautious and skeptical when using kNN is its HIGH sensitivity to the parameter k . To understand better, consider this dummy 2D dataset below. The red data point is a test instance we intend to generate a prediction for using kNN.



Say we set the value of $k=7$. The prediction for the red instance is generated in two steps:

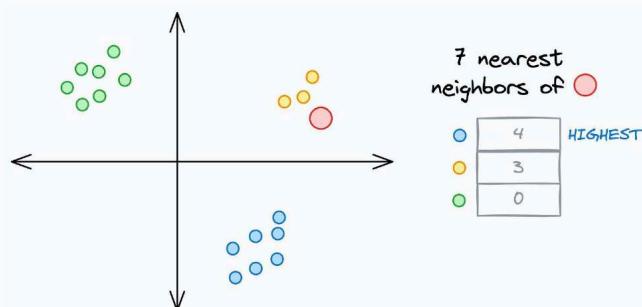
- First, we count the 7 nearest neighbors of the red data point.
- Next, we assign it to the class with the highest count among those 7 nearest neighbors.

This is depicted below:



The problem is that step 2 is entirely based on the notion of class contribution — the class that maximally contributes to the k nearest neighbors is assigned to the data point. But this can miserably fail at times, especially when we have a class with few samples.

For instance, as shown in this image, with $k=7$, the red data point can NEVER be assigned to the yellow class, no matter how close it is to that cluster.



While it is easy to tweak the hyperparameter k visually in the above demo, this approach is infeasible in high-dimensional datasets. There are two ways to address this.

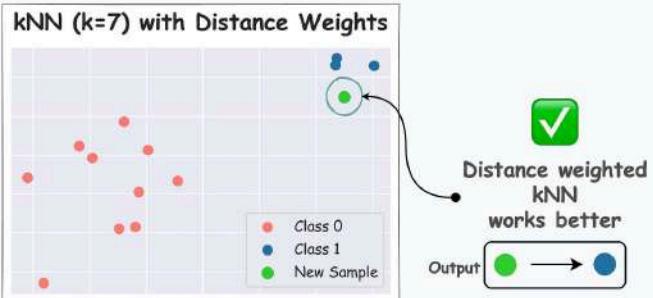
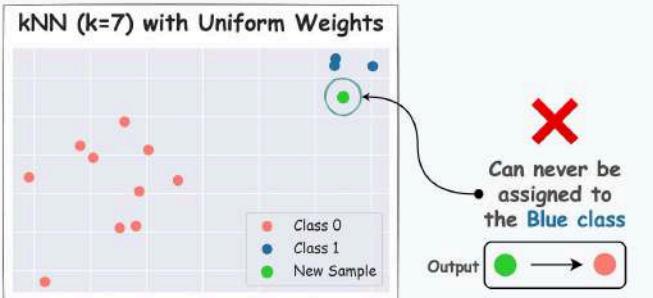
Solution #1: Used distance-weighted kNN

Distance-weighted kNNs are a much more robust alternative to traditional kNNs. As the name suggests, in step 2, they consider the distance to the nearest neighbor.

As a result, the closer a specific neighbor is, the more impact it will have on the final prediction. Its effectiveness is evident from this image.

An Underrated Technique to Supercharge kNNs

 blog.DailyDoseofDS.com



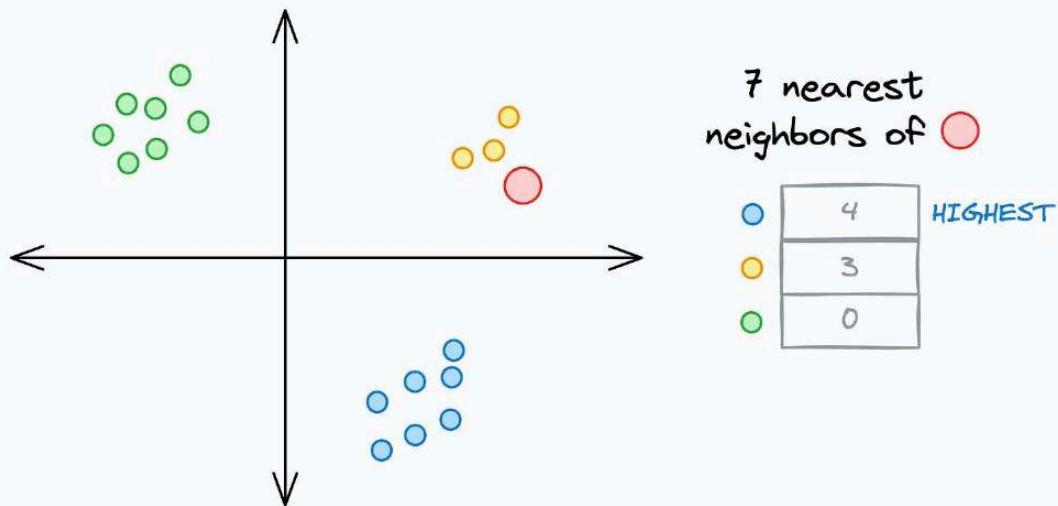
- In the first plot, traditional kNN (with $k=7$) can never predict the blue class.
- In the second plot, distance-weighted kNN is found to be more robust in its prediction.

As per my observation, a distance-weighted kNN typically works better than a traditional kNN. Yet, this may go unnoticed because, by default, the kNN implementation of sklearn considers “uniform” weighting.



Solution #2: Dynamically update the hyperparameter k

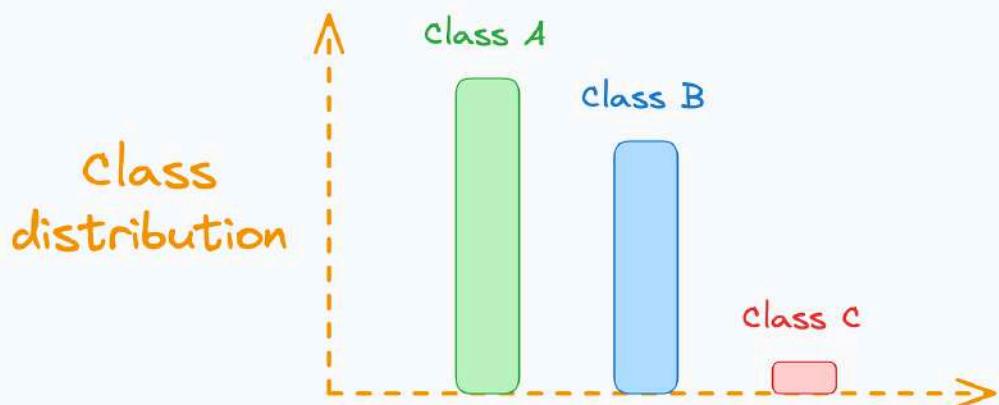
Recall the above demo again:



Here, one may argue that we must refrain from setting the hyperparameter k to any value greater than the minimum number of samples that belong to a class in the dataset.

Of course, I agree with this to an extent.

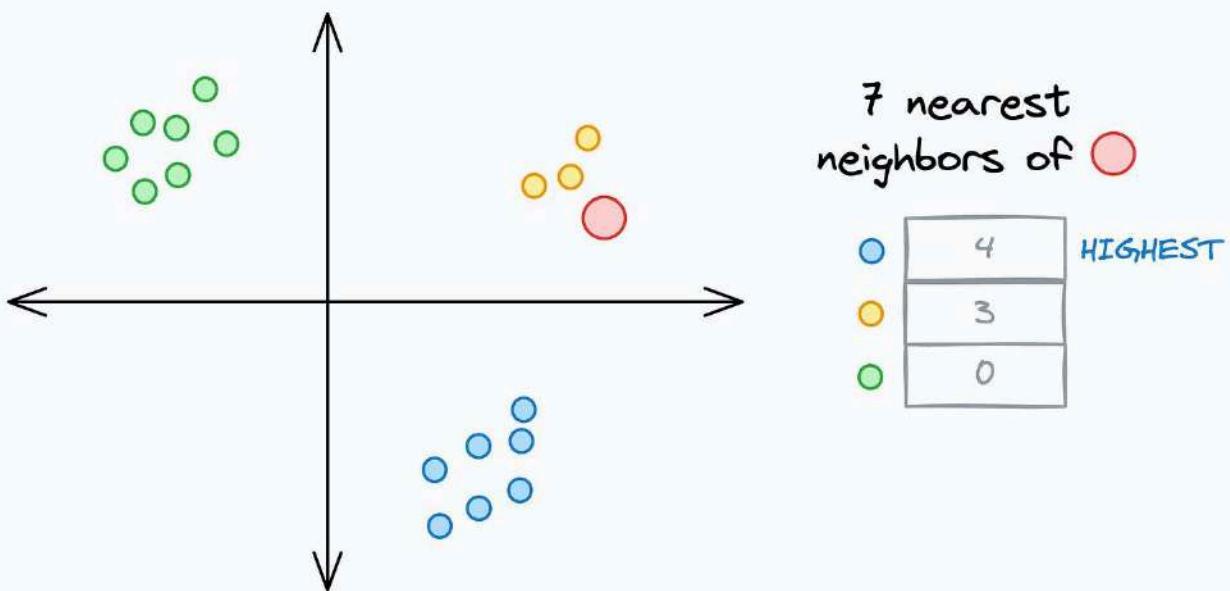
But let me tell you the downside of doing that. Setting a very low value of k can be highly problematic in the case of extremely imbalanced datasets.



To give you more perspective, I have personally used kNN on datasets that had merely one or two instances for a particular class in the training set. And I

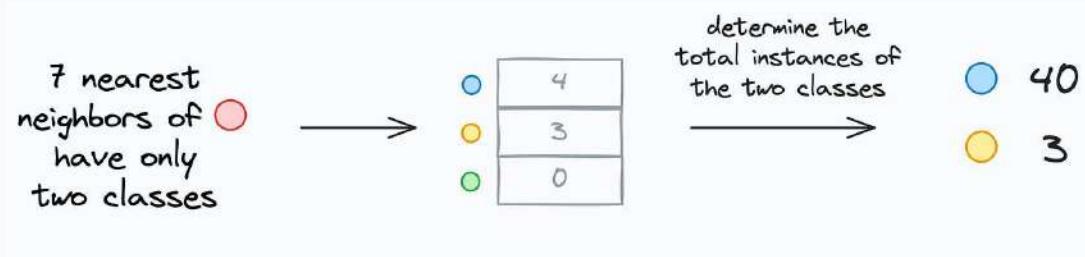
discovered that setting a low of k (say, 1 or 2) led to suboptimal performance because the model was not as holistically evaluating the nearest neighbor patterns as it was when a large value of k was used.

In other words, setting a relatively larger value of k typically gives more informed predictions than using lower values. But we just discussed above that if we set a large value of k, the majority class can dominate the classification result:



To address this, I found dynamically updating the hyperparameter k to be much more effective. More specifically, there are three steps in this approach. For every test instance:

1. Begin with a standard value of k as we usually would and find the k nearest neighbors.
2. Next, update the value of the k as follows:
 - a. For all unique classes that appear in the k nearest neighbor, find the total number of training instances they have.



b. Update the value of k to:

update k value

$$k' = \min(k, 40, 3)$$

3. Now perform majority voting only on the first k' neighbors only.

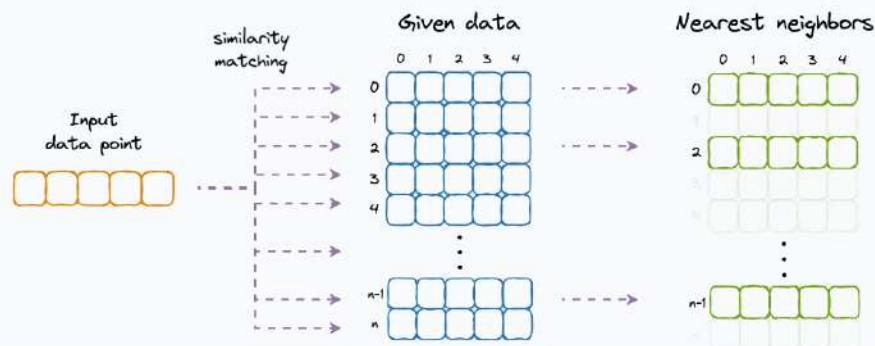
This makes an intuitive sense as well:

- If a minority class appears in the top k nearest neighbor, we must reduce the value of k so that the majority class does not dominate.
- If a minority class DOES NOT appear in the top k nearest neighbor, we will likely not update the value of k and proceed with a holistic classification.

The only shortcoming is that you wouldn't find this approach in any open-source implementations. In fact, in my projects as well, I had to write a custom implementation, so take that into account.

Approximate Nearest Neighbor Search Using Inverted File Index

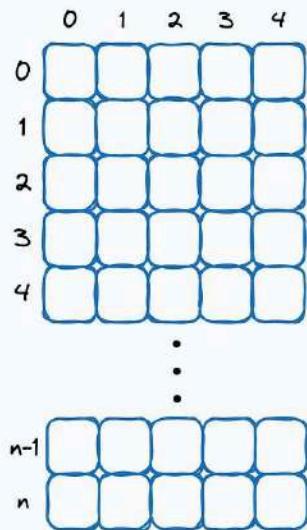
One of the biggest issues with nearest neighbor search using kNN is that it performs an exhaustive search.



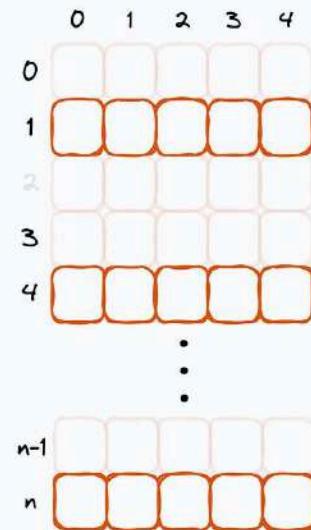
In other words, the query data point must be matched across all data points to find the nearest neighbor(s). This is highly inefficient, especially when we have many data points and a near-real-time response is necessary.

That is why approximate nearest neighbor search algorithms are becoming increasingly popular. The core idea is to narrow down the search space using indexing techniques, thereby improving the overall run-time performance.

kNN searches in this



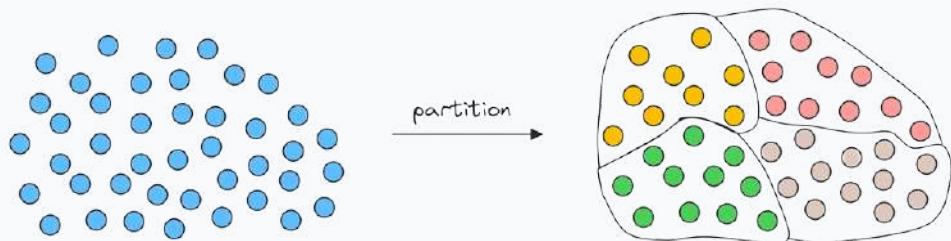
ANN searches in this



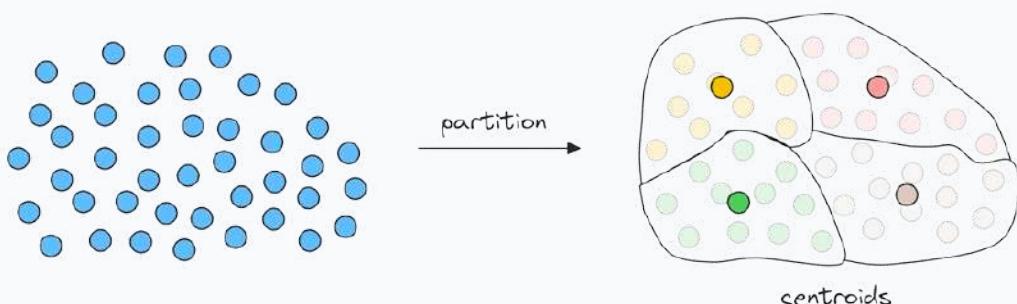
Inverted File Index (IVF) is possibly one of the simplest and most intuitive techniques here, which you can immediately start using.

Here's how it works

Given a set of data points in a high-dimensional space, the idea is to organize them into different partitions, typically using clustering algorithms like k-means.



As a result, each partition has a corresponding centroid, and every data point gets associated with only one partition corresponding to its nearest centroid.



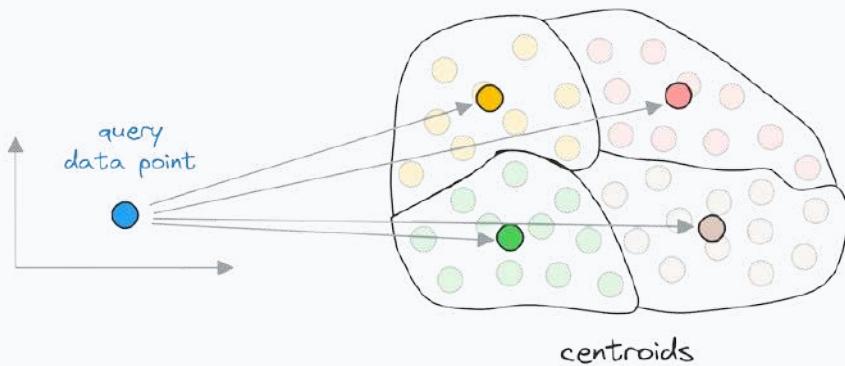
Also, every centroid maintains information about all the data points that belong to its partition.



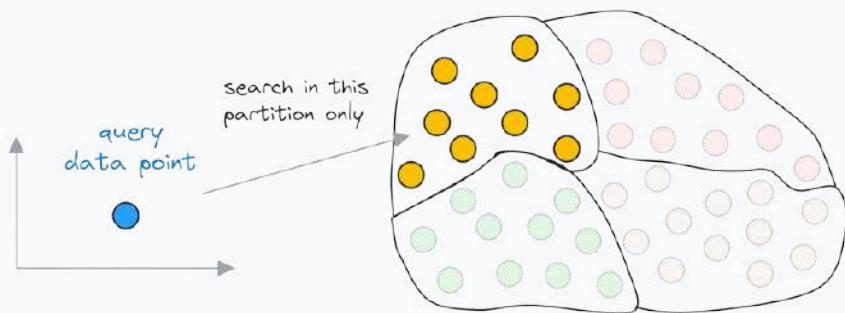
Indexing done!

Here's how we search.

When searching for the nearest neighbor(s) to the query data point, instead of searching across the entire dataset, we first find the closest centroid to the query:



Once we find the nearest centroid, the nearest neighbor is searched in only those data points that belong to the closest partition found:



Let's see how the run-time complexity stands in comparison to traditional kNN. Consider the following:

- There are N data points
- Each data point is D dimensional
- We create K partitions.
- Lastly, for simplicity, let's assume that each partition gets equal data points.

In kNN, the query data point is matched to all N data points, which makes the time complexity $\rightarrow O(ND)$.

In IVF, however, there are two steps:

1. Match to all centroids $\rightarrow O(KD)$.
2. Find the nearest neighbor in the nearest partition $\rightarrow O(ND/K)$.

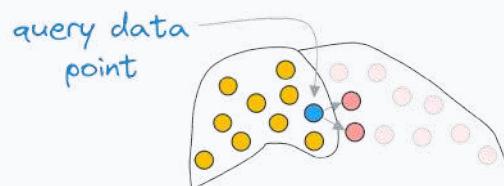
The final time complexity comes out to be the following:

$$O(KD + \frac{ND}{K})$$

...which is significantly lower than that of kNN. To get some perspective, assume we have 10M data points. The search complexity of kNN will be proportional to 10M. But with IVF, say we divide the data into 100 centroids, and each partition gets roughly 100k data points.

Thus, the time complexity comes out to be proportional to $100 + 100k = 100100$, which is nearly 100 times faster.

Of course, it is essential to note that if some data points are actually close to the input data point but still happen to be in the neighboring partition, we will miss them during the nearest neighbor search, as shown in this image.



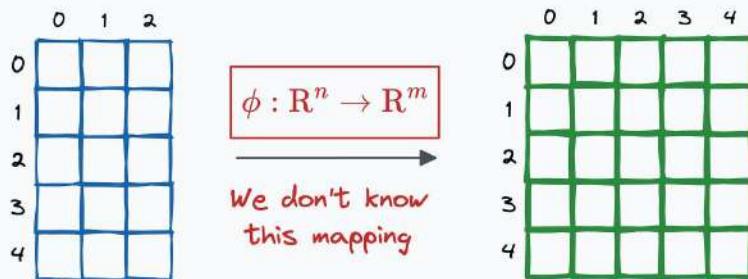
But this accuracy tradeoff is something we willingly accept for better run-time performance, which is precisely why these techniques are called “approximate nearest neighbors search.”

Kernels

Why is Kernel Trick Called a "Trick"?

So many ML algorithms use kernels for robust modeling, like SVM, KernelPCA, etc.

In a gist, a kernel function lets us compute dot products in some other feature space (mostly high-dimensional) without even knowing the mapping from the current space to the other space.



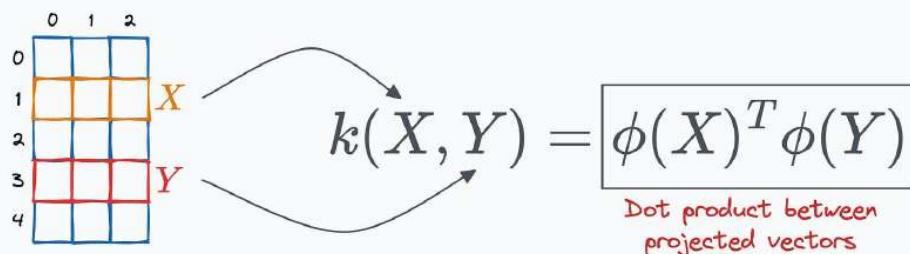
But how does that even happen?

Let's understand!

The objective

Firstly, it is important to note that the kernel provides a way to compute the dot product between two vectors, X and Y , in some high-dimensional space without projecting the vectors to that space.

This is depicted below, where the output of the kernel function is expected to be the same as the dot product between projected vectors:



The key advantage is that the kernel function is applied to the vectors in the original feature space.

However, that equals the dot product between the two vectors when projected into a higher-dimensional (yet unknown) space.

If that is a bit confusing, let me give an example.

A motivating example

Let's assume the following polynomial kernel function:

Polynomial kernel

$$k(X, Y) = (1 + X^T Y)^2$$

For simplicity, let's say both X and Y are two-dimensional vectors:

$$\begin{array}{c} \boxed{} \\ \boxed{} \end{array} \quad \begin{array}{c} \boxed{} \\ \boxed{} \end{array}$$

$$X = (x_1, x_2) \quad Y = (y_1, y_2)$$

Simplifying the kernel expression above, we get the following:

$$\begin{aligned} k(X, Y) &= (1 + X^T Y)^2 \\ &= (1 + x_1 y_1 + x_2 y_2)^2 \end{aligned}$$

Expanding the square term, we get:

$$\begin{aligned} k(X, Y) &= (1 + x_1 y_1 + x_2 y_2)^2 \\ &\quad \text{expand the square term} \\ &= 1 + x_1^2 y_1^2 + x_2^2 y_2^2 + \\ &\quad 2x_1 y_1 + 2x_2 y_2 + 2x_1 y_1 x_2 y_2 \end{aligned}$$

Now notice the final expression:

$$k(X, Y) = 1 + x_1^2 y_1^2 + x_2^2 y_2^2 + 2x_1 y_1 + 2x_2 y_2 + 2x_1 y_1 x_2 y_2$$

The above expression is the dot product between the following 6-dimensional vectors:

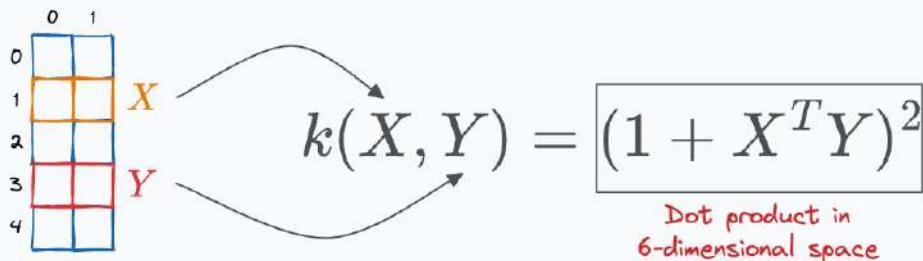
$$\begin{aligned} X' &= (1, x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2) \\ Y' &= (1, y_1^2, y_2^2, \sqrt{2}y_1, \sqrt{2}y_2, \sqrt{2}y_1 y_2) \end{aligned}$$

Thus, our projection function comes out to be:

Projection function

$$\phi(X) = (1, x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2)$$

This shows that the kernel function we chose earlier computes the dot product in a 6-dimensional space without explicitly visiting that space.

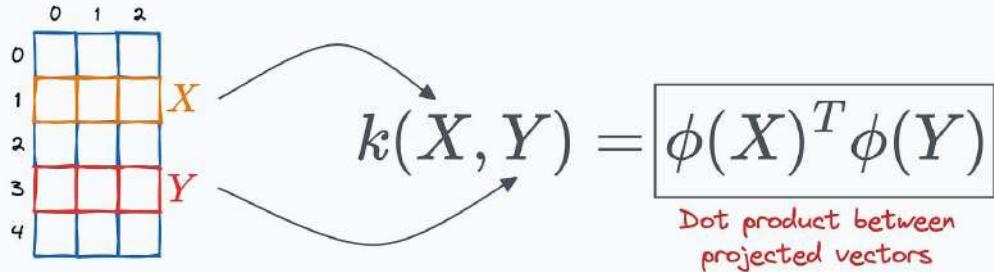


And that is the primary reason why we also call it the “kernel trick.”

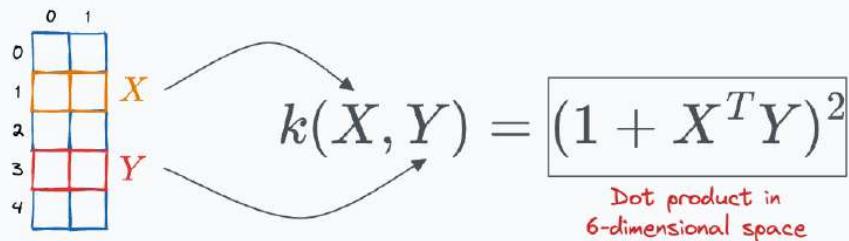
More specifically, it’s framed as a “trick” since it allows us to operate in high-dimensional spaces without explicitly computing the coordinates of the data in that space.

The Mathematics Behind RBF Kernel

As we saw above, a kernel provides a way to compute the dot product between two vectors, X and Y, in some high-dimensional space without projecting the vectors to that space.



In that post, we looked at the polynomial kernel and saw that it computes the dot product of a 2-dimensional vector in a 6-dimensional space without explicitly visiting that space.



Next, let's talk about the RBF kernel, another insanely powerful kernel, which is also the default kernel in a support vector classifier class implemented by sklearn:

[Home](#) > API Reference > `sklearn.svm > SVC`

SVC RBF is default kernel

```
class sklearn.svm.SVC(*, C=1.0, kernel='rbf', degree=3, gamma='scale',
coef0=0.0, shrinking=True, probability=False, tol=0.001, cache_size=200,
class_weight=None, verbose=False, max_iter=-1, decision_function_shape='ovr',
break_ties=False, random_state=None)
[source]
```

C-Support Vector Classification.

To begin, the mathematical expression of the RBF kernel is depicted below (and consider that we have just a 1-dimensional feature vector):

$$k(x, y) = e^{(-\gamma \|x - y\|^2)}$$

You may remember from high school mathematics that the exponential function is defined as follows:

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} + \dots$$

Expanding the square term in the RBF kernel expression, we get:

$$\begin{aligned} k(x, y) &= e^{(-\gamma \|x - y\|^2)} \\ &= e^{(-\gamma(x^2 + y^2 - 2xy))} \end{aligned}$$

Distributing the gamma term and expanding the exponential term using the exponent rule, we get:

$$\begin{aligned} k(x, y) &= e^{(-\gamma(x^2 + y^2 - 2xy))} \\ &= e^{-\gamma(x^2 + y^2) + 2\gamma xy} \\ &= e^{-\gamma(x^2 + y^2)} \cdot e^{2\gamma xy} \end{aligned}$$

Next, we apply the exponential expansion to the last term and get the following:

$$\begin{aligned} k(x, y) &= e^{-\gamma(x^2 + y^2)} \cdot e^{2\gamma xy} \\ &= e^{-\gamma(x^2 + y^2)} \cdot \boxed{(1 + 2\gamma xy + \frac{(2\gamma xy)^2}{2!} + \frac{(2\gamma xy)^3}{3!} + \dots)} \end{aligned}$$

Notice closely that the exponential expansion above can be rewritten as the dot product between the following two vectors:

$$(1 + 2\gamma xy + \frac{(2\gamma xy)^2}{2!} + \frac{(2\gamma xy)^3}{3!} + \dots)$$

Vector #1 ↓ Vector #2

$(1, \sqrt{2\gamma} x, \sqrt{\frac{(2\gamma)^2}{2!}} x^2, \sqrt{\frac{(2\gamma)^3}{3!}} x^3, \dots)$
only includes x

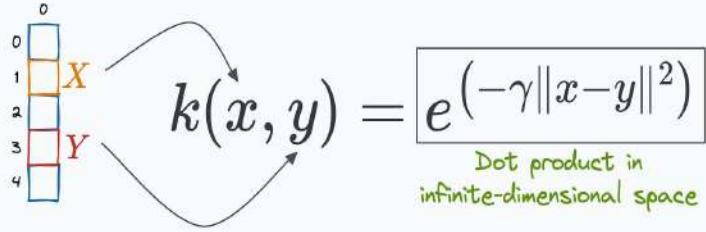
$(1, \sqrt{2\gamma} y, \sqrt{\frac{(2\gamma)^2}{2!}} y^2, \sqrt{\frac{(2\gamma)^3}{3!}} y^3, \dots)$
only includes y

And there we get our projection function:

$$\phi(x) = (1, \sqrt{2\gamma} x, \sqrt{\frac{(2\gamma)^2}{2!}} x^2, \sqrt{\frac{(2\gamma)^3}{3!}} x^3, \dots)$$

It is evident that this function maps the 1-dimensional input to an infinite-dimensional feature space.

This shows that the RBF kernel function we chose earlier computes the dot product in an infinite-dimensional space without explicitly visiting that space.



This is why the RBF kernel is considered so powerful, allowing it to easily model highly complex decision boundaries.

Here, I want to remind you that even though the kernel is equivalent to the dot product between two infinite-dimensional vectors, we NEVER compute that dot product, so the computation complexity is never compromised.

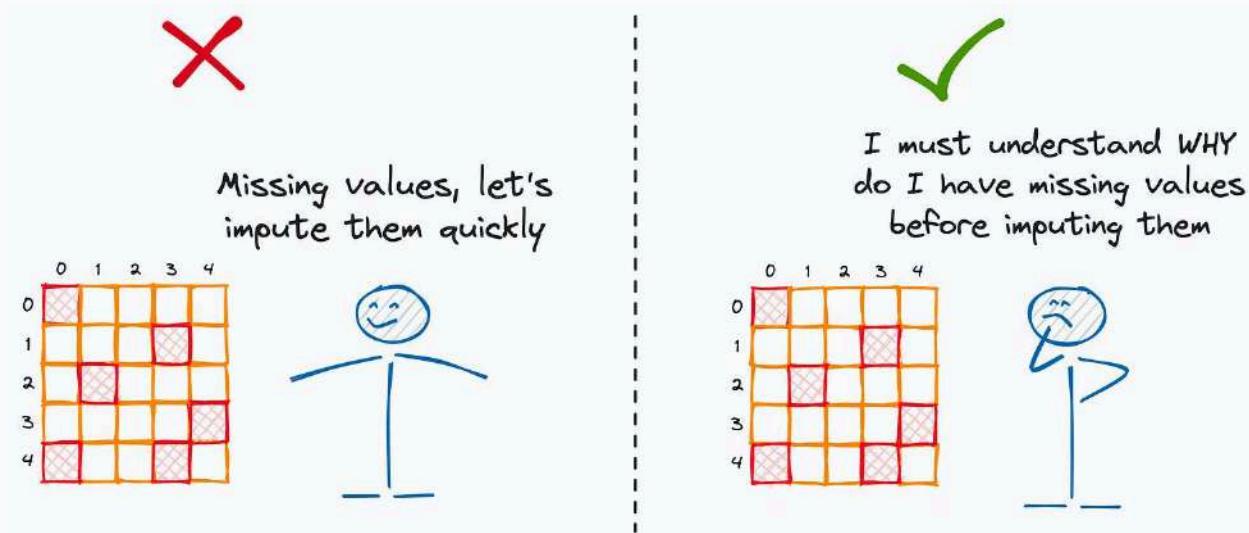
That is why the kernel trick is called a “trick.” In other words, it allows us to operate in high-dimensional spaces without explicitly computing the coordinates of the data in that space.

Missing Data

3 Types of Missing Values

After seeing missing values in a dataset, most people jump directly into imputing them. But as counterintuitive as it may sound, the first step towards imputing missing data should NEVER be imputation.

Instead, the first step must be to understand the reason behind data missingness. This is because the data imputation strategy largely depends on WHY data is missing.



More specifically, missingness could be of three types:

- Missing completely at random (MCAR).
- Missing at random (MAR).
- Missing not at random (MNAR).

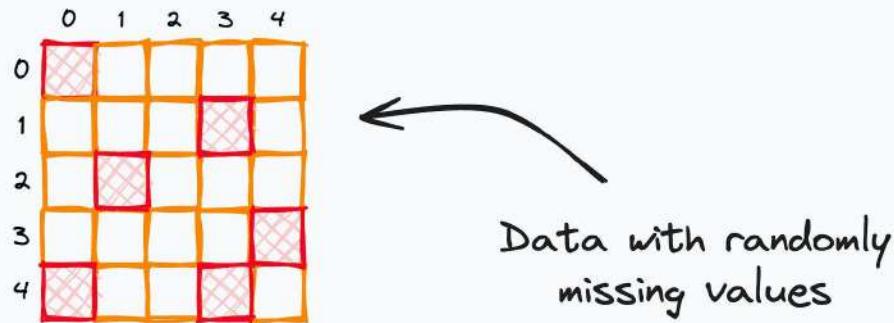
Only when we determine the reason can we proceed with the appropriate imputation techniques.

What are MCAR, MAR, and MNAR?

Let's understand them in this chapter in detail.

1) Missing completely at random (MCAR)

MCAR is a situation in which the data is genuinely missing at random and has no relation to any observed or unobserved variables. In other words, the missing data points follow no discernable pattern.



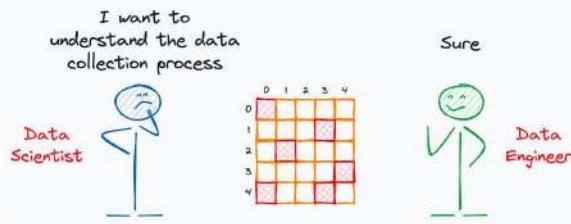
For instance, in survey responses with missing values, assuming MCAR would mean that some participants may have unintentionally skipped to answer any random question.

However, in my experience, MCAR has mostly been an unrealistic assumption for missingness, i.e., data is usually not MCAR for most real-world datasets with missing values. This is because, in real-world scenarios, the occurrence of missing data often tends to be influenced by some factor, either observed or unobserved.

Human behavior, survey administration, or any external events can be motivating factors for missing values to appear in a dataset. For instance, participants may selectively omit sensitive information, or certain groups may be more prone to non-response, causing missingness.

That is why assuming MCAR for missing data is not a sensible assumption unless you know the end-to-end data collection process and/or have domain expertise in.

Here, it becomes essential for data scientists to talk to data engineers and understand the data collection process. In fact, this is not just about MCAR but applicable to the other two situations as well, which we shall discuss shortly.



Understanding the data collection process will NEVER hurt.

That said, nothing stops us from assuming that missing values are MCAR if that appears to be a fair thing to do based on the analysis and/or input from the domain experts and/or after understanding the data collection mechanism. You can proceed with the simplest univariate imputation techniques.

2) Missing at random (MAR)

MAR is a situation in which the missingness of one feature can be explained by other observed features in the dataset.

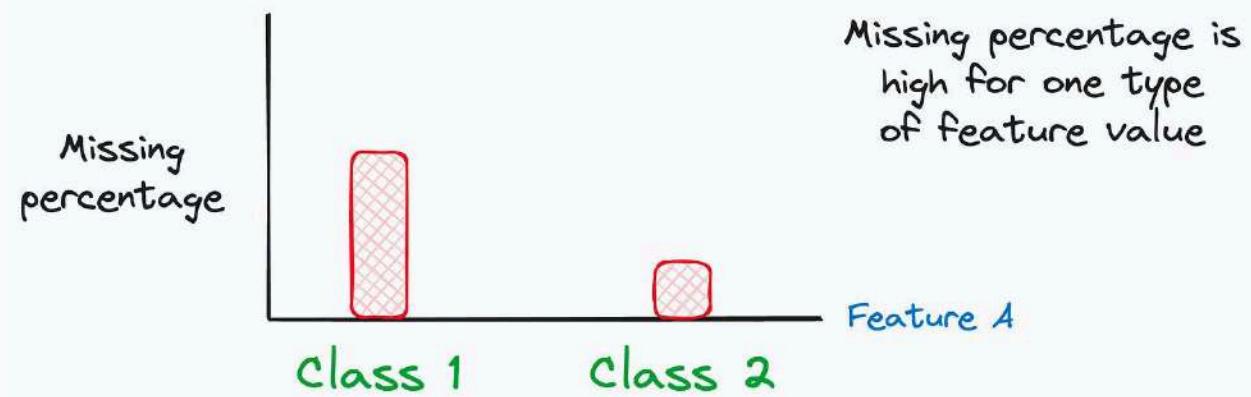


In contrast to MCAR, MAR is more practically observed.

In this case, the missingness can be accounted for through appropriate statistical methods with reasonable accuracy.

Thus, even though the data is missing, its occurrence can still be (somewhat) estimated based on the information available in the dataset.

A common way to determine MAR is by conditioning on another observed feature and noticing any increase in the probability of missingness.



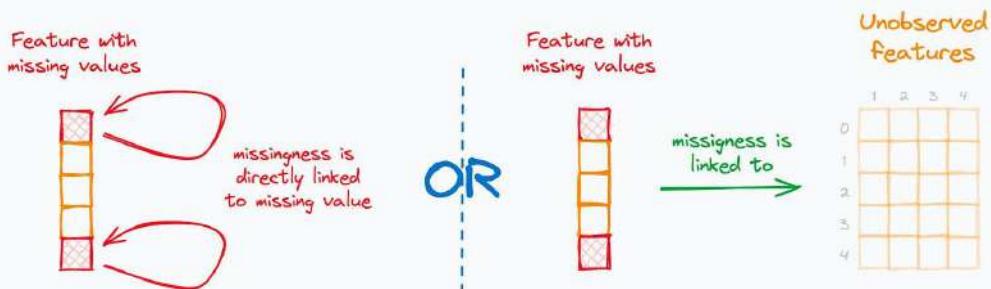
For instance, in an academic performance survey, students with higher grades might be less likely to disclose information about the number of hours they study (due to competition, maybe).

This is where techniques like kNN imputation, Miss Forest, etc., are quite effective. They use other observed features to impute the missing feature.

If you want to learn about these techniques, the chapter is precisely about this topic.

3) Missing not at random (MNAR)

MNAR is the most complicated situation of all three. In MNAR, missingness is either attributed to the missing value itself or the feature(s) that we didn't collect data for.



In other words, within MNAR, there is a definite pattern in missing variables. However, the pattern is unrelated to any observed feature(s). This is different from MCAR where there is no definite pattern.

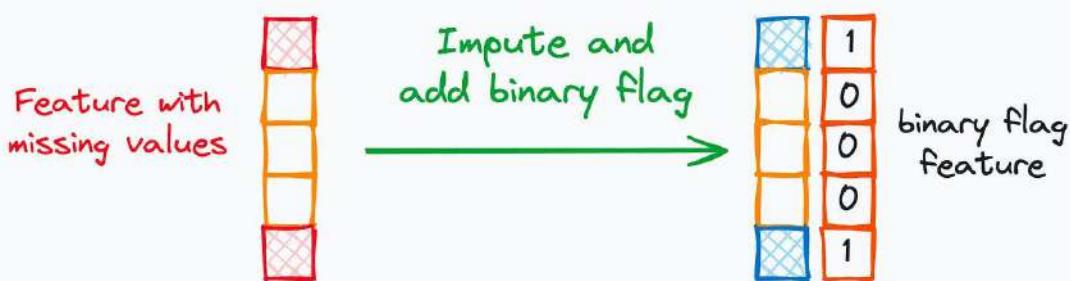
For instance, in a health survey, participants with very high stress levels might consciously choose not to disclose their stress level due to stigma, or fear of judgment. As a result, the missing data on stress level is not random; it is influenced by the stress level itself.

So, in a way, the higher the stress level, the less likely one will disclose it, and the more likely the value will be missing from the collected dataset. Thus, the missingness is directly dependent on the very variable that is missing in the first place. That's tricky, isn't it?

There's not much we can do to address this, except for collecting more data/features. What's more, domain expertise become extremely important to smartly tackle MNAR and improve the data collection process.

At times, I have also preferred proceeding with typical imputation techniques (used in MCAR and MAR) because further data collection can be infeasible in most cases. But as discussed above, there is a definite missingness pattern in MNAR, which can be important. But direct imputation will discard that information.

One way to retain that missingness pattern is by adding a binary feature, indicating whether the feature was imputed.

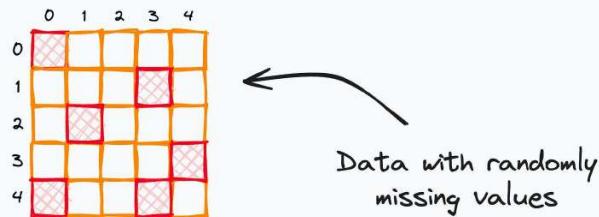


This way, the ML algorithm can still access and learn the missing data patterns.

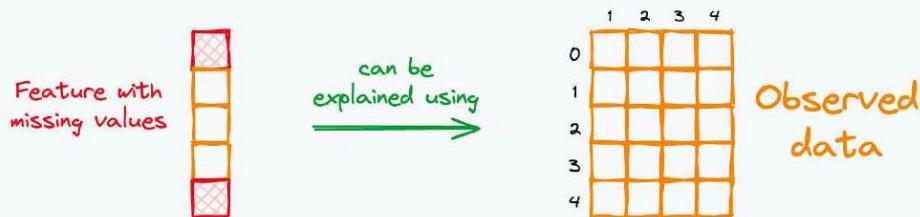
MissForest and kNN Imputation

The imputation strategy for missing data largely depends on the type of missingness, which can be of three types:

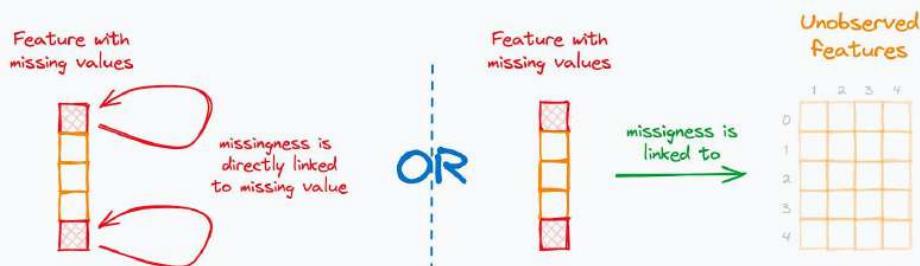
- Missing completely at random (MCAR): Data is genuinely missing at random and has no relation to any observed or unobserved variables.



- Missing at random (MAR): The missingness of one feature can be explained by other observed features in the dataset.



- Missing not at random (MNAR): Missingness is either attributed to the missing value itself or the feature(s) that we didn't collect data for.



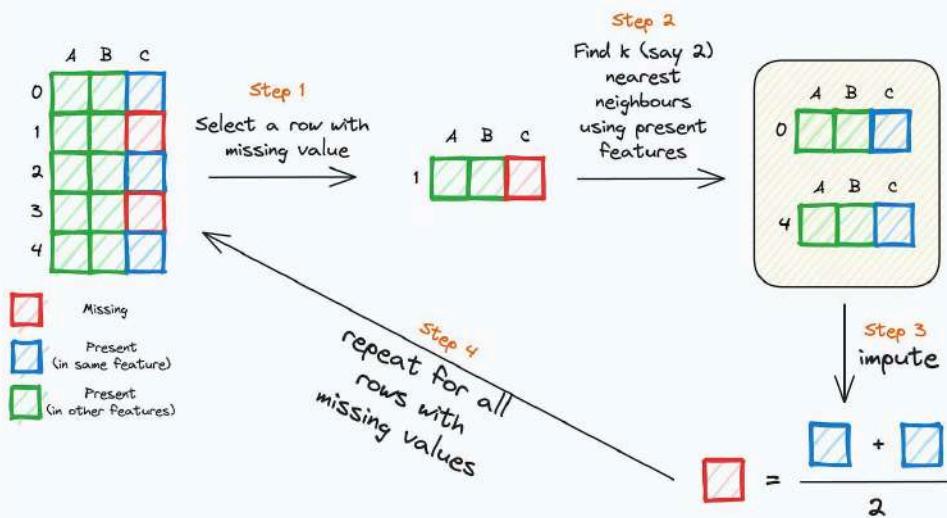
I have observed missing at random (MAR) appearing relatively much more in practical situations than the other two, so in this chapter, I want to share two imputation strategies I typically prefer to use in such a case.

#1) kNN imputation

As the name suggests, it imputes missing values using the k-nearest neighbors algorithm.

More specifically, missing features are imputed by running a kNN on non-missing feature values.

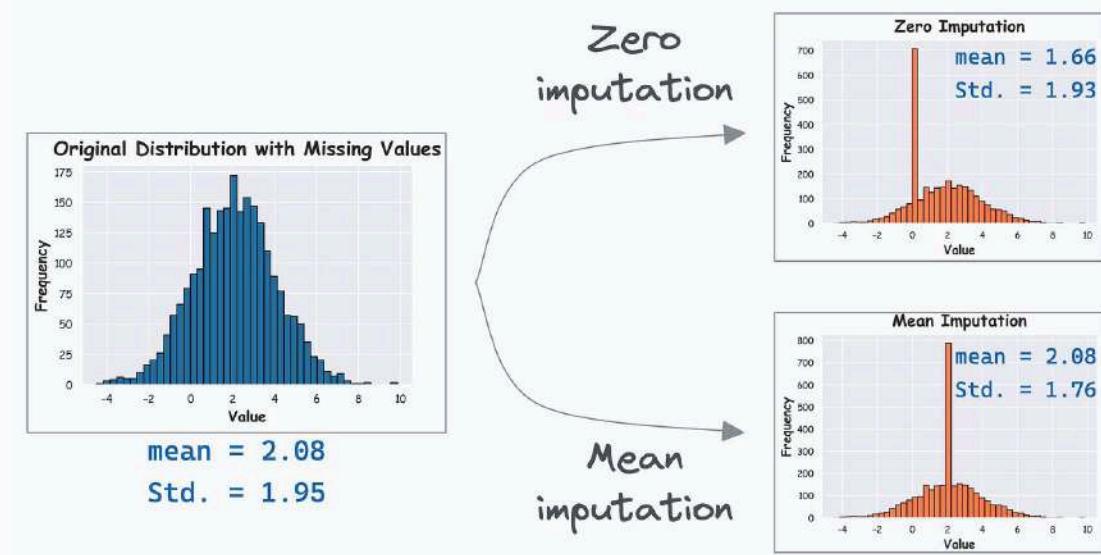
The following image depicts how it works:



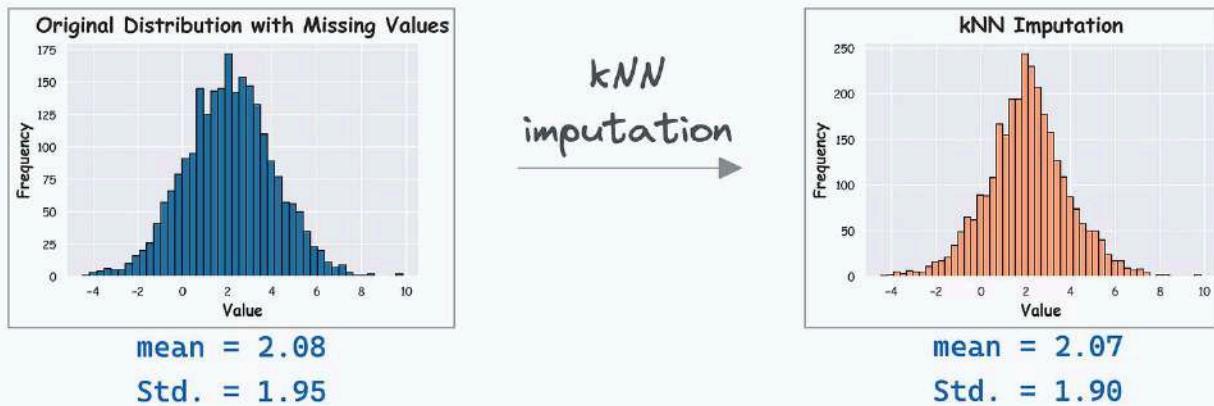
- Step 1: Select a row (r) with a missing value.
- Step 2: Find its k nearest neighbors using the non-missing feature values.
- Step 3: Impute the missing feature of the row (r) using the corresponding non-missing values of k nearest neighbor rows.
- Step 4: Repeat for all rows with missing values.

Its effectiveness over mean/zero imputation is evident from the demo below.

On the left, we have a feature distribution with missing values. Assuming we have already validated that the data was missing at random (MAR), using mean/zero alters the summary statistics and distribution, as depicted below.



However, as depicted below, kNN imputer appears to be more reliable, and it preserves the summary statistics:

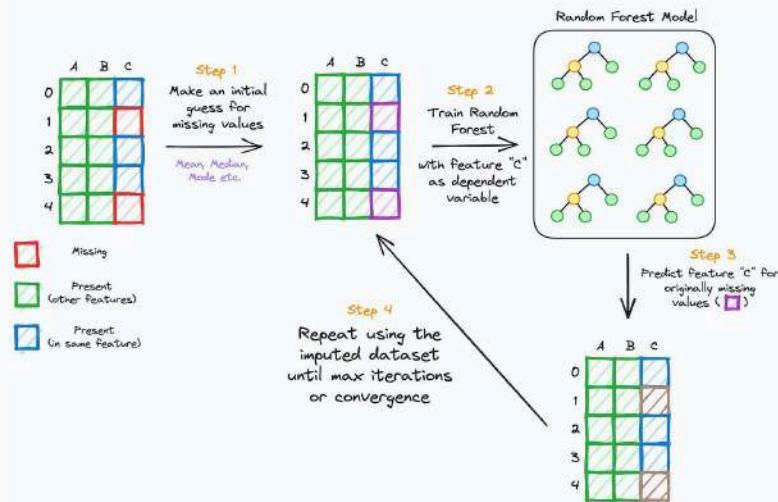


#2) MissForest

Some major issues with kNN imputation are:

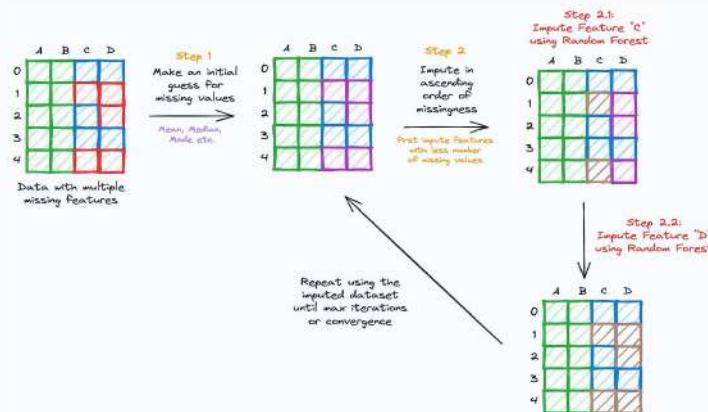
1. It has a high run-time for imputation — especially for high-dimensional datasets.
2. It raises issues with distance calculation in the case of categorical non-missing features.
3. It requires feature scaling, etc.

MissForest imputer is another reliable choice for missing value imputation when data is missing at random (MAR). As the name suggests, it imputes missing values using the Random Forest algorithm.



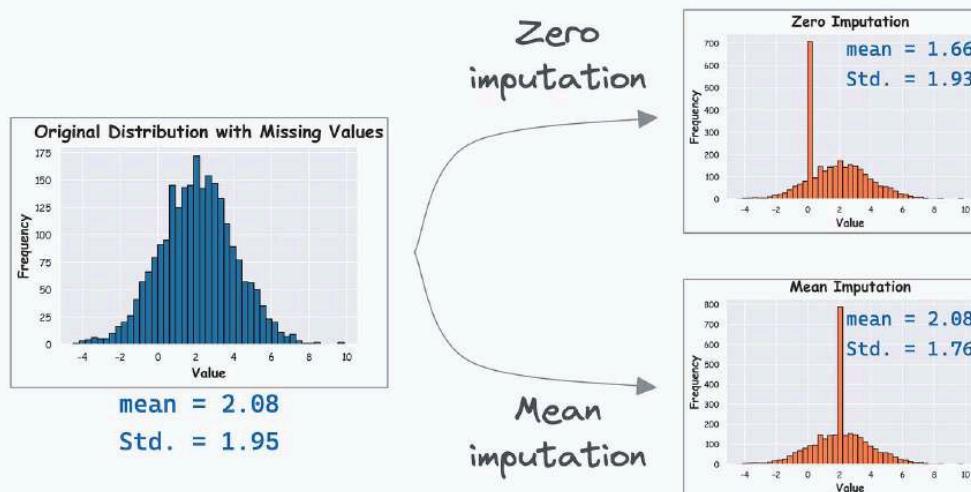
1. To begin, impute the missing feature with a random guess — Mean, Median, etc.
2. Model the missing feature using Random Forest.
3. Impute ONLY originally missing values using Random Forest's prediction.
4. Back to Step 2. Use the imputed dataset from Step 3 to train the next Random Forest model.
5. Repeat until convergence (or max iterations).

In case of multiple missing features, the idea (somewhat) stays the same:

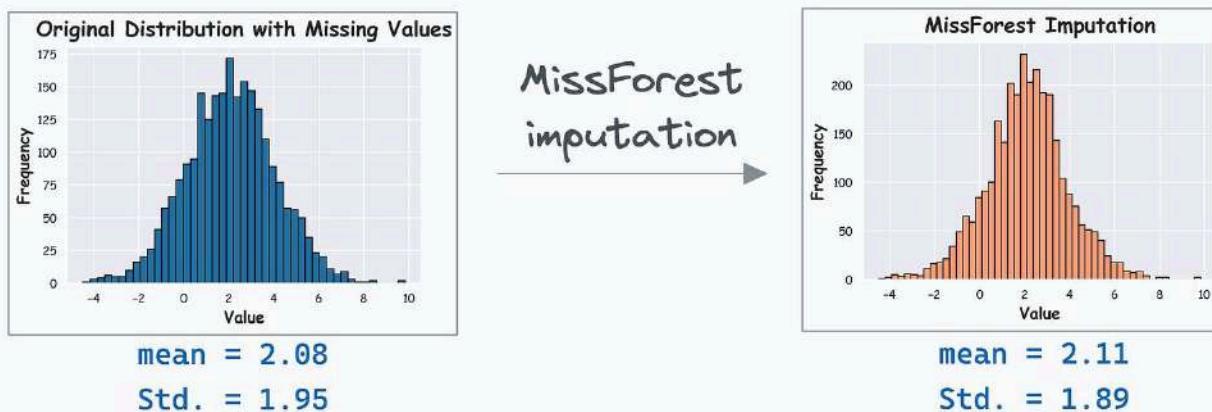


Impute features sequentially in increasing order missingness — features with fewer missing values are imputed first.

Like before, on the left, we have a feature distribution with missing values. Again, assuming we have already validated that the data was missing at random (MAR), using mean/zero alters the summary statistics and distribution:



However, as depicted below, MissForest appears more reliable and preserves the summary statistics.

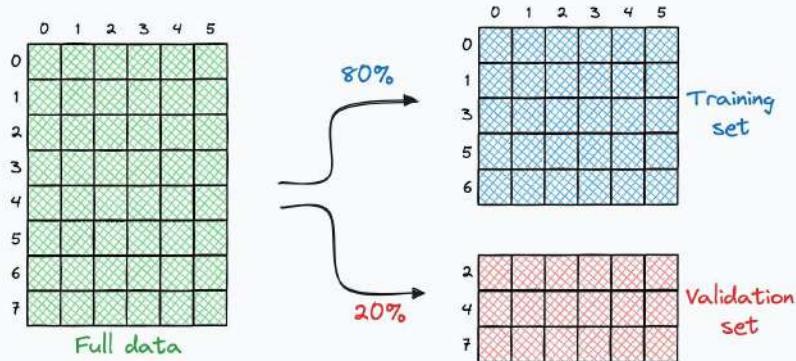


kNN imputation and MissForest Jupyter notebook: <https://bit.ly/3RW3i8k>.

Pitfalls and Misconceptions

When is Random Splitting Fatal for ML Models?

One thing we hear in almost all introductory ML lessons is to split the given data RANDOMLY into train and validation sets.



Note: Ignore the test set for now.

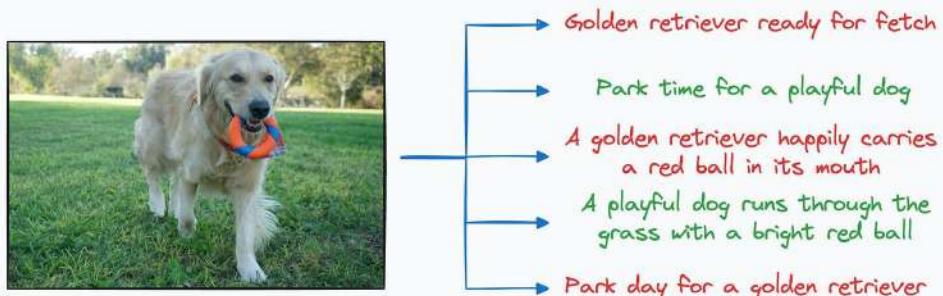
Random splitting makes sense because it ensures that the data is divided without any bias.

However, I have come across many situations where random splitting is fatal for model building. Yet, many people don't realize it.

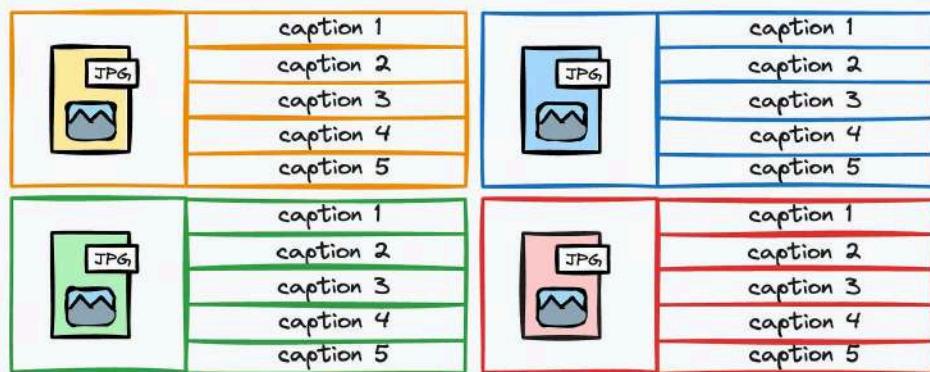
And I am not talking about temporal datasets here.

Scenario

Consider you are building a model that generates captions for images.

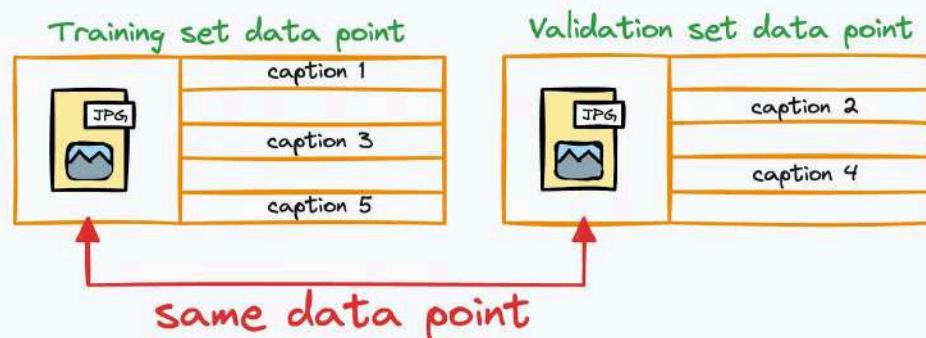


Due to the inherent nature of language, every image can have many different captions:



Now, you might realize what would happen if we randomly split this dataset into train and validation sets.

During the random split, the same data point (image) will be available in the train and validation sets.



This is a typical example of data leakage, which results in high overfitting!

This type of leakage is also known as group leakage.

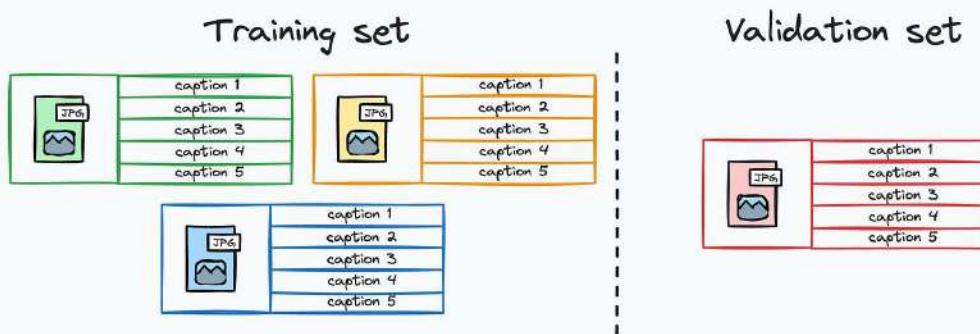
Solution

From the above discussion, it is clear that random splitting is the cause of the problem.

Group shuffle split helps us solve this.

There are two steps:

- Group all training instances corresponding to one image (or features that may result in leakage, any other grouping criteria, etc.).
- After grouping, the whole group must be randomly sent to either the training set or the validation set.



This will prevent the group leakage we witnessed earlier and prevent overfitting.

One thing to note here is that in the above example, all features in the dataset, i.e., the image pixels, contributed to the grouping criteria.

But more generally speaking, there could only be a subset of features that must be grouped together for data splitting.

For instance, consider a dataset containing medical imaging data. Each sample consists of multiple images (e.g., different views of the same patient's body part), and the model is intended to detect the severity of a disease.

In this case, it is crucial to group all images corresponding to the same patient together and then perform data splitting. Otherwise, it will result in data leakage and the model will not generalize well to new patients.

Demo

If you use Sklearn, the `GroupShuffleSplit` implements this idea.

`sklearn.model_selection.GroupShuffleSplit`

```
class sklearn.model_selection.GroupShuffleSplit(n_splits=5, *, test_size=None, train_size=None,
random_state=None)
```

[source]

Shuffle-Group(s)-Out cross-validation iterator.

Provides randomized train/test indices to split data according to a third-party provided group. This group information can be used to encode arbitrary domain specific stratifications of the samples as integers.

For instance the groups could be the year of collection of the samples and thus allow for cross-validation against time-based splits.

Consider we have the following dataset:

	x1	x2	y	group
0	1	2	1	A
1	3	4	2	
2	5	6	3	B
3	7	8	4	
4	9	10	5	C

- x_1 and x_2 are the features.
- y is the target variable.
- group denotes the grouping criteria.

First, we import the `GroupShuffleSplit` from sklearn and instantiate the object:

Define group splitter

```
from sklearn.model_selection import GroupShuffleSplit
gss = GroupShuffleSplit(n_splits=1, test_size=0.2)
```

The `split()` method of this object lets us perform group splitting:



This returns a generator, and we can unpack it to get the following output:

A screenshot of a Jupyter Notebook cell. The code cell contains a `for` loop:

```
for train_index, test_index in gss.split(X, y, groups):
    X_train, X_test = X[train_index], X[test_index]
    y_train, y_test = y[train_index], y[test_index]
```

Below the code cell are two tables illustrating the data splitting. The left table shows the training set (`x_train` and `y_train`) and the group (`group`) for each row index (0, 1, 4). The right table shows the validation/test set (`x_test` and `y_test`) and the group (`group`) for each row index (2, 3).

	<code>x_train</code>	<code>y_train</code>	<code>group</code>
0	1 2	1	A
1	3 4	2	
4	9 10	5	

	<code>x_test</code>	<code>y_test</code>	<code>group</code>
2	5 6	3	B
3	7 8	4	

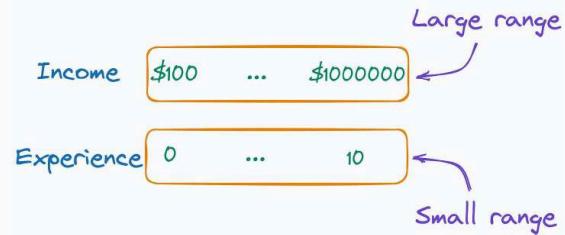
As demonstrated above:

- The data points in groups “A” and “C” are together in the training set.

The data points in group “B” are together in the validation/test set.

Feature Scaling is NOT Always Necessary

Feature scaling is commonly used to improve the performance and stability of ML models. This is because it scales the data to a standard range. This prevents a specific feature from having a strong influence on the model's output.



For instance, in the image above, the scale of income could massively impact the overall prediction. Scaling both features to the same range can mitigate this and improve the model's performance. I am sure you already know this, so we won't get into more detail here. However, have you ever wondered the following:

Is feature scaling always necessary when our dataset's features have a diverse range?

While feature scaling is often crucial, we often overlook whether it is even needed or not. This is because many ML algorithms are unaffected by scale.

This is evident from the image, which depicts the test accuracy of some classification algorithms with and without feature scaling.

Logistic regression (trained using SGD), SVM Classifier, MLP, and kNN do better with feature scaling.

Decision trees, Random forests, Naive bayes, and Gradient boosting are unaffected by scale.

Feature Scaling is NOT Always Necessary

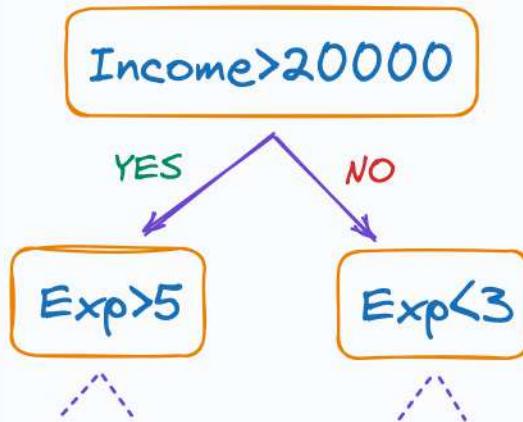


Test Set Classification Performance			
Algorithm	Without Feature Scaling	With Feature Scaling	Scaling Required?
Logistic Regression	0.53	0.70	YES
Support Vector Classifier	0.72	0.94	YES
MLP Classifier	0.73	0.89	YES
kNN Classifier	0.66	0.93	YES
Decision Tree	0.83	0.83	NO
Random Forest	0.91	0.91	NO
Gradient Boosting	0.86	0.86	NO
Naive Bayes	0.75	0.75	NO

Better performance with feature scaling

Same performance irrespective of feature scaling

To understand better, consider a decision tree. It splits the data based on thresholds determined solely by the feature values, regardless of their scale.



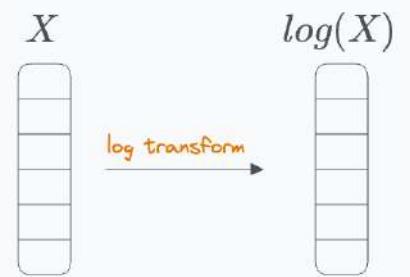
As a result, its performance is unaffected by the scale. This makes intuitive sense as well.

Thus, the takeaway is that when we do feature scaling, it's important to understand not just the nature of our data but also the algorithm we intend to use. We may never need feature scaling if the algorithm is insensitive to the scale of the data.

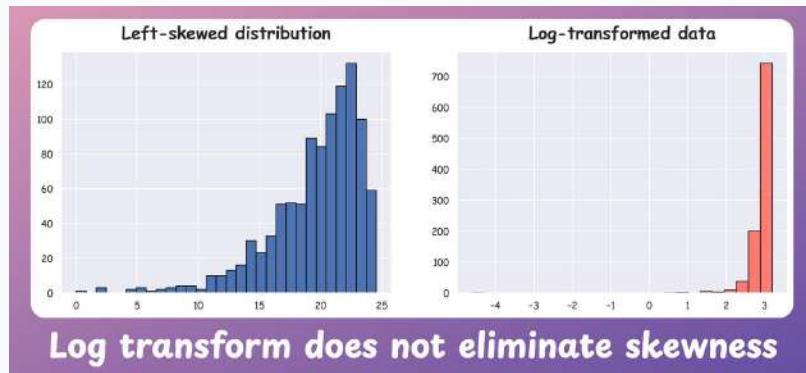
A Misconception About Log Transform

Log transform is commonly used to eliminate skewness in data. Yet, it is not always the ideal solution for eliminating skewness. It is important to note that log transform:

- DOES NOT eliminate left-skewness.
- Only works for right-skewness, that too when the values are small and positive.

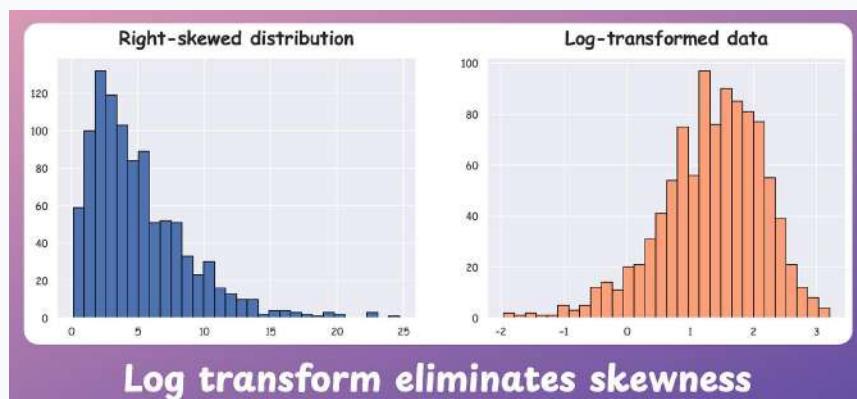


To understand better, consider the left-skewed distribution and its log transform below:



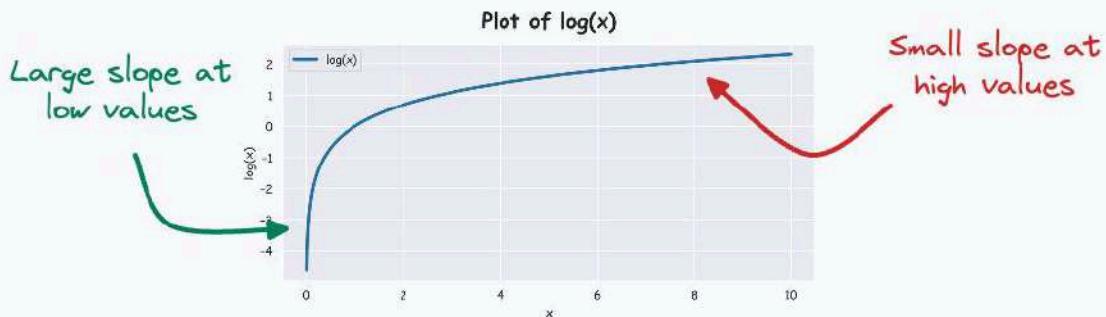
It is clear that log transform did not eliminate skewness.

However, now consider a right-skewed distribution:

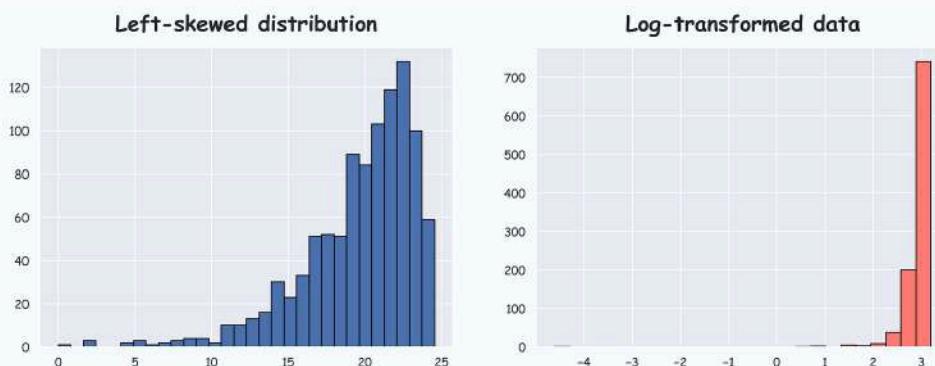


This time, it eliminates the skewness.

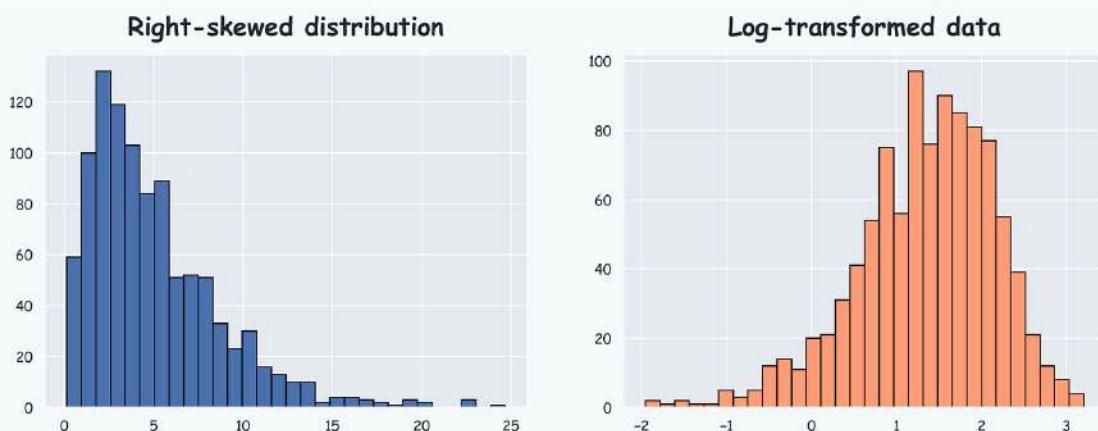
This happens because the log function grows faster for lower values. Thus, it stretches out the lower values more than the higher values.



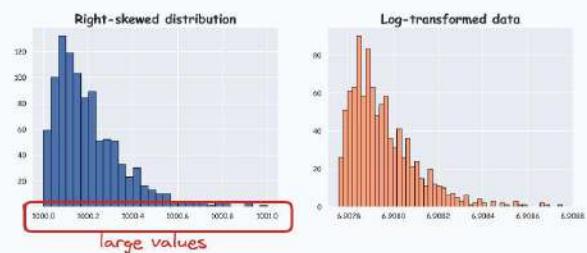
More specifically, in the case of left-skewness, the tail exists to the left, which gets stretched more than the majority of the probability mass that exists to the right. Thus, skewness isn't affected much.



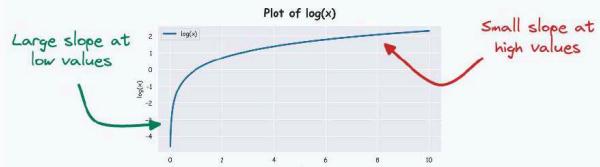
However, in the case of right-skewness, the majority of probability mass exists to the left, which gets stretched out more.



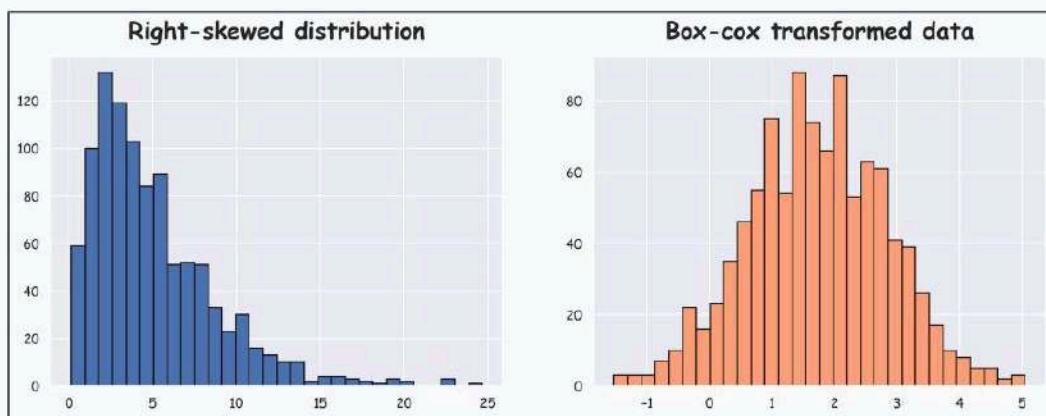
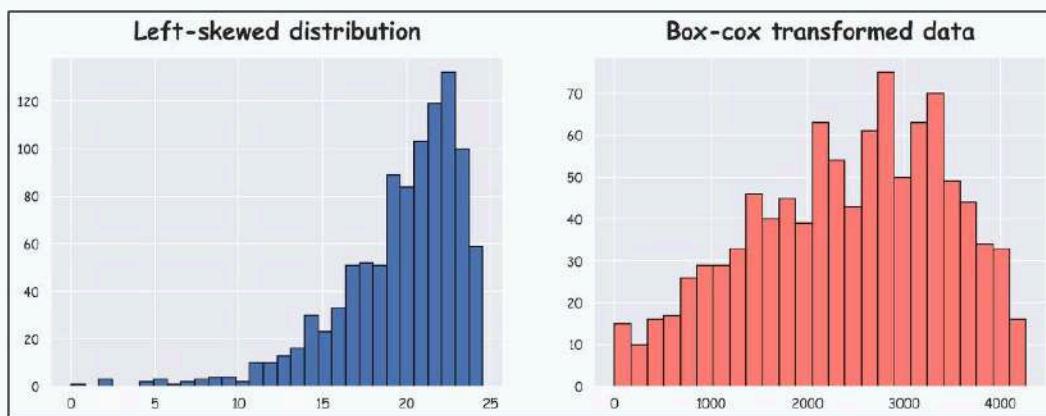
That said, it is important to note that even if we have a right-skewed distribution, log transform will not be effective if the values are large.



This happens because the stretch effect of the log function diminishes at large values:



I have often found the box-cox transform to be quite effective at eliminating both right-skewness and left-skewness, as depicted below:



The True Purpose of Feature Scaling and Standardization

Feature scaling and standardization are two common ways to alter a feature's range. For instance:

- MinMaxScaler changes the range of a feature to [0,1]:

$$x' = \frac{x - x_{min}}{x_{max} - x_{min}}$$

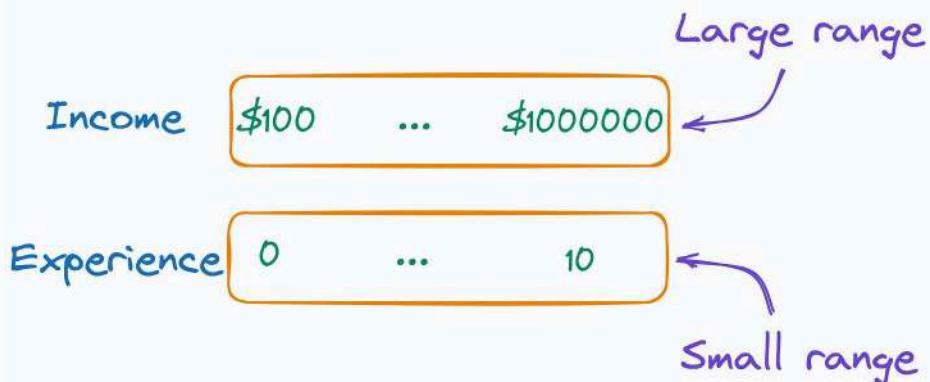
- Standardization makes a feature's mean zero and standard deviation one:

$$z = \frac{x - \mu}{\sigma}$$

As you may already know, these operations are necessary because:

- They prevent a specific feature from strongly influencing the model's output.
- They ensure that the model is more robust to wide variations in the data.

For instance, in the image below, the scale of "income" could massively impact the overall prediction.



Scaling (or standardizing) the data can mitigate this and improve the model's performance. We also learned this in the [feature scaling](#) chapter:

Feature Scaling is NOT Always Necessary



blog.DailyDoseofDS.com

Test Set Classification Performance			
Algorithm	Without Feature Scaling	With Feature Scaling	Scaling Required?
Logistic Regression	0.53	0.70	YES
Support Vector Classifier	0.72	0.94	YES
MLP Classifier	0.73	0.89	YES
kNN Classifier	0.66	0.93	YES
Decision Tree	0.83	0.83	NO
Random Forest	0.91	0.91	NO
Gradient Boosting	0.86	0.86	NO
Naive Bayes	0.75	0.75	NO

Better performance with feature scaling

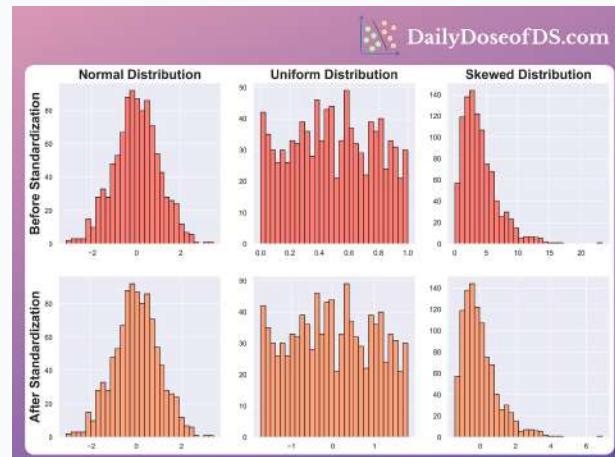
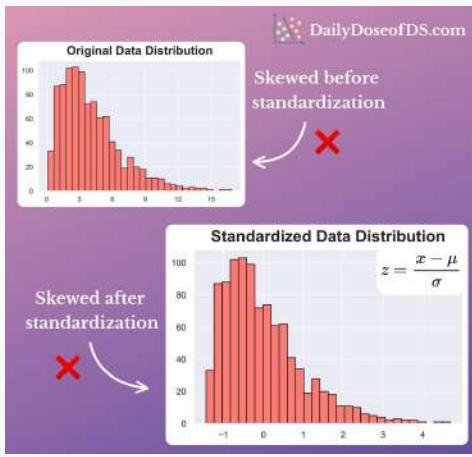
Same performance irrespective of feature scaling

As depicted above, feature scaling is necessary for the better performance of many ML models. So while the importance of feature scaling and standardization is pretty clear and well-known, I have seen many people misinterpreting them as techniques to eliminate skewness.

But contrary to this common belief, feature scaling and standardization NEVER change the underlying distribution. Instead, they just alter the range of values. Thus, after scaling (or standardization):

- Normal distribution → stays Normal
- Uniform distribution → stays Uniform
- Skewed distribution → stays Skewed
- and so on...

We can also verify this from the two illustrations below:



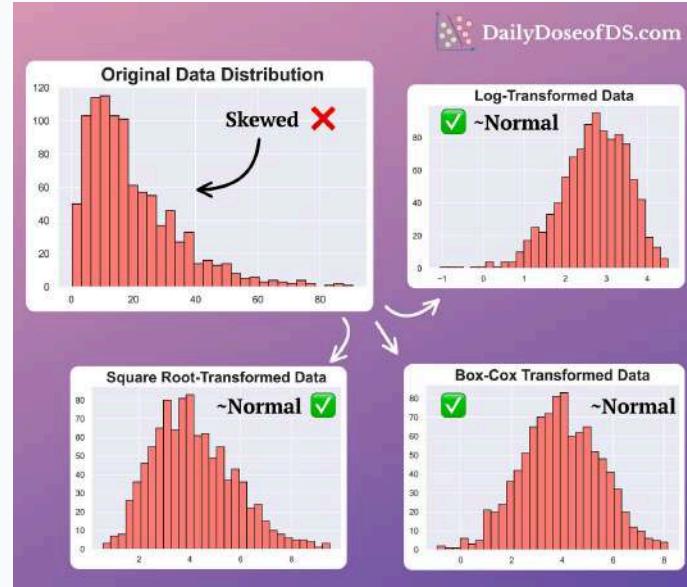
It is clear that scaling and standardization have no effect on the underlying distribution. Thus, always remember that if you intend to eliminate skewness, scaling/standardization will never help. Try feature transformations instead. There are many of them, but the most commonly used transformations are:

- Log transform
- Sqrt transform
- Box-cox transform

Their effectiveness is evident from the image below:

As depicted in this image, applying these operations transforms the skewed data into a (somewhat) normally distributed variable.

Before I conclude, please note that while log transform is commonly used to eliminate data skewness, it is not always the ideal solution.



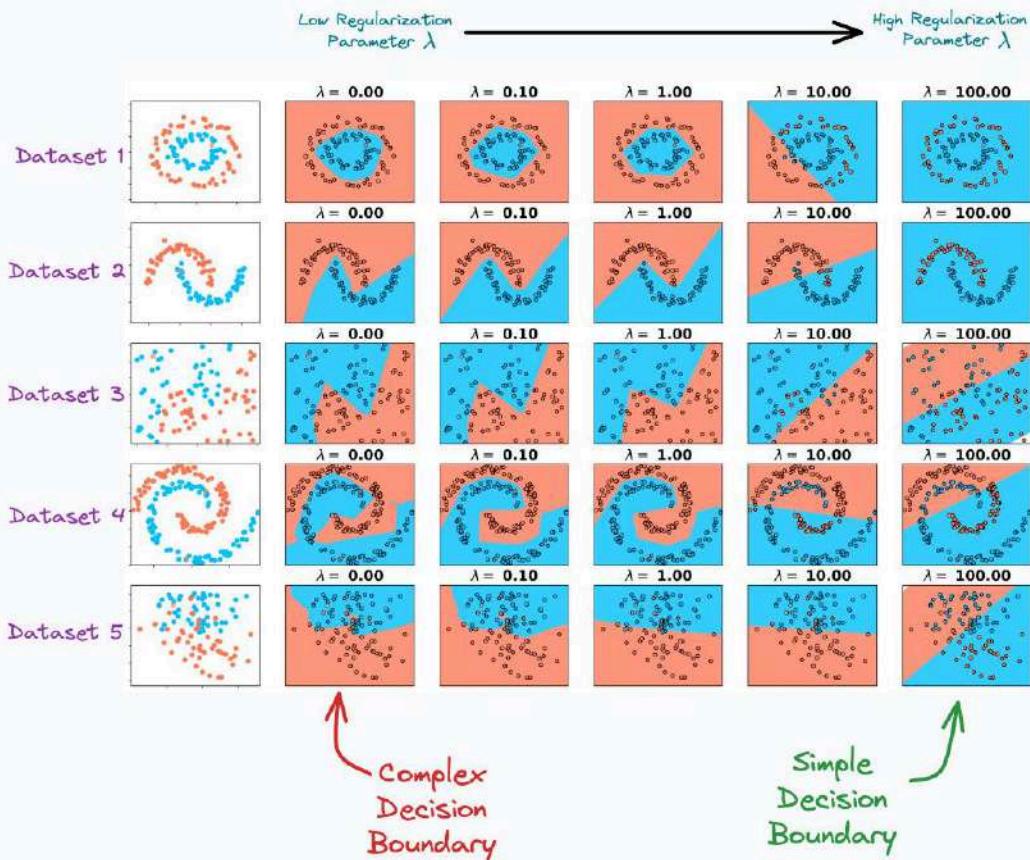
L2 Regularization is Not Just Used for Regularization

Almost every tutorial/course/blog mentioning L2 regularization I have seen talks about just one thing:

L2 regularization is a machine learning technique that avoids overfitting by introducing a penalty term into the model's loss function based on the squares of the model's parameters.

In classification tasks, for instance, increasing the effect of regularization will produce simpler decision boundaries.

Of course, the above statements are indeed correct, and I am not denying that. In fact, we can also verify this from the diagram below:

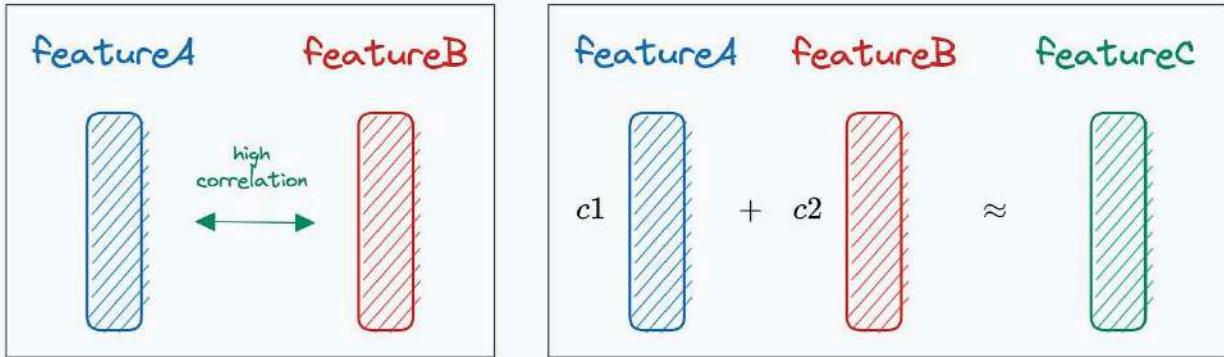


In the image above, as we move to the right, the regularization parameter increases, and the model creates a simpler decision boundary on all 5 datasets.

Coming back to the topic...

However, what disappoints me the most is that most resources don't point out that L2 regularization is a great remedy for multicollinearity.

Multicollinearity arises when two (or more) features are highly correlated OR two (or more) features can predict another feature:



When we use L2 regularization in linear regression, the algorithm is also called Ridge regression.

Ridge regression →
$$\sum_{i=1}^N (y_i - \theta^T X)^2 + (\theta_1^2 + \theta_2^2)$$

But how does L2 regularization eliminate multicollinearity?

In this chapter, let's build a demonstrative intuition into this topic, which will also explain why “ridge regression” is called “ridge regression.”

Dummy dataset

For demonstration purposes, consider this dummy dataset of two features:

```

import numpy as np

# data features
featureA = np.random.rand(100)
featureB = 2*featureA + 0.5*np.random.randn(100)

# output variable (a linear combination of features)
output = 2*featureA + 3*featureB + np.random.randn(100)

```

As shown above, we have intentionally made featureB highly correlated with featureA. This gives us a dummy dataset to work with.

Going ahead, we shall be ignoring any intercept term for simplicity.

Linear regression without L2 penalty

During regression modeling, the goal is to determine those specific parameters (θ_1, θ_2), which minimizes the residual sum of squares (RSS):

Find parameters θ that minimize this quantity $\rightarrow \sum_{i=1}^N (y_i - \theta^T X)^2$

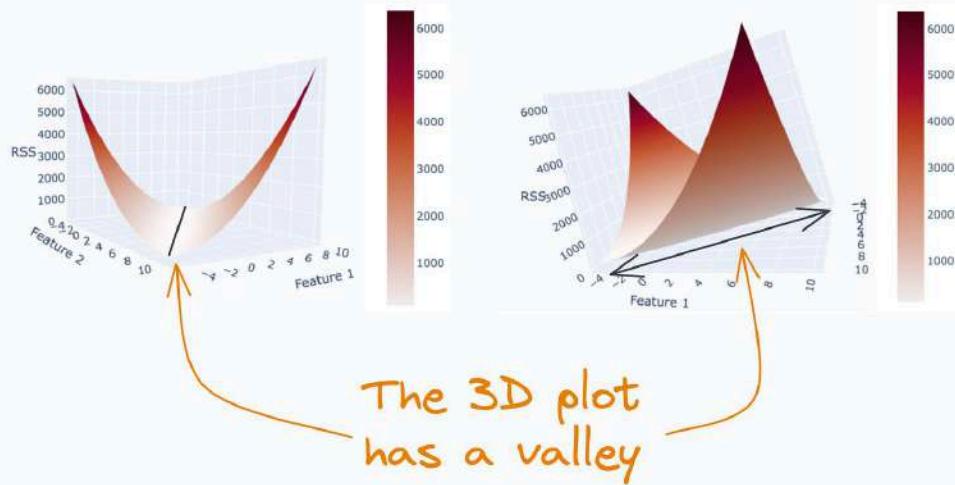
residual sum of squares

So how about we do the following:

- We shall plot the RSS value for many different combinations of (θ_1, θ_2) parameters. This will create a 3D plot:
 - x-axis $\rightarrow \theta_1$
 - y-axis $\rightarrow \theta_2$
 - z-axis \rightarrow RSS value
- Next, we shall visually assess this plot to locate those specific parameters (θ_1, θ_2) that minimize the RSS value.

Let's do this.

Without the L2 penalty, we get the following plot (it's the same plot but viewed from different angles):



Did you notice something?

The 3D plot has a valley. There are multiple combinations of parameter values (θ_1 , θ_2) for which RSS is minimum.

Thus, obtaining a single value for the parameters (θ_1 , θ_2) that minimize the RSS is impossible.

Linear regression with L2 penalty

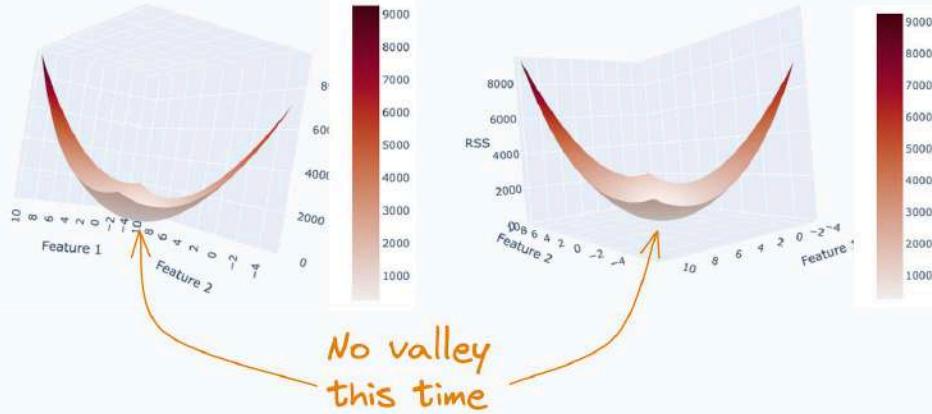
When using an L2 penalty, the goal is to minimize the following:

Find parameters θ that minimize this quantity \rightarrow

$$\sum_{i=1}^N (y_i - \theta^T X)^2 + (\theta_1^2 + \theta_2^2)$$

residual sum of squares L2 penalty

Creating the same plot again as we did above, we get the following:



Did you notice something different this time?

As depicted above, using L2 regularization removes the valley we saw earlier and provides a global minima to the RSS error.

Now, obtaining a single value for the parameters (θ_1, θ_2) that minimizes the RSS is possible.

Out of nowhere, L2 regularization helped us eliminate multicollinearity.

Why the name “ridge regression”?

In fact, this is where “ridge regression” also gets its name from — it eliminates the ridge in the likelihood function when the L2 penalty is used.

This is
valley



This is
ridge

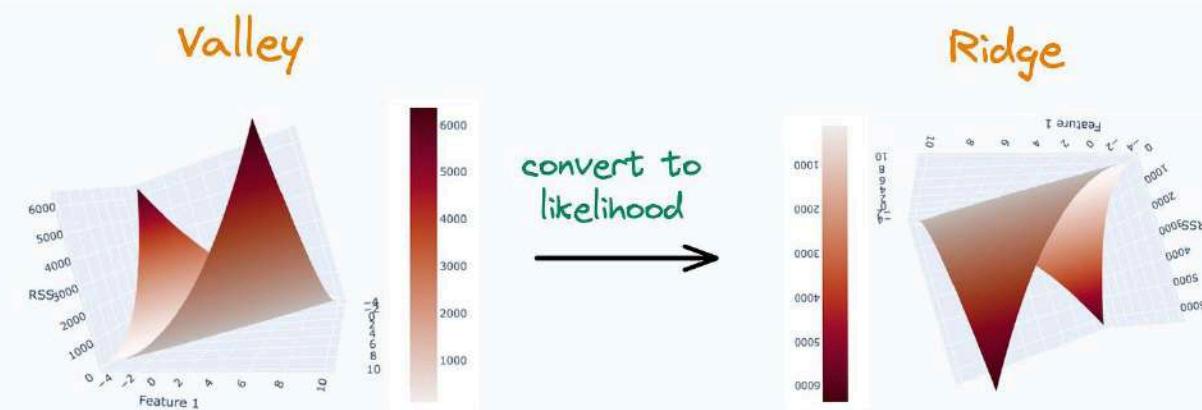


Of course, in the demonstrations we discussed earlier, we noticed a valley, not a ridge.

However, in that case, we were considering the residual sum of error — something which is minimized to obtain the optimal parameters.

Thus, the error function will obviously result in a valley.

If we were to use likelihood instead — something which is maximized, it would (somewhat) invert the graph upside down and result in a ridge instead:



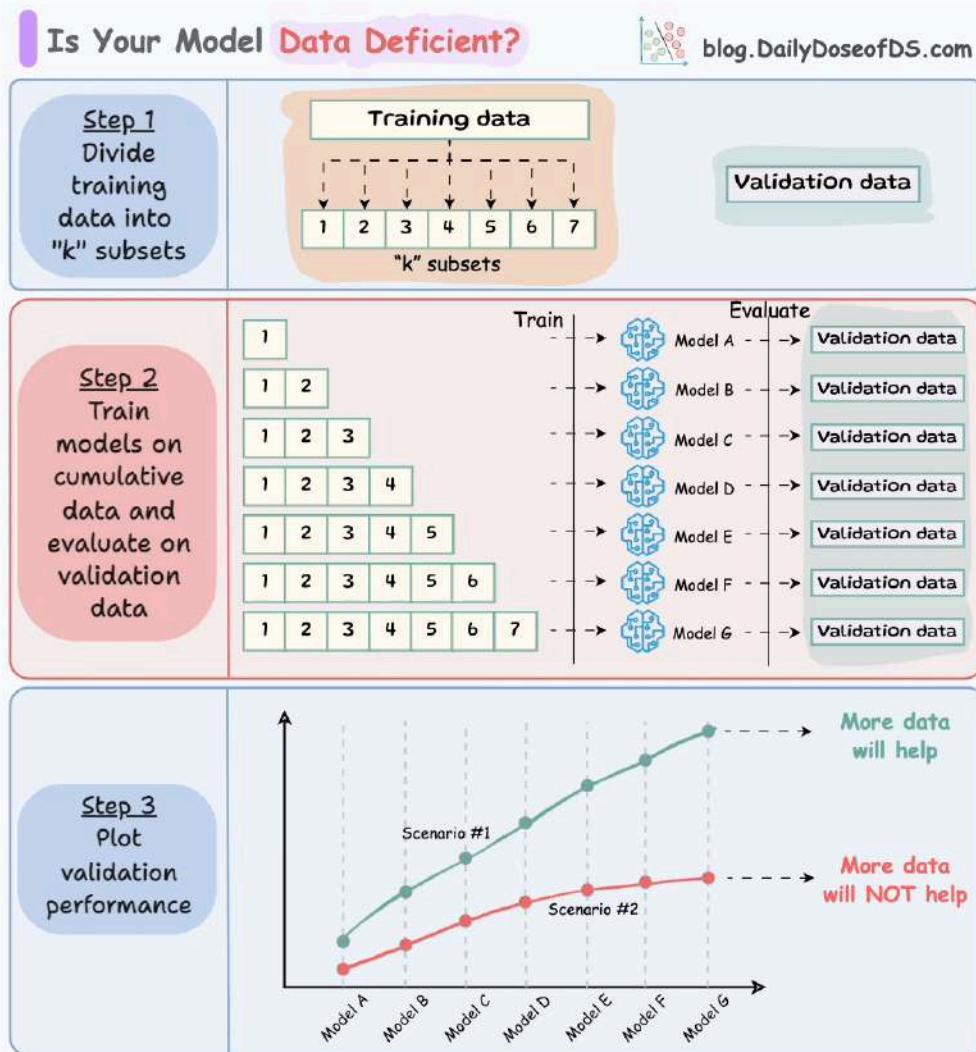
Apparently, while naming the algorithm, the likelihood function was considered.

And that is why it was named “ridge regression.”

When I first learned about this some years back, I literally had no idea that such deep thought went into naming ridge regression.

Miscellaneous

Is Your Model Data Deficient?

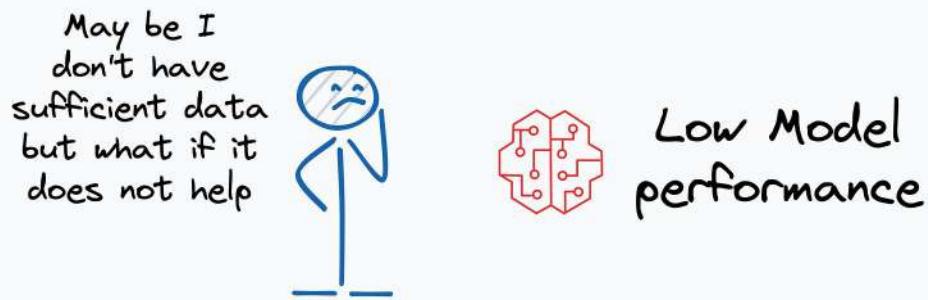


During model development, many people find themselves in situations where, no matter how much they try, the model performance barely improves:

- Feature engineering is giving marginal improvement.
- Trying different models does not produce satisfactory results either.
- and more...

This is usually (not always) an indicator that the model is data deficient.

In other words, we don't have enough data to work with.



However, gathering new data can be a time-consuming and tedious process. So before venturing into that direction, it would be good to get some insights about whether new data will help. Here's a trick I have often used to determine this.

Data subsetting and Model building

Let's say this is your full training and validation set:

Full training dataset

Val. set

Divide the training dataset into “k” equal parts. The validation set remains as is.

Full training dataset

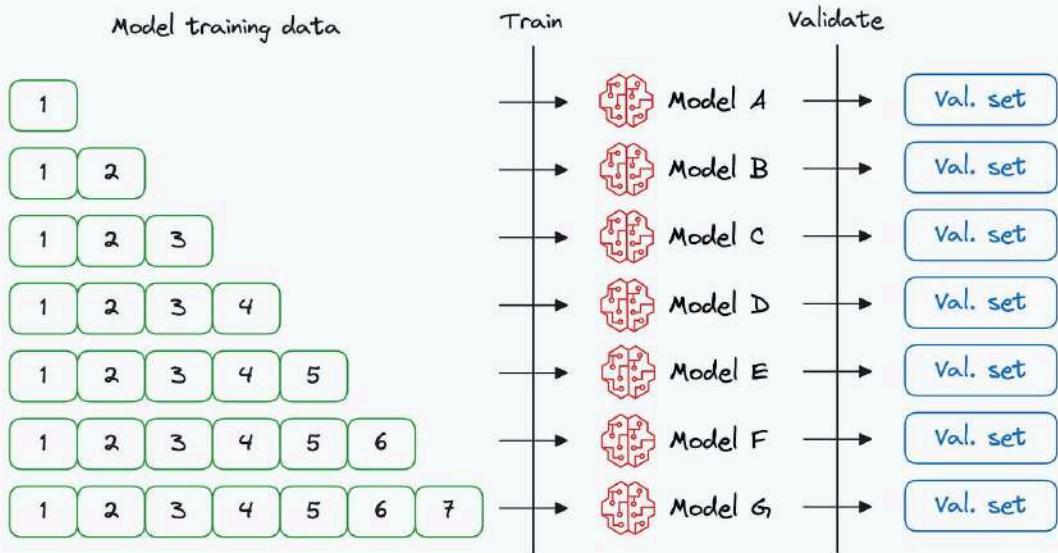
Val. set

“k” subsets



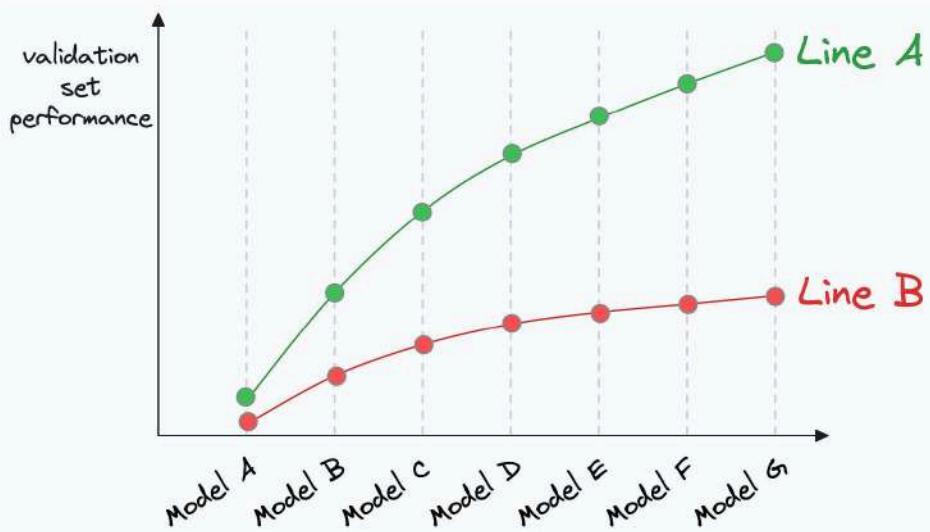
“k” does not have to be super large. Any number between 7 to 12 is fine depending on how much data you have. If there's plenty of data, setting a low value in this range is recommended (you will understand why shortly).

Next, train models cumulatively on the above subsets and measure their performance on the validation set:



- Train a model on the first subset and evaluate the validation set.
- Train a model on the first two subsets and evaluate the validation set.
- Train a model on the first three subsets and evaluate the validation set.
- And so on...

Plotting the validation performance of these models (in order of increasing training data) is likely to produce two types of lines:



- Line A conveys that adding more data is likely to increase model performance.
- Line B conveys that the model performance has already saturated. Adding more data will most likely not result in any considerable gains.

Now, you might also understand why I mentioned this above: “If there’s plenty of data, setting a low value in this range is recommended.”

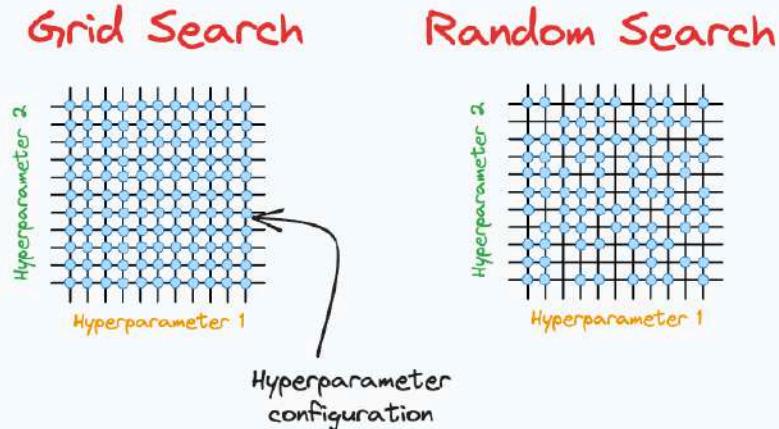
- Because we train multiple models, setting a high value of “k” means more subsets, which in turn means more models.

This way, you can determine whether the model is data deficient and whether gathering data will be helpful or not.

Bayesian Optimization

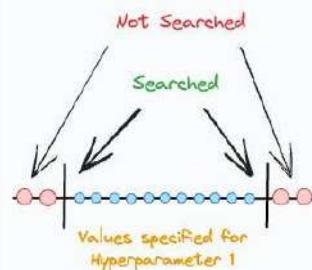
Hyperparameter tuning is a tedious and time-consuming task in training ML models.

Typically, we use two common approaches for this – Grid search and Random search.



But they have many limitations. For instance:

- Grid search performs an exhaustive search over all combinations. This is computationally expensive.
- Grid search and random search are restricted to the specified hyperparameter range. Yet, the ideal hyperparameter may exist outside that range.
- They can ONLY perform discrete searches, even if the hyperparameter is continuous.



To this end, Bayesian Optimization is a highly underappreciated yet immensely powerful approach for tuning hyperparameters.

It uses Bayesian statistics to estimate the distribution of the best hyperparameters. Here's how it differs from Grid search and Random Search:

Both Grid search and Random Search evaluate every hyperparameter configuration independently. Thus, they iteratively explore all hyperparameter configurations to find the most optimal one.

However, Bayesian Optimization takes informed steps based on the results of the previous hyperparameter configurations. This lets it confidently discard non-optimal configurations. Consequently, the model converges to an optimal set of hyperparameters much faster. The efficacy of Bayesian Optimization is evident from the image.

Limitation of Grid Search and Random Search

 blog.DailyDoseofDS.com

	Grid Search	Random Search	Bayesian Optimization	
Total Iterations	720	360	100	7x less iterations <input checked="" type="checkbox"/>
Total Run-time	242s	113s	48s	5x less run-time <input checked="" type="checkbox"/>
Best Trial Index	667	221	83	Optimal config. found earlier <input checked="" type="checkbox"/>
Best F1 score	0.93	0.93	0.93	Same score <input checked="" type="checkbox"/>

Bayesian optimization leads the model to the same F1 score but:

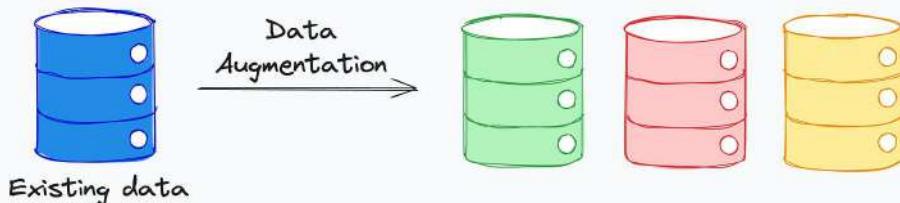
- it takes 7x fewer iterations
- it executes 5x faster
- it reaches the optimal configuration earlier

Pretty cool, isn't it?

Train and Test-time Data Augmentation

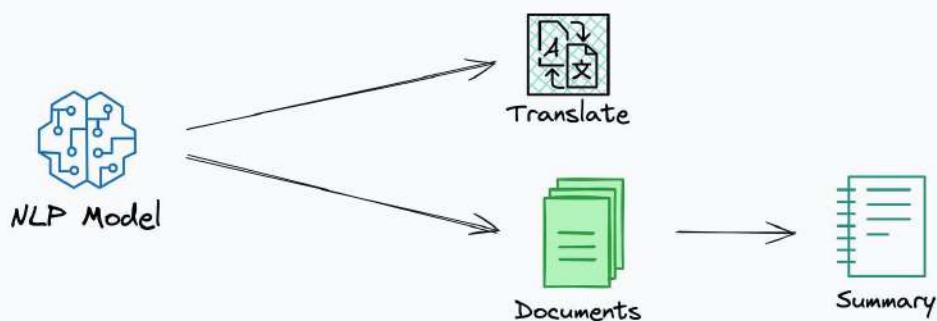
Data augmentation strategies are typically used during training time.

The idea is to use some clever techniques to create more data from existing data, which is especially useful when you don't have much data to begin with:



Let me give you an example.

These days, language-related ML models have become quite advanced and general-purpose. The same model can translate, summarize, identify speech tags (nouns, adjectives, etc.), and much more.



But earlier, models used to be task-specific (we have them now as well, but they are fewer than we used to have before).

- A dedicated model that would translate.
- A dedicated model that would summarize, etc.

In one particular use case, I was building a named entity recognition (NER) model, and the objective was to identify named entities.

An example is shown below:



I had minimal data — around 8-10k labeled sentences. The dataset was the CoNLL 2003 NER dataset if you know it.

Here's how I approached data augmentation in this case.

Observation: In NER, the factual correctness of the sentences does not matter.

Revisiting the above example, it would not have mattered if I had the following sentence in the training data:



The sentence is factually incorrect, of course, but that does not matter.

The only thing that matters to the model is that the output labels (named entity tags in this case) must be correct.

So using this observation, I created many more sentences by replacing the named entities in an existing sentence with other named entities in the whole dataset:



For such substitutions, I could have used named entities from outside. However, it was important to establish a fair comparison with other approaches.

This technique (along with a couple more architectural tweaks) resulted in state-of-the-art performance.

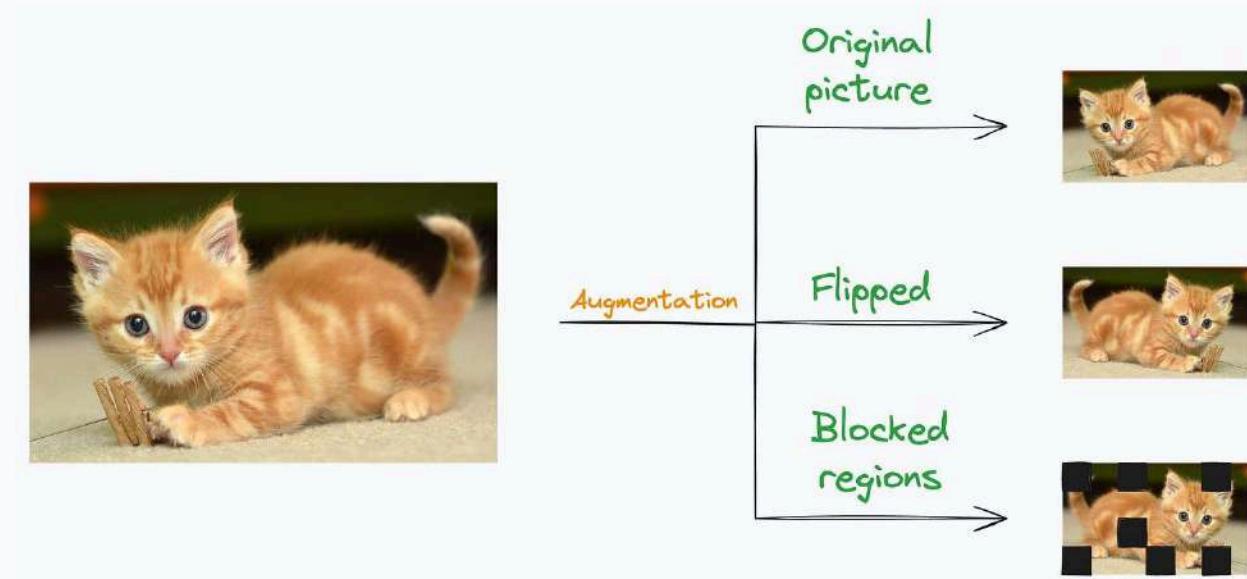
Moving on...

The above discussion was about training data augmentation.

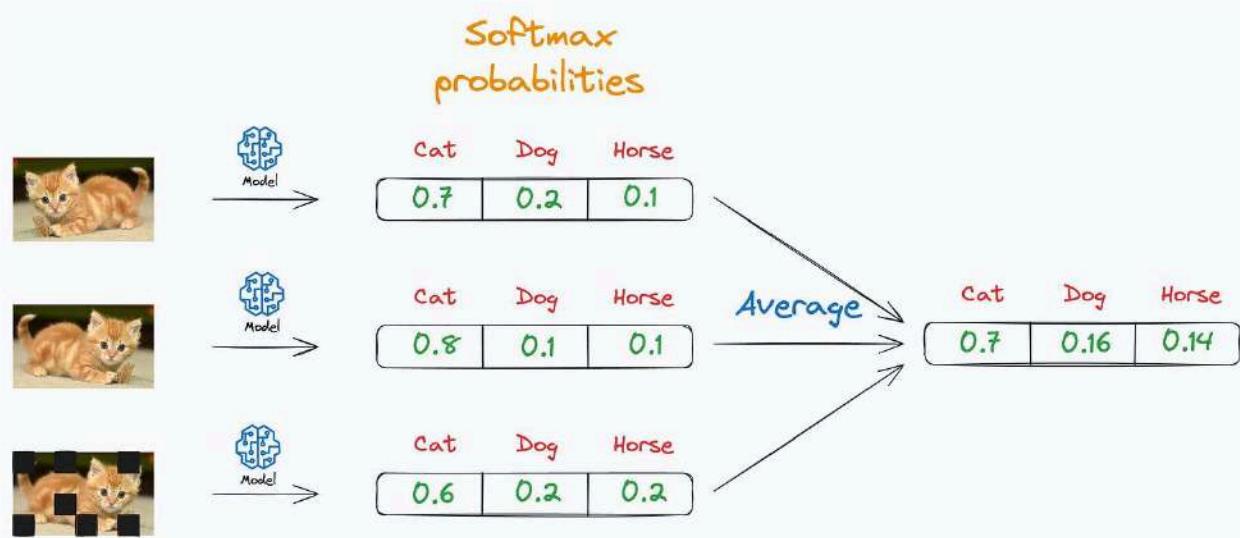
But there's also test-time augmentation.

Test Time Augmentation (TTA) is when we apply data augmentation during testing.

More specifically, instead of showing just one test example to the model, we show multiple versions of the test example by applying different operations.



The model makes probability predictions for every version of the test example, which are then averaged to generate the final prediction:



Essentially, TTA creates an ensemble of predictions by considering multiple augmented versions of the same input, which leads to a more robust final prediction.

In fact, in this paper (<https://arxiv.org/html/2402.06892v1>), the authors proved that the average model error with TTA never exceeds the average error of the original model, which is great.

As you may have guessed, the only catch is that it increases the inference time.

- Data augmentation takes some time.
- Generating multiple predictions increases the overall prediction run-time.

So, when a low inference time is important to you, think twice about TTA.

To summarize, if you can compromise a bit on inference time, TTA can be a powerful way to improve predictions from an existing model without having to engineer a better model.

Data Analysis

15 Pandas ↔ Polars ↔ SQL ↔ PySpark Translations

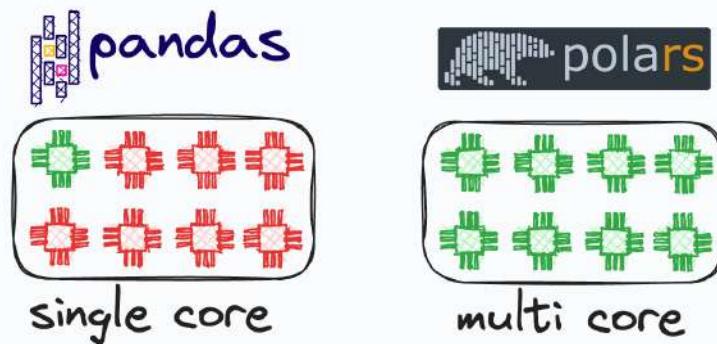
The following visual depicts the 15 most common tabular operations in Pandas and their corresponding translations in SQL, Polars, and PySpark.

Pandas ↔ Polars ↔ SQL ↔ PySpark				blog.DailyDoseofDS.com
Operation	Pandas	Polars	SQL	PySpark
Import	<code>import pandas as pd</code>	<code>import polars as pl</code>	-	<code>from pyspark.sql import SparkSession spark = SparkSession.builder.appName("ABCD")</code>
Read CSV	<code>df = pd.read_csv(file)</code>	<code>df = pl.read_csv(file)</code>	<code>LOAD DATA INFILE 'data.csv' INTO TABLE table FIELDS TERMINATED BY ',' LINES TERMINATED BY '\n' IGNORE 1 ROWS;</code>	<code>df = spark.read .csv("data.csv")</code>
Print first 10 (or k) rows	<code>df.head(10)</code>	<code>df.head(10)</code>	<code>SELECT * FROM table LIMIT 10;</code>	<code>df.show(10)</code>
Dimensions	<code>df.shape</code>	<code>df.shape</code>	<code>SELECT count(*) FROM table; SELECT count(*) FROM INFORMATION_SCHEMA.COLUMNS where TABLE_NAME = 'table';</code>	<code>df.count()</code> <code>len(df.columns)</code>
Datatype	<code>df.dtypes</code>	<code>df.dtypes</code>	<code>DESCRIBE table;</code>	<code>df.printSchema()</code>
Select column(s)	<code>df[["col1", "col2"]]</code>	<code>df[["col1", "col2"]]</code>	<code>SELECT column FROM table;</code>	<code>df.select("col1", "col2")</code>
Filter Data	<code>df[df.column > 10]</code>	<code>df[df.column > 10]</code> <code>df.filter(pl.col("column") > 10)</code>	<code>SELECT * FROM table where column>10;</code>	<code>df.filter(df["column"]>10)</code>
Sort	<code>df.sort_values("column")</code>	<code>df.sort("column")</code>	<code>SELECT * FROM table ORDER BY column;</code>	<code>df.orderBy("column")</code>
Fill NaN	<code>df.column.fillna(0)</code>	<code>df.column.fill_nan(0)</code>	<code>UPDATE table SET column=0 WHERE column IS NULL;</code>	<code>df.na.fill(0)</code>
Join	<code>pd.merge(df1, df2, on ="col", how="inner")</code>	<code>df1.join(df2, on="col", how="inner")</code>	<code>SELECT * FROM table1 JOIN table2 ON (table1.col = table2.col);</code>	<code>df1.join(df2, on="col", how="inner")</code>
Concatenate	<code>pd.concat((df1, df2))</code>	<code>pl.concat((df1, df2))</code>	<code>SELECT * FROM table1 UNION ALL table2;</code>	<code>df1.union(df2)</code>
Group	<code>df.groupby("column"). agg_col.mean()</code>	<code>df.groupby("column"). agg(pl.mean("agg_col"))</code>	<code>SELECT column, avg(agg_col) FROM table GROUP BY column;</code>	<code>df.groupBy("column"). agg(avg("agg_col"))</code>
Unique values	<code>df.column.unique()</code>	<code>df.column.unique()</code>	<code>SELECT DISTINCT column FROM table;</code>	<code>df.select("column"). distinct()</code>
Rename column	<code>df.rename(columns = {"old_name": "new_name"})</code>	<code>df.rename(mapping = {"old_name": "new_name"})</code>	<code>ALTER TABLE table RENAME COLUMN old_name TO new_name;</code>	<code>df.withColumnsRenamed({"old_name": "new_name"})</code>
Delete column	<code>df.drop(columns = ["column"])</code>	<code>df.drop(name = ["column"])</code>	<code>ALTER TABLE table DROP COLUMN column;</code>	<code>df.drop("col1", "col2")</code>

While the motivation for Pandas and SQL is clear and well-known, let me tell you why you should care about Polars and PySpark.

Why Polars?

Pandas has many limitations, which Polars addresses, such as:



- Pandas always adheres to single-core computation → Polars is multi-core.
- Pandas offers no lazy execution → Polars does.
- Pandas creates bulky DataFrames → Polars' DFs are lightweight.
- Pandas is slow on large datasets → Polars is remarkably efficient.

In fact, if we look at the run-time comparison on some common operations, it's clear that Polars is much more efficient than Pandas.

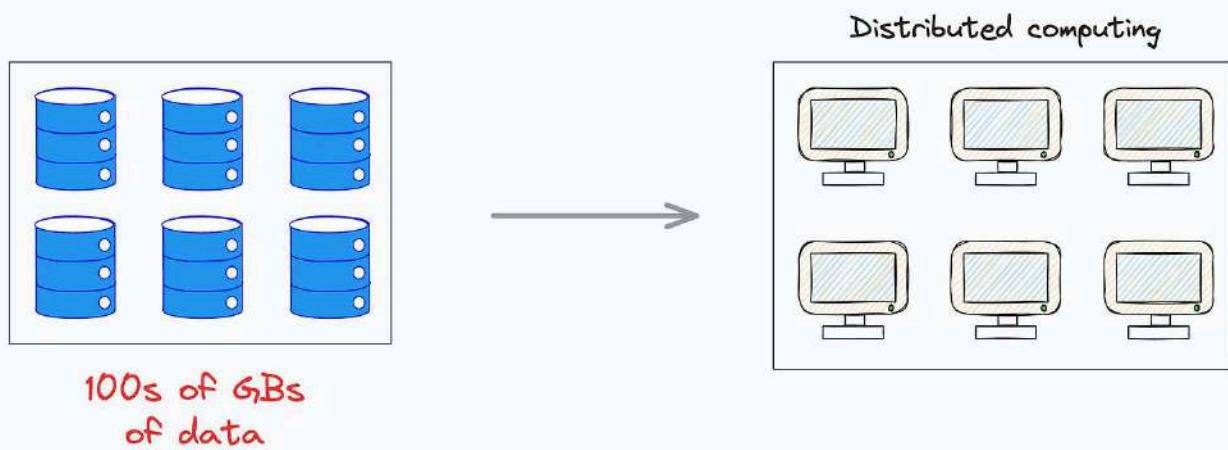
The table compares the performance of Pandas and Polars across several operations. The operations listed are Memory Usage, Read CSV, Save to CSV, Selecting Columns, Filtering, Grouping, and Sorting. For each operation, the time taken by Pandas is shown in orange and the time taken by Polars is shown in red. A note indicates that the number of records is 4M. The table shows that Polars consumes 3x less memory and is significantly faster for all operations compared to Pandas.

Method	Pandas	Polars	*Records: 4M
Memory Usage	1.7 GBs	0.63 GBs	Polars consumes 3x Less memory
Read CSV	4.52 s	0.72 s	Polars is 6x Faster
Save to CSV	15.3 s	3.25 s	Polars is 5x Faster
Selecting Columns	36.5 ms	2.2 μs	Polars is 16500x Faster
Filtering	130 ms	90 ms	Polars is 1.5x Faster
Grouping	217 ms	56 ms	Polars is 4x Faster
Sorting	1.1 s	1.04s	Similar Performance

Why Spark?

While tabular data space is mainly dominated by Pandas and Sklearn, one can hardly expect any benefit from them beyond some GBs of data due to their single-node processing.

A more practical solution is to use distributed computing instead — a framework that disperses the data across many small computers.



Spark is among the best technologies used to quickly and efficiently analyze, process, and train models on big datasets.

That is why most data science roles at big tech demand proficiency in Spark. It's that important

2 Alternatives to Pandas' Describe

Probably the first (or second) thing I do when I load any Pandas or Polars DataFrame is describe it, using the df.describe() method. However, I always find its output to be pretty naive and almost of no use. In other words, it hardly highlights any key information about the data.

	sepal_length	sepal_width	petal_length	petal_width
count	150.000000	150.000000	150.000000	150.000000
mean	5.843333	3.057333	3.758000	1.199333
std	0.828066	0.435866	1.765298	0.762238
min	4.300000	2.000000	1.000000	0.100000
25%	5.100000	2.800000	1.600000	0.300000
50%	5.800000	3.000000	4.350000	1.300000
75%	6.400000	3.300000	5.100000	1.800000
max	7.900000	4.400000	6.900000	2.500000

This is so hard to interpret



But some time back, I came across two pretty cool libraries that IMMENSELY supercharge this DataFrame summary. Since then, I don't think I have ever used the describe() method.

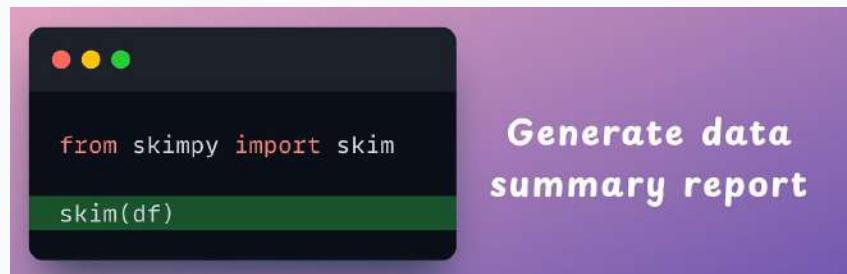
#1) Skimpy

It is a Jupyter-based tool that provides a standardized and comprehensive data summary.

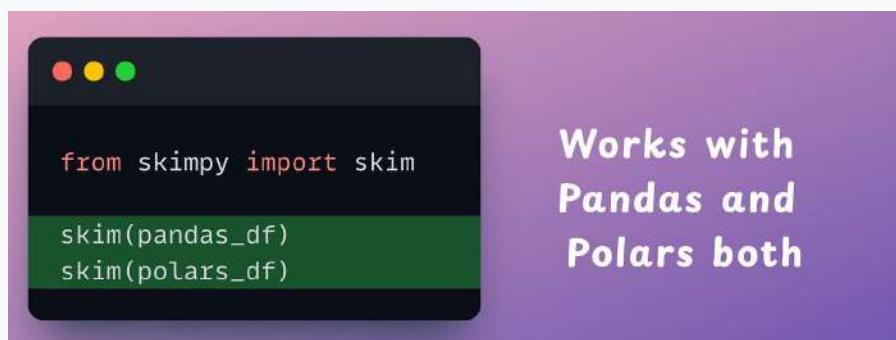
This includes data shape, column data types, column summary statistics, distribution charts, missing stats, etc.:

Data Summary		Data Types		skimpy summary		Categories			
dataframe	Values	Column Type	Count			Categorical Variables			
Number of rows	1000	float64	3						
Number of columns	10	category	2						
		datetime64	2						
		int64	1						
		bool	1						
		string	1						
number									
column_name	NA	NA %	mean	sd	p0	p25	p75	p100	hist
length	0	0	0.5	0.36	1.6e-06	0.13	0.86	1	
width	0	0	2	1.9	0.0021	0.6	3	14	
depth	0	0	10	3.2	2	8	12	20	
rnd	120	12	-0.02	1	-2.8	-0.74	0.66	3.7	
category									
column_name	NA	NA %		ordered		unique			
class	0	0		False		2			
location	1	0.1		False		5			
datetime									
column_name	NA	NA %	first	last		frequency			
date	0	0	2018-01-31	2101-04-30	M	None			
date_no_freq	3	0.3	1992-01-05	2023-03-04					
string									
column_name	NA	NA %	words per row		total words				
text	6	0.6	5.8		5800				
bool									
column_name	true		true rate		hist				
booly_col	520		0.52						
End									

What's more, the summary is grouped by datatypes for faster analysis. This is the code to use Skimpy:

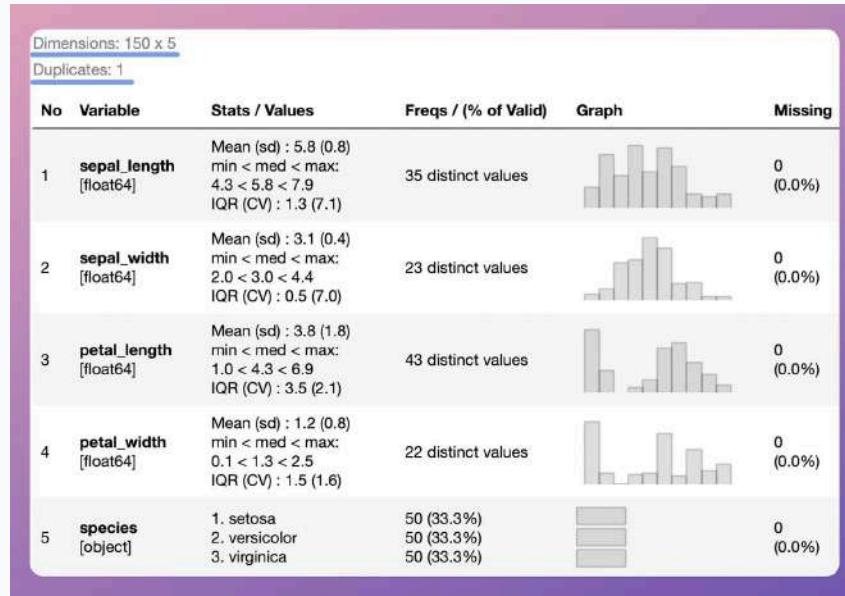


One thing I really love about Skimpy is that it works seamlessly with Polars, which I have started using more often than I use Pandas these days.

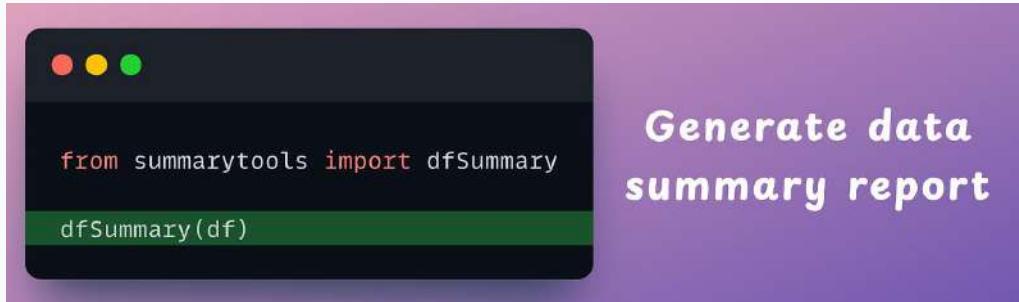


#2) SummaryTools

The second one is [SummaryTools](#), which does almost the exact same thing as Skimpy, i.e., it generates a standardized report:



This is the code to use SummaryTools:

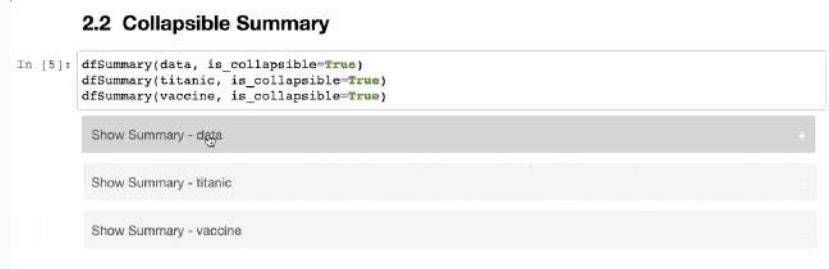


```
from summarytools import dfSummary
dfSummary(df)
```

Generate data summary report

Two pretty cool things about SummaryTools are that it can create:

1. A collapsible summary of the dataset, as illustrated below:



2.2 Collapsible Summary

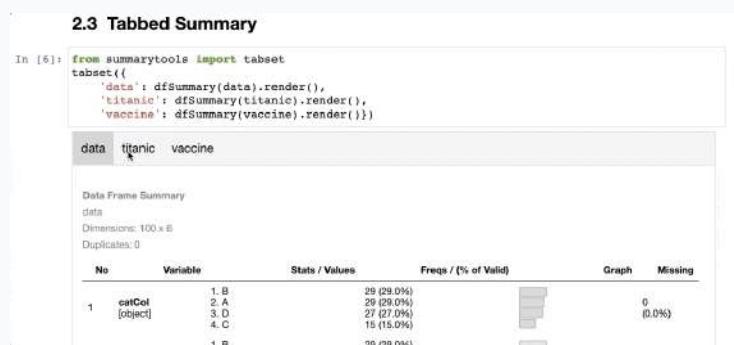
```
In [5]: dfSummary(data, is_collapsible=True)
dfSummary(titanic, is_collapsible=True)
dfSummary(vaccine, is_collapsible=True)
```

Show Summary - data

Show Summary - titanic

Show Summary - vaccine

2. A tabbed summary of the dataset, as shown below:



2.3 Tabbed Summary

```
In [6]: from summarytools import tabset
tabset({
    'data': dfSummary(data).render(),
    'titanic': dfSummary(titanic).render(),
    'vaccine': dfSummary(vaccine).render()})
```

Data Frame Summary

data

Dimensions: 100 x 6
Duplicates: 0

No	Variable	Stats / Values	Freqs / (% of Valid)	Graph	Missing
1	<code>catCol</code> [object]	1. B 2. A 3. D 4. C 1. R	29 (29.0%) 29 (29.0%) 27 (27.0%) 15 (15.0%) 29 (29.0%)		0 (0.0%)

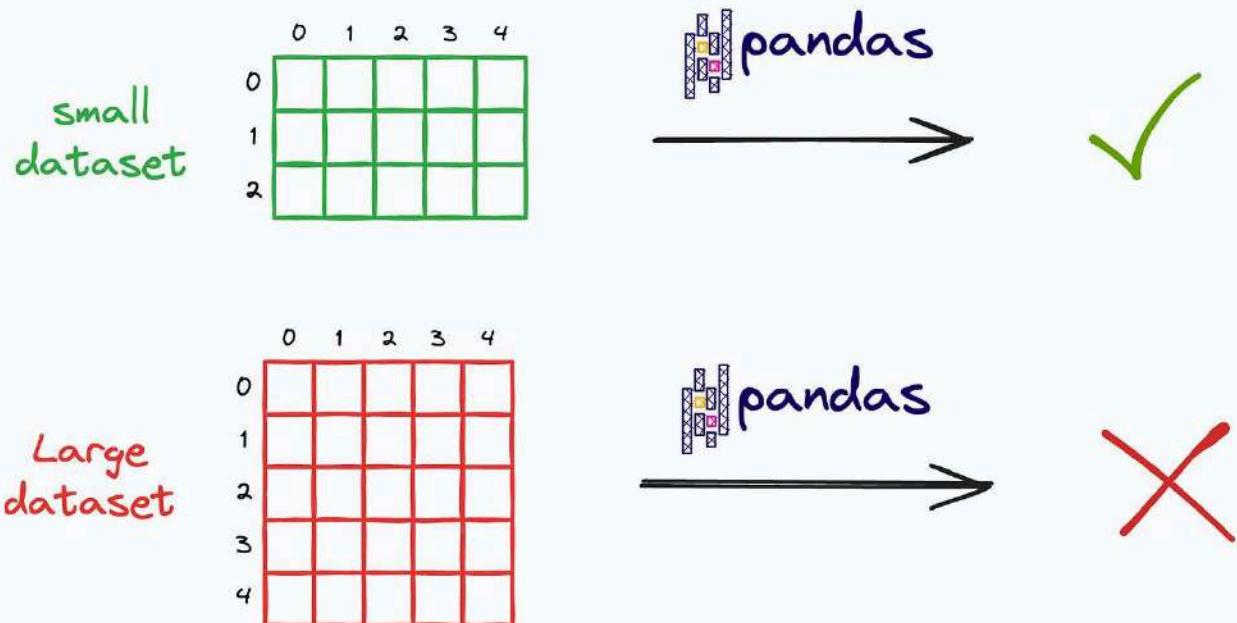
The only thing I don't like about SummaryTools is that it is not compatible with Polars (yet). Nonetheless, I find both of them extremely promising for understanding my dataset with more granularity than Pandas' `describe()` method.

Try by downloading this Jupyter notebook: <https://bit.ly/45Rheq1>.

Accelerate Pandas with GPU Using RAPIDS cuDF

Two of the biggest problems with Pandas is that:

- It always adheres to a single-core computation on a CPU.
- It creates bulky DataFrames.



While many libraries (Polars, for instance) do address these limitations, they are still limited to CPU-driven computations.

NVIDIA's RAPIDS cuDF library allows Pandas users to supercharge their Pandas workflow with GPUs.

How to use it?

Within a GPU runtime, do the following :

- Load the extension: `%load_ext cudf.pandas`
- Import Pandas: `import pandas as pd`

Done! Use Pandas' methods as you usually would.

Just loading the extension provides immense speedups.



As per NVIDIA's official release, this can be as fast as 150x. In my personal experimentation, however, I mostly observed it to range between 50-70x, which is still pretty good.

The good thing is that the extension accelerates most Pandas' methods. Yet, if needed, it can automatically fall back to the CPU.

How does it work?

Whenever `cudf.pandas` is enabled, the `import pandas as pd` statement does not import the original Pandas library which we use all the time. Instead, it imports another library that contains GPU-accelerated implementations of all Pandas methods.

```
In [1]: import pandas as pd
In [2]: pd
Out[2]: <module 'pandas' from '/Users/avi/opt/anaconda3/...'

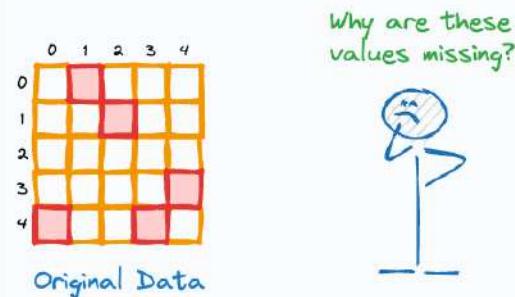
In [1]: %load_ext cudf.pandas
In [2]: import pandas as pd
In [3]: pd
Out[3]: <module 'pandas' (ModuleAccelerator(fast=cudf, slow=pandas))>
```

This alternative implementation preserves the entire syntax of Pandas. So if you know Pandas, you already know how to use cuDF's Pandas.

You can find the code here: <https://bit.ly/4cOAPZW>.

Missing Data Analysis with Heatmaps

Real-world datasets almost always have missing values. In most cases, it is unknown to us beforehand why values are missing.



But it's good to know that there could be multiple reasons for missing values:

- Missing Completely at Random (MCAR): The value is genuinely missing by itself and has no relation to that or any other observation.
- Missing at Random (MAR): Data is missing due to another observed variable. For instance, we may observe that the percentage of missing values differs significantly based on other variables.
- Missing NOT at Random (MNAR): This one is tricky. MNAR occurs when there is a definite pattern in the missing variable. However, it is unrelated to any feature we can observe in our data. In fact, this may depend on an unobserved feature.

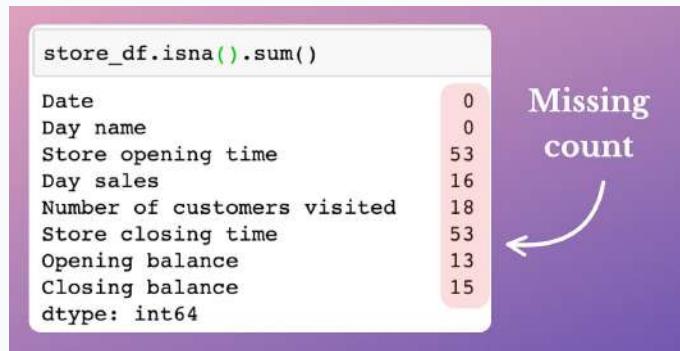
And identifying the reason for missingness can be extremely useful for further analysis, imputation, and modeling.

Consider we have a daily sales dataset of a store that has the following information:

	Date	Day name	Store opening time	Day sales	Number of customers visited	Store closing time	Opening balance	Closing balance
0	2023-01-02	Monday	09:00 AM	845.912507	102	08:00 PM	1431.060983	2276.973489
1	2023-01-03	Tuesday	08:00 AM	710.915666	77	08:00 PM	1696.911885	2407.827551
2	2023-01-04	Wednesday	10:00 AM	977.149863	133	06:00 PM	1247.242954	2224.392818
3	2023-01-05	Thursday	08:00 AM	690.447059	85	07:00 PM	1354.447274	2044.894333
4	2023-01-06	Friday	10:00 AM	Nan	43	07:00 PM	1264.830272	1785.676885
5	2023-01-07	Saturday	09:00 AM	108.363621	119	06:00 PM	1808.420181	1916.783802
6	2023-01-08	Sunday	Nan	0	0	Nan	1916.783802	1916.783802

- Day and Date
- Store opening and closing time
- Number of customers
- Total sales
- Account balance at open and close time

The reason for missing values is unknown to us. Here, when doing EDA, many folks compute the column-wise missing frequency as follows:



store_df.isna().sum()	
Date	0
Day name	0
Store opening time	53
Day sales	16
Number of customers visited	18
Store closing time	53
Opening balance	13
Closing balance	15
dtype: int64	

Missing count

The above table just highlights the number of missing values in each column.

More specifically, we get to know that:

- Missing values are relatively high in two columns compared to others.
- Missing values in the opening and closing time columns are the same (53).

That's the only info it provides. However, the problem with this approach is that it hides many important details about missing values, such as:

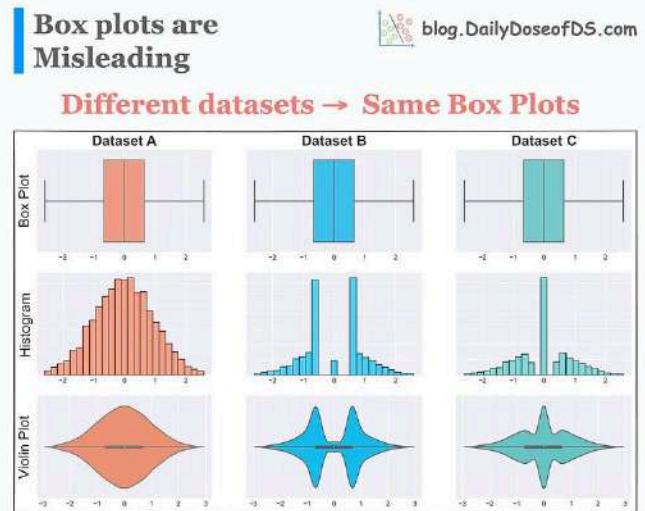
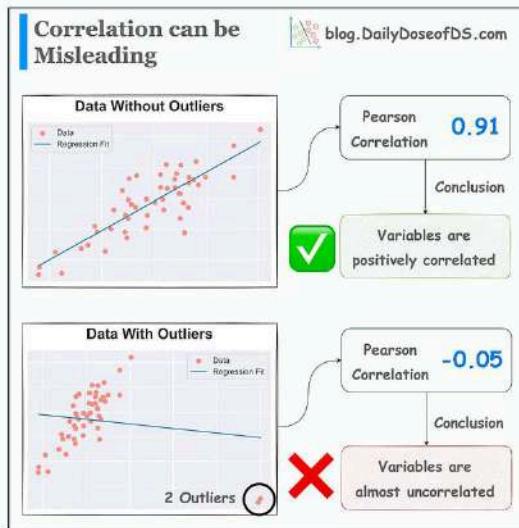
- Their specific location in the dataset.
- Periodicity of missing values (if any).
- Missing value correlation across columns, etc.

...which can be extremely useful to understand the reason for missingness.

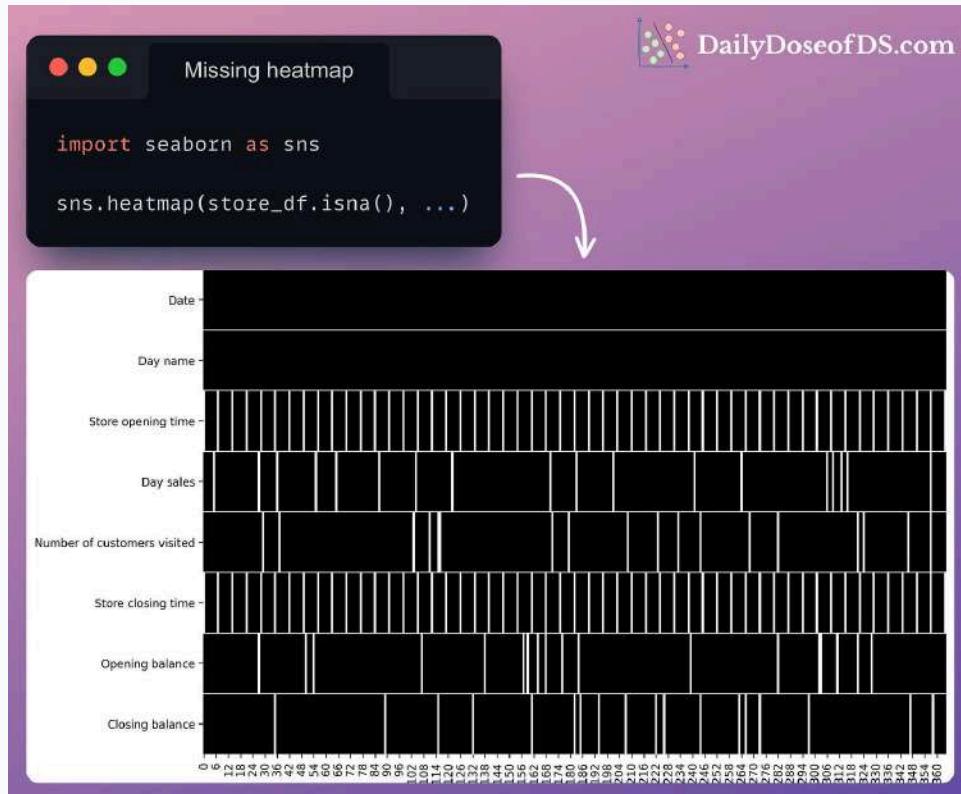
To put it another way, the above table is more like summary statistics, which rarely depict the true picture.

Why?

We discussed this in the statistics section of this book.



But here's how I often enrich my missing value analysis with heatmaps. Compare the missing value table we discussed above with the following heatmap of missing values:



The white vertical lines depict the location of missing values in a specific column.

Now, it is immediately clear that:

- Values are periodically missing in the opening and closing time columns.
- Missing values are correlated in the opening and closing time columns.
- The missing values in other columns appear to be (not necessarily though) missing completely at random.

Further analysis of the opening time lets us discover that the store always remains closed on Sundays:

```
store_df[store_df["Store opening time"].isna()]
```

Date	Day name	Store opening time	Day sales	Number of customers visited	Store closing time	Opening balance	Closing balance
0 2023-01-01	Sunday	NaN	0	0	NaN	10000	10000
7 2023-01-08	Sunday	NaN	0	0	NaN	1916.783802	1916.783802
14 2023-01-15	Sunday	NaN	0	0	NaN	1687.342885	1687.342885
21 2023-01-22	Sunday	NaN	0	0	NaN	2374.595033	2374.595033
28 2023-01-29	Sunday	NaN	0	0	NaN	2084.814504	2084.814504

Store remains
closed on
Sunday

Now, we know why the opening and closing times are missing in our dataset. This information can be beneficial during its imputation.

This specific situation is “Missing at Random (MAR).” Essentially, as we saw above, the missingness is driven by the value of another observed column.

As we know the reason, we can use either the kNN imputation or the MissForest techniques to impute these values, which we discussed in the ML section of this book ([MissForest and kNN Imputation](#)).

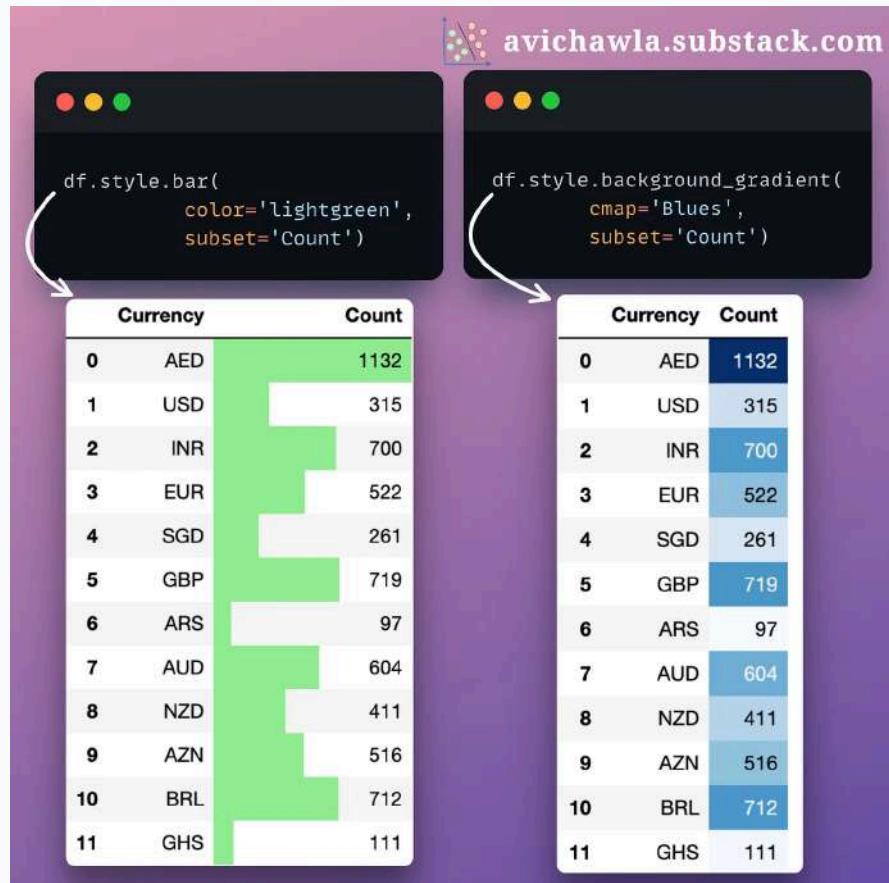
DataFrame Styling

Instead of previewing raw DataFrames, styling can make data analysis much easier and faster.

Here's how.

Jupyter is a web-based IDE. Anything you print is rendered using HTML and CSS.

This means you can style your output in many different ways.



To style Pandas DataFrames, use its Styling API (**df.style**).

As a result, the DataFrame is rendered with the specified styling.

Pandas documentation page on Styling: <https://bit.ly/3zDmxgE>.

8 Automated EDA Tools

Below are 8 powerful EDA tools that automate many redundant EDA steps and help you profile your data quickly.

8 Automated EDA Tools



DailyDoseofDS.com

1 SweetViz



2 ydata-Profilin



3 DataPrep



4 AutoViz



5 D-Tale



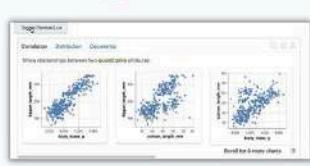
6 dabl



7 QuickDA



8 Lux



Please note that these tools are not the ultimate EDA alternatives that will answer all your questions about the dataset.

But given that the preliminary EDA steps in almost all projects are the same — plotting the response variable, checking imbalance, running correlation analysis, missing value analysis, and more, these tools pretty well automate these steps in my opinion.

Also, at times, manual EDA can be prone to human errors and one may miss out on checking a few things. Automated tools eliminate these risks and provide a standardized report across all projects.

Full summary about each tool is available here: <https://bit.ly/3xZP4fK>.

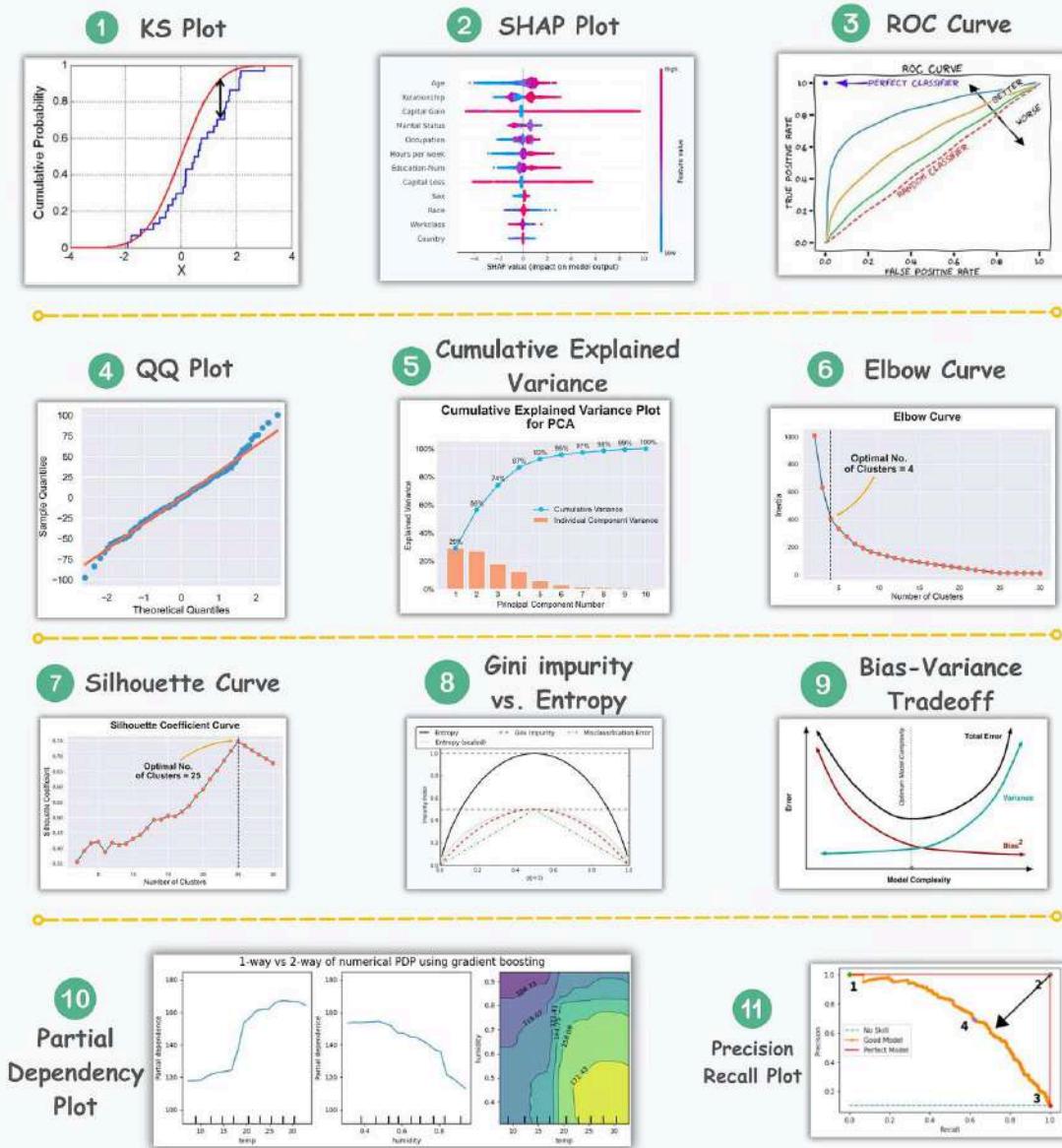
Data Visualisation

Most Important Plots in Data Science

The visual below depicts the 11 most important and must-know plots in data science:

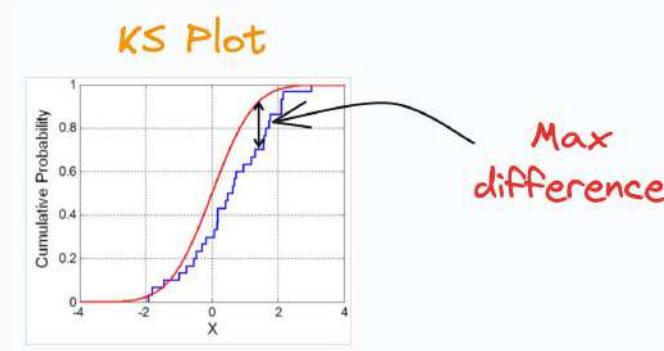
11 Most Important Plots in Data Science

 blog.DailyDoseofDS.com



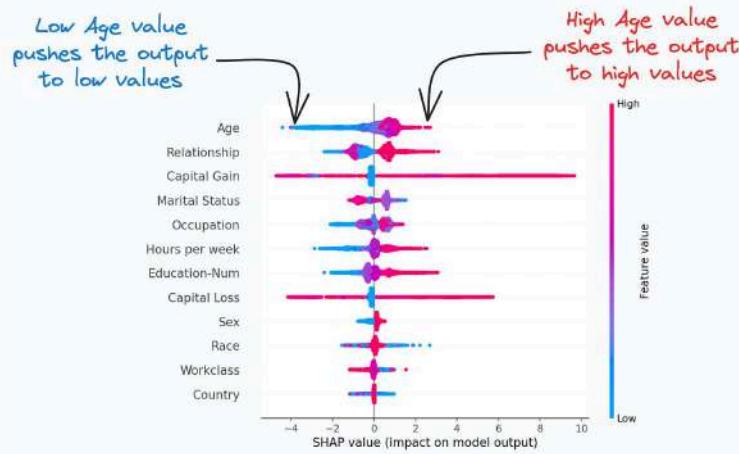
Let's understand them briefly and how they are used.

KS Plot:



- It is used to assess the distributional differences.
- The core idea is to measure the maximum distance between the cumulative distribution functions (CDF) of two distributions.
- The lower the maximum distance, the more likely they belong to the same distribution.
- Thus, instead of a “plot”, it is mainly interpreted as a “statistical test” to determine distributional differences.

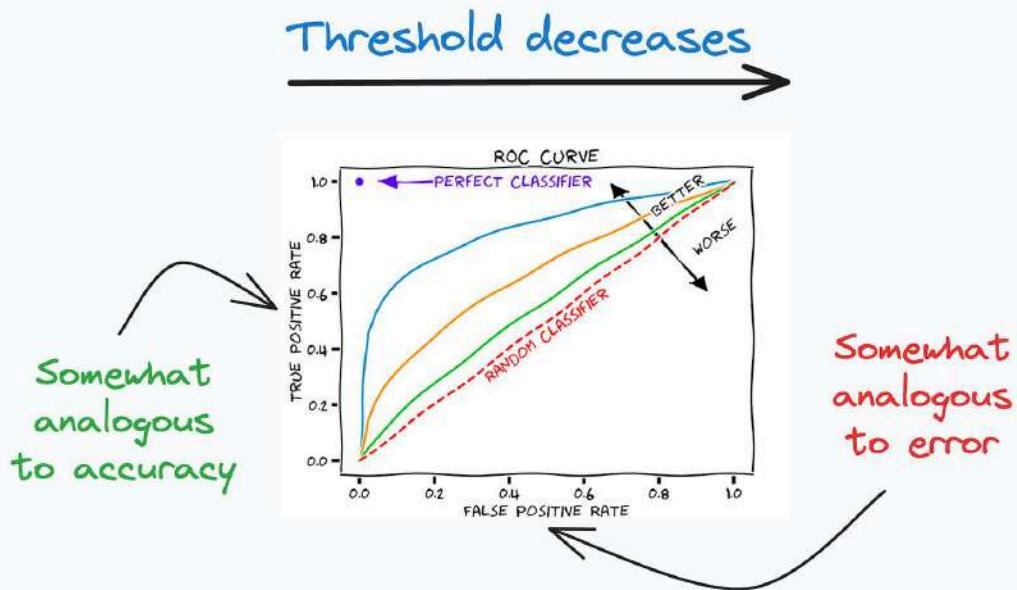
SHAP Plot:



- It summarizes feature importance to a model's predictions by considering interactions/dependencies between them.

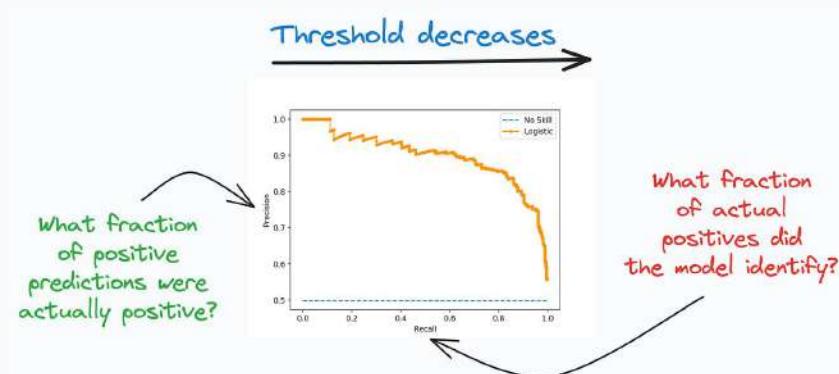
- It is useful in determining how different values (low or high) of a feature affect the overall output.

ROC Curve:



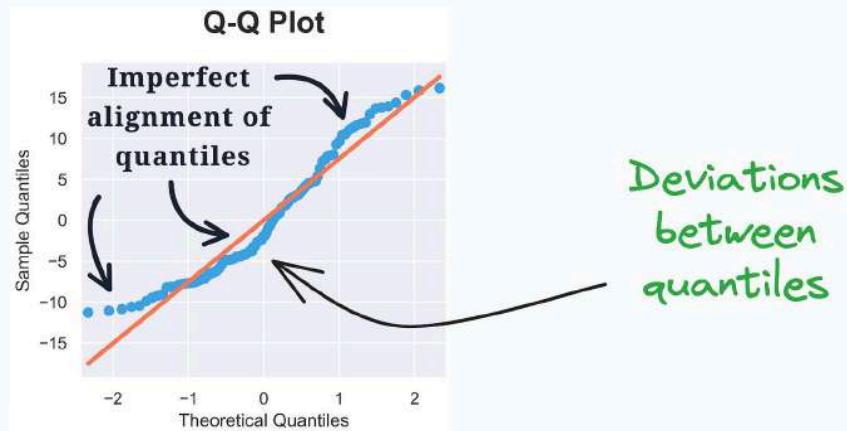
- It depicts the tradeoff between the true positive rate (good performance) and the false positive rate (bad performance) across different classification thresholds.
- The idea is to balance TPR (good performance) vs. FPR (bad performance).

Precision-Recall Curve:



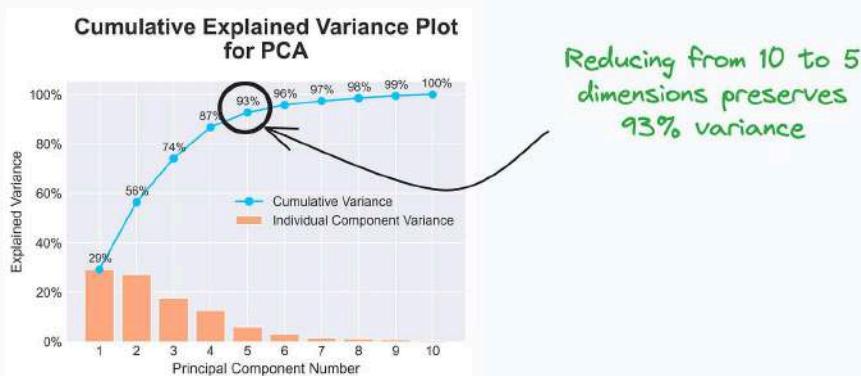
- It depicts the tradeoff between Precision and Recall across different classification thresholds.

QQ Plot:



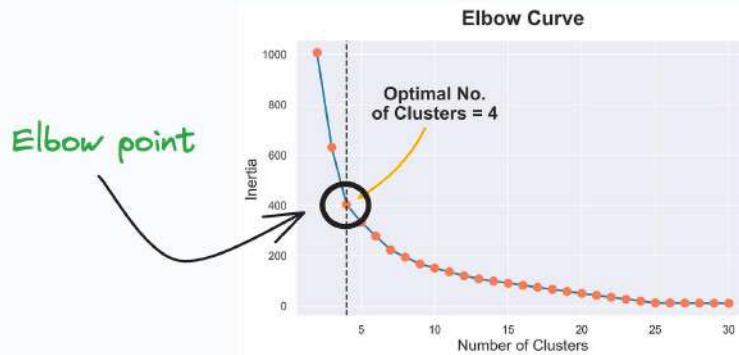
- It assesses the distributional similarity between observed data and theoretical distribution.
- It plots the quantiles of the two distributions against each other.
- Deviations from the straight line indicate a departure from the assumed distribution.

Cumulative Explained Variance Plot:



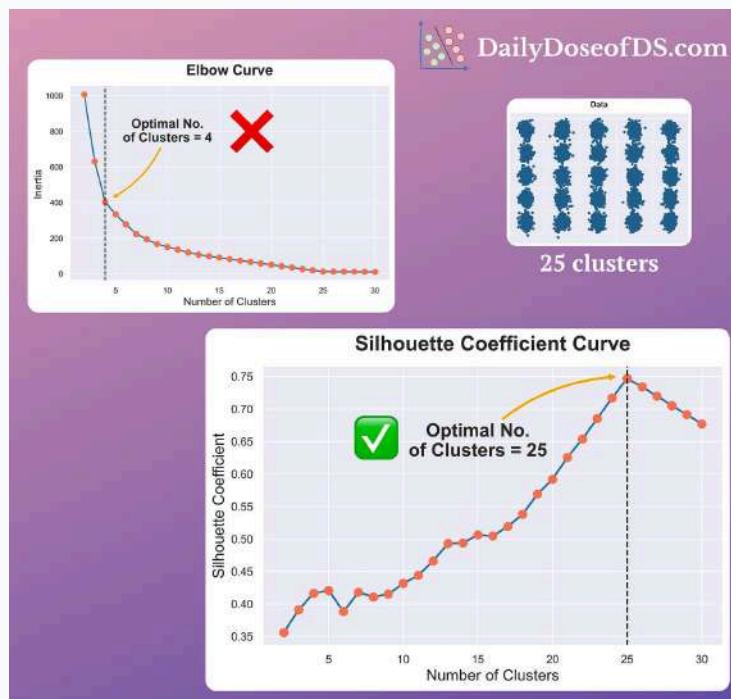
- It is useful in determining the number of dimensions we can reduce our data to while preserving max variance during PCA.

Elbow Curve:



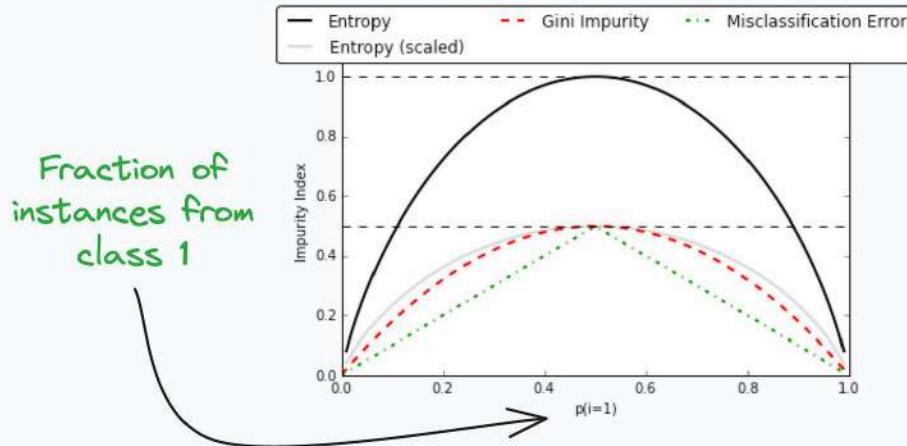
- The plot helps identify the optimal number of clusters for the k-means algorithm.
- The point of the elbow depicts the ideal number of clusters.

Silhouette Curve:



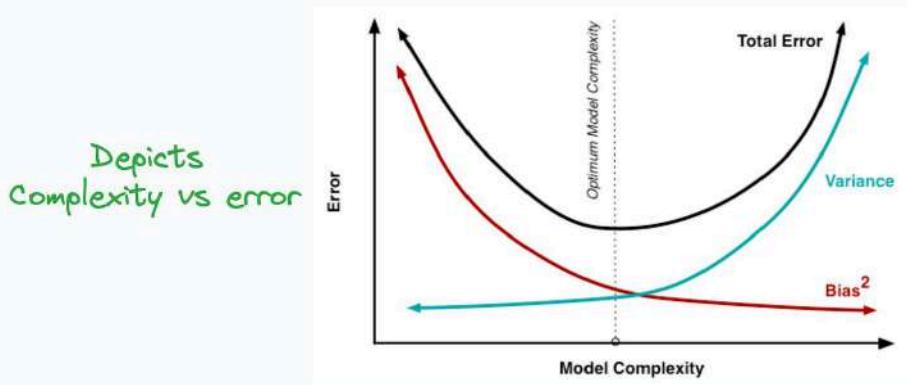
- The Elbow curve is often ineffective when you have plenty of clusters.
- Silhouette Curve is a better alternative, as depicted above.

Gini-Impurity and Entropy:



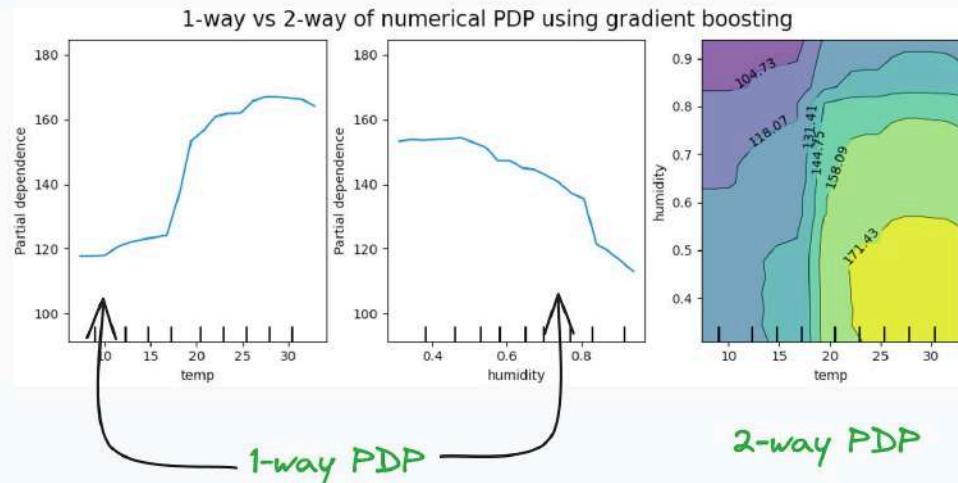
- They are used to measure the impurity or disorder of a node or split in a decision tree.
- The plot compares Gini impurity and Entropy across different splits.
- This provides insights into the tradeoff between these measures.

Bias-Variance Tradeoff:



- It's probably the most popular plot on this list.
- It is used to find the right balance between the bias and the variance of a model against complexity.

Partial Dependency Plots:



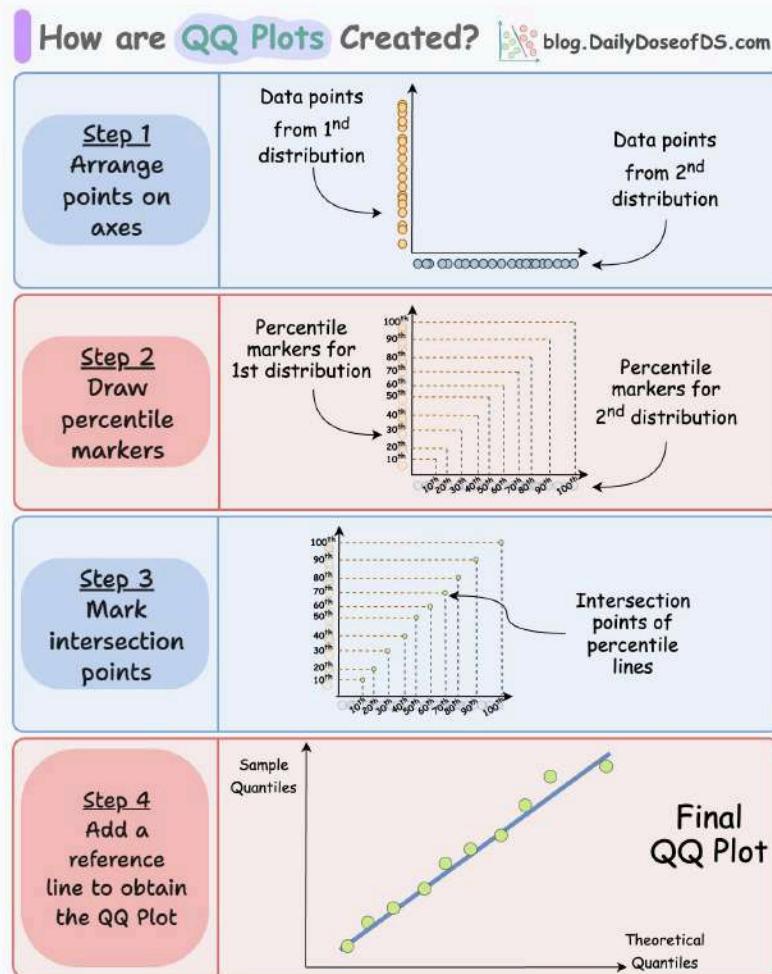
- Depicts the dependence between target and features.
- A plot between the target and one feature forms → 1-way PDP.
- A plot between the target and two feature forms → 2-way PDP.
- In the leftmost plot, an increase in temperature generally results in a higher target value.

How are QQ Plots Created?

A QQ plot is a great way to visually assess the similarity between two distributions.



It does this by plotting the quantiles of the two distributions against each other. The deviations from the straight line indicate the differences between the two distributions. The following visual depicts how it is created:

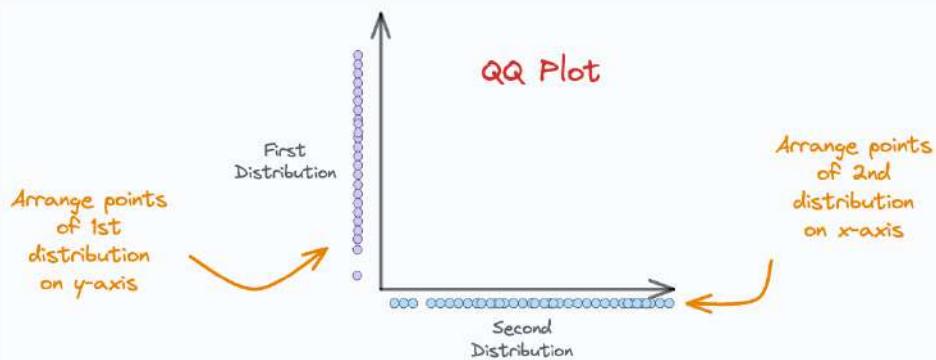


Let's discuss it in more detail.

Consider we have two distributions, D1 and D2.

Step 1) Arrange points on axes:

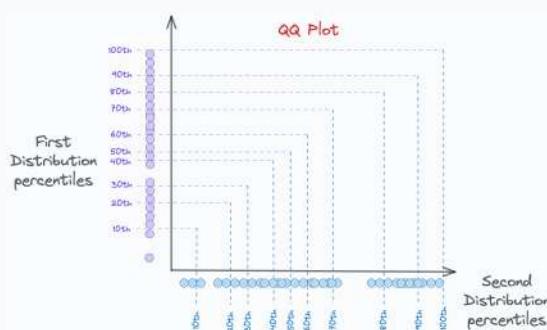
As shown below, we arrange points of D1 on the y-axis and D2 on the x-axis.



Step 2) Draw percentile lines

Next, for both distributions, we create some percentile lines.

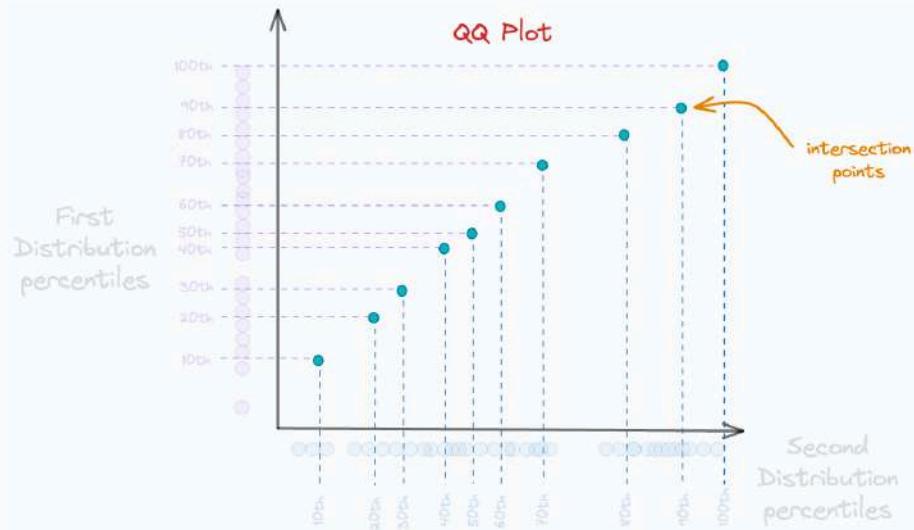
For instance, on both axes, we can mark the points of 10th percentile, 20th percentile, 30th percentile, etc., from both distributions. This is shown below:



We mark the percentile locations for both distributions and intersect the corresponding lines.

- 10th percentile of D1 is intersected with 10th percentile of D2.
- 20th percentile of D1 is intersected with 20th percentile of D2.
- and so on.

The intersection points of these percentile lines gives us the points we typically see in a QQ plot:



Now, we can get rid of the percentile marker lines.



In a gist, the above plot gives us the location where the corresponding percentiles of the two distributions match.

Step 3) Add the reference line

Finally, we must add a reference line to determine the deviations between the two distributions.

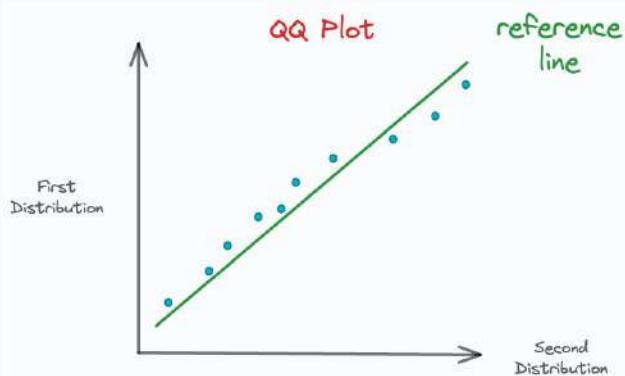
There are many ways to do this.

For instance:

- The line connecting the 25th and 75th percentiles of both distributions can be considered as a reference line.
- The regression fit on the above scatter plot can be considered as a reference line.

Typically, the line connecting the 25-75th percentile is preferred because the regression fit can be influenced by outliers.

After adding the reference line, we get our QQ plot:



The deviations from this reference line indicate that the two distributions differ from each other. In other words, the deviations mean that the corresponding percentiles do not align.



This becomes an indicator of distributional dissimilarities.

And, of course, the more percentiles we plot, the better and more accurate the QQ plot will be.

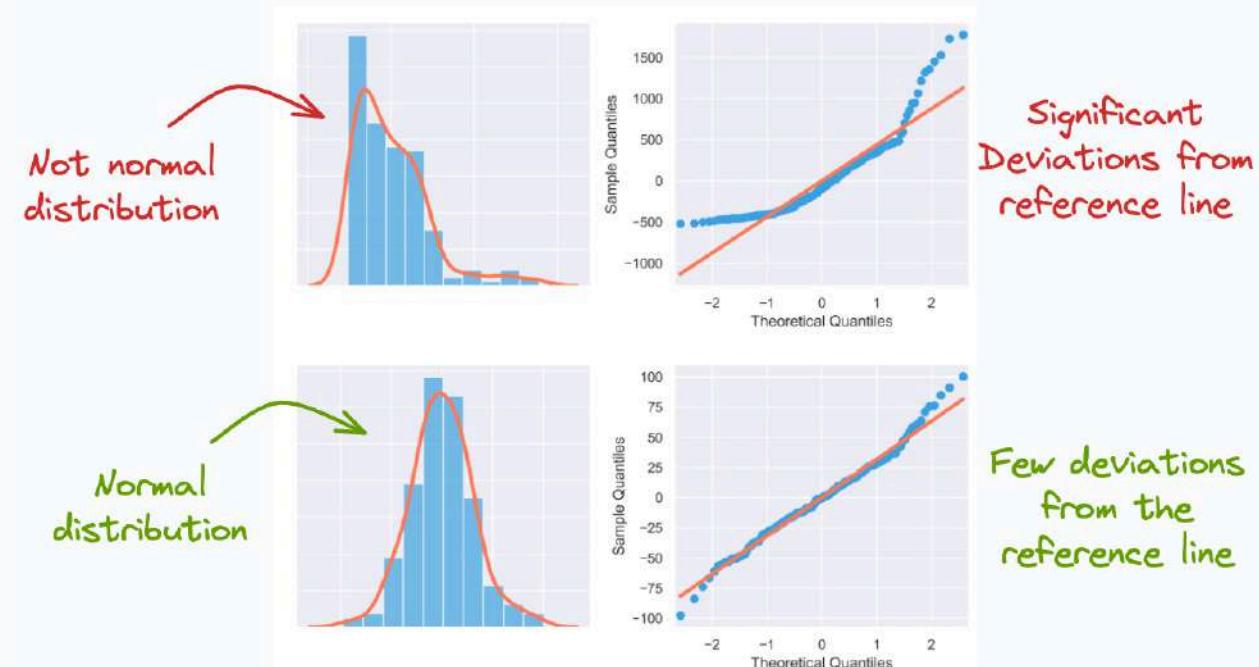
There are many applications of the QQ plot.

For instance, say we have an observed distribution and want to determine if it resembles a normal distribution.

We can use a QQ plot for this:

- D1: The observed distribution
- D2: Normal distribution.

If the percentile points lie closer to the reference line, this would mean that the observed distribution is more like a normal distribution. This is depicted below:



8 Elegant Alternatives to Traditional Plots

Scatter plots, bar plots, line plots, box plots, and heatmaps are the most frequently used plots for data visualization. Although they are simple and known to almost everyone, I believe they are not the right choice to cover every possible scenario.

Instead, many other plots originate from these standard plots that can be much more suitable, if used appropriately. In this chapter, let's discuss a few alternatives to these popular plots. We will also understand specific situations where they can be more useful over standard plots.

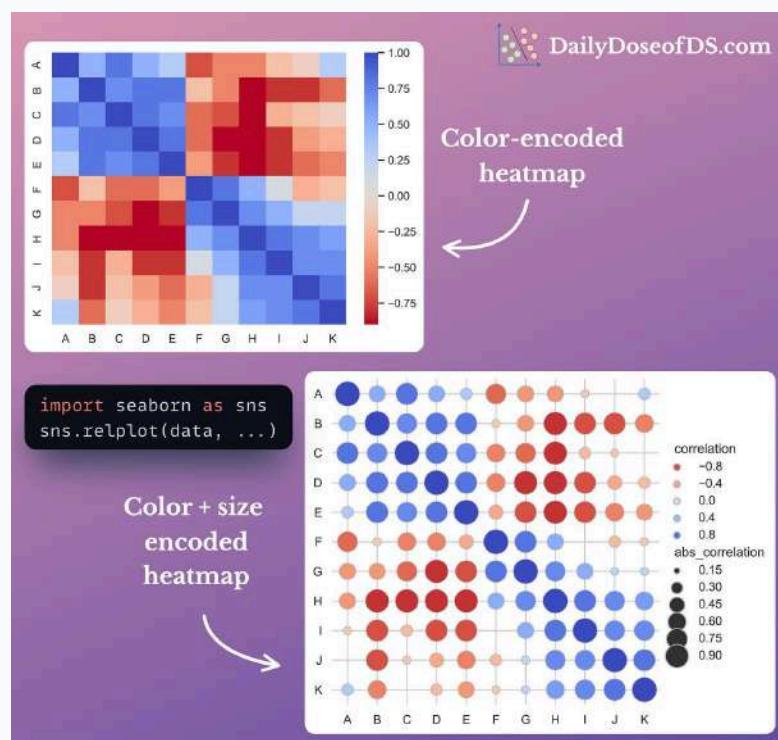
This chapter is not intended to discourage the use of these traditional plots. They will always have their place.

Instead, it is to highlight specific situations where they can be replaced with better plotting ideas.

#1) Size-encoded heatmaps

A traditional heatmap represents the values using a color scale. Yet, mapping the cell color to exact numbers is still challenging.

Embedding a size component to heatmaps can be extremely helpful in such cases. In essence, the bigger the size, the higher the absolute value



This is especially useful to make heatmaps cleaner, as many values nearer to zero will immediately shrink.

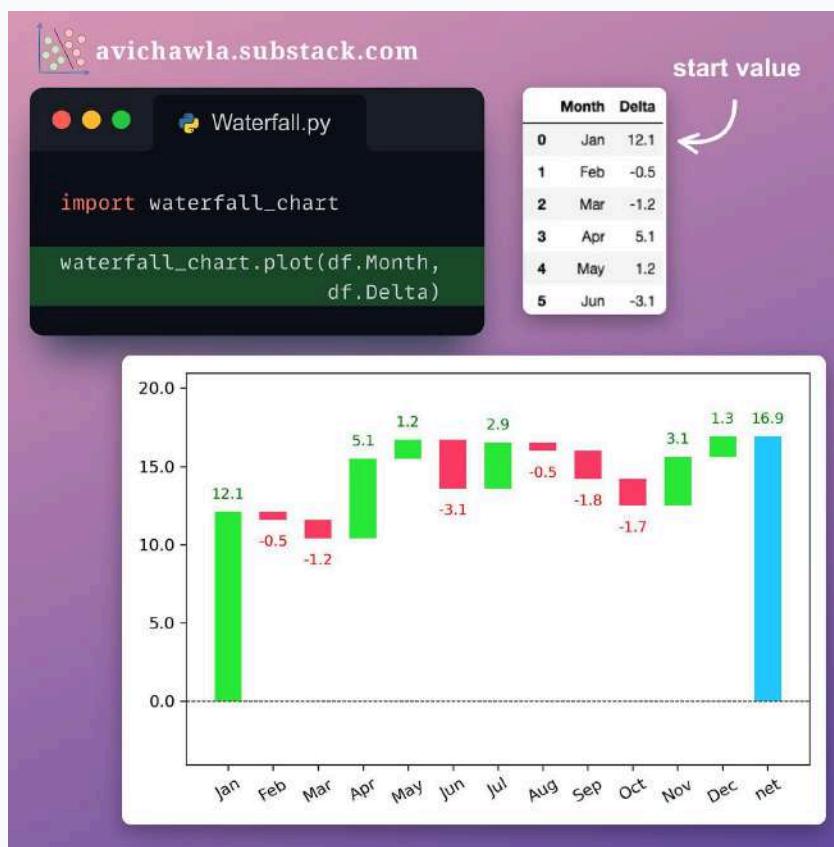
#2) Waterfall charts

To visualize the change in value over time, a line (or bar) plot may not always be an apt choice.

This is because a line plot (or bar plot) depicts the actual values in the chart. Thus, it is difficult to visually estimate the scale and direction of incremental changes.

Instead, you can use a waterfall chart. It elegantly depicts these rolling differences, as depicted in this image.

Here, the start and final values are represented by the first and last bars.



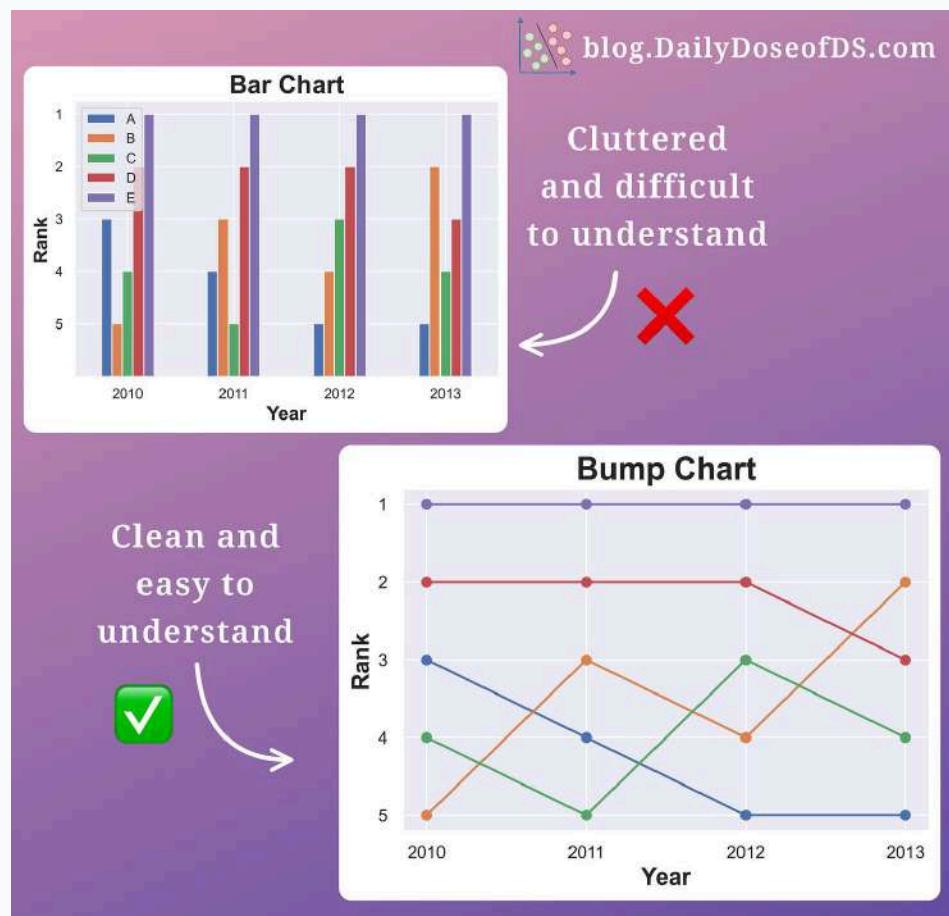
Also, the consecutive changes are automatically color-coded, making them easier to interpret.

#3) Bump charts

When visualizing the change in rank over time of multiple categories, using a bar chart may not be appropriate.

This is because bar charts quickly become cluttered with many categories.

Instead, try Bump Charts. They are specifically used to visualize the rank of different items over time.



Comparing the bar chart and bump chart above, it is far easier to interpret the change in rank with a bump chart rather than a bar chart.

#4) Raincloud Plots

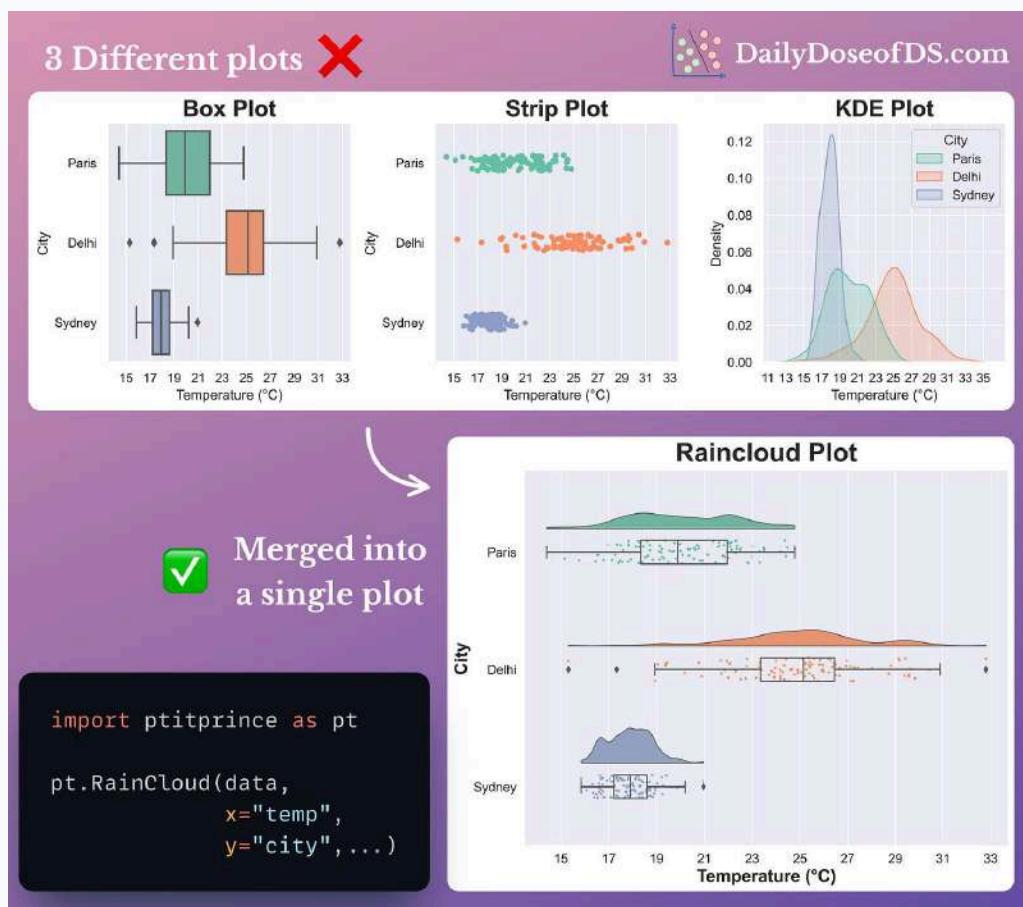
Visualizing data distributions using box plots and histograms can be misleading at times. This is because:

- It is possible to get the same box plot with entirely different data.
- Altering the number of bins changes the shape of a histogram.

Thus, to avoid misleading conclusions, it is always recommended to plot the data distribution as precisely as possible. Raincloud plots provide a concise way to combine and visualize three different types of plots together.

These include:

- Box plots for data statistics.
- Strip plots for data overview.
- KDE plots for the probability distribution of data.



With Raincloud plots, you can:

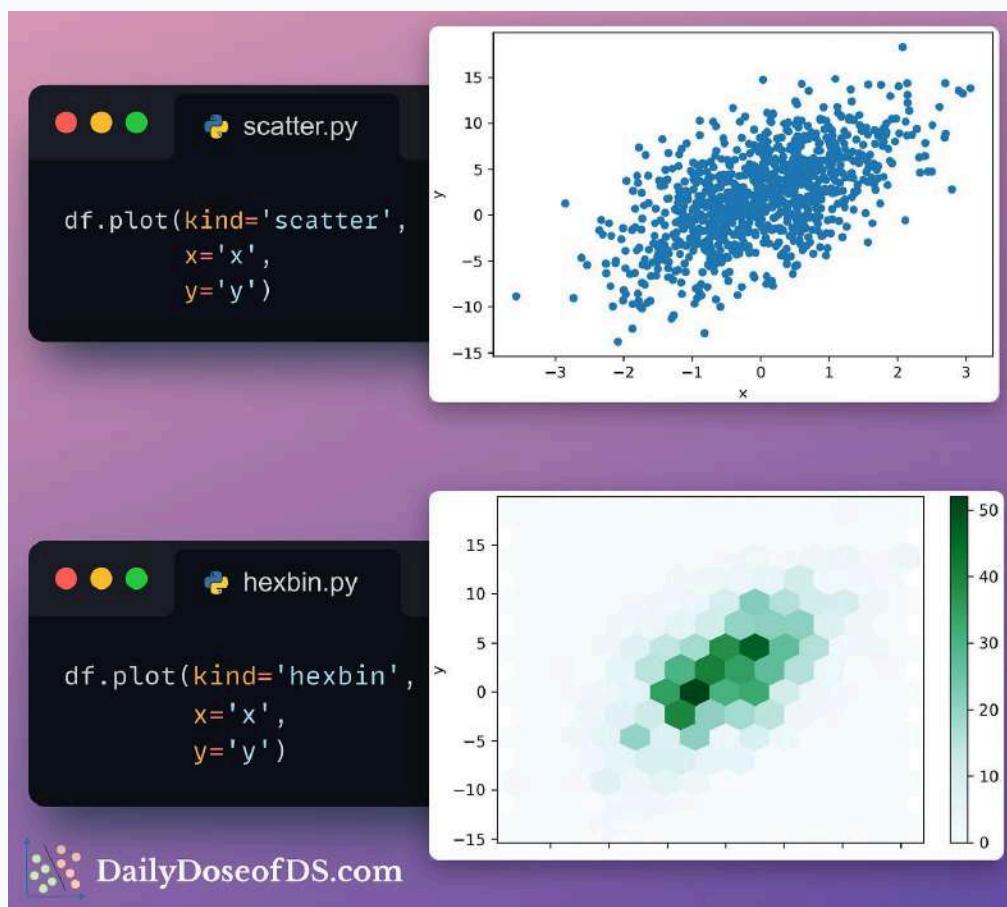
- Combine multiple plots to prevent incorrect/misleading conclusions
- Reduce clutter and enhance clarity
- Improve comparisons between groups
- Capture different aspects of the data through a single plot

#5-6) Hexbin and Density Plots

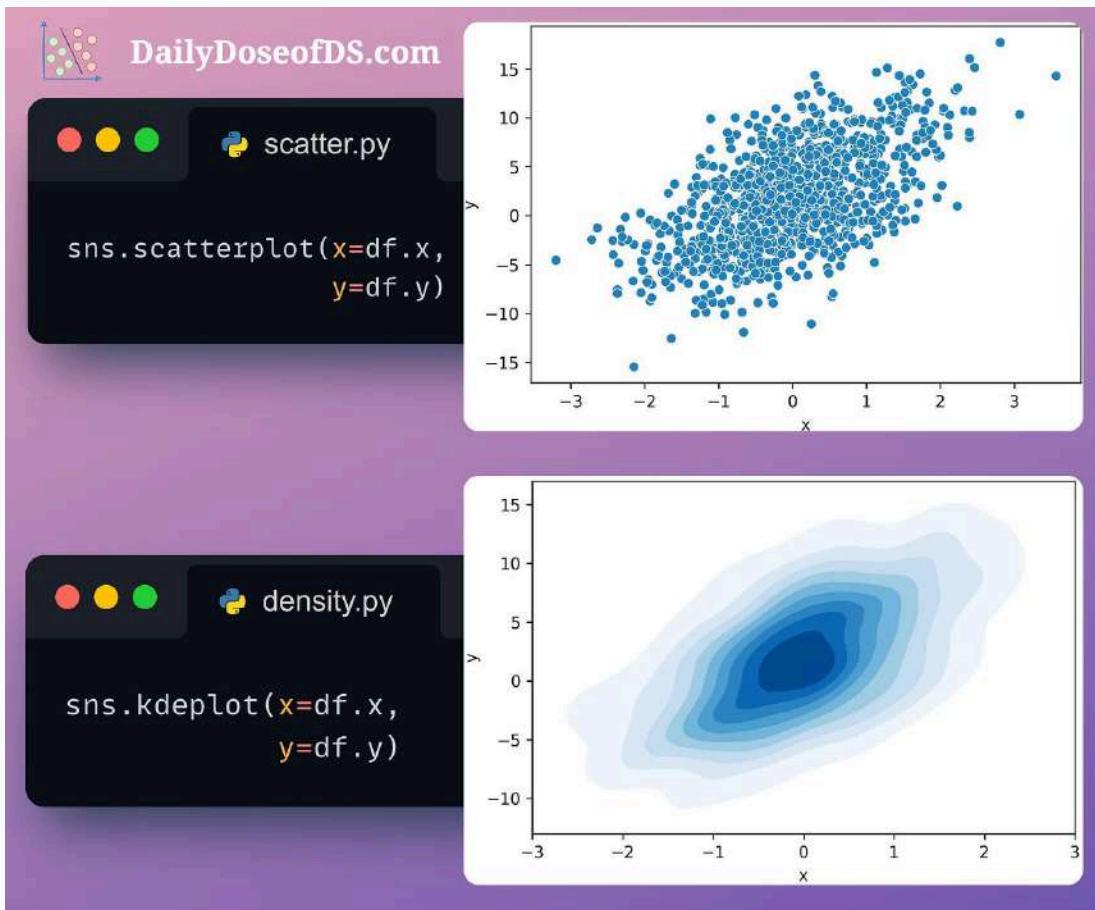
Scatter plots can get too dense to interpret when you have thousands of data points.

Instead, you can replace them with Hexbin plots.

Hexbin plots bin the area of a chart into hexagonal regions. Each region is assigned a color intensity based on the method of aggregation used (the number of points, for instance).



Another choice is a density plot, which illustrates the distribution of points in a two-dimensional space.



A contour is created by connecting points of equal density. In other words, a single contour line depicts an equal density of data points.

#7-8) Bubble charts and Dot plots

As discussed above, bar plots quickly get messy and cluttered as the number of categories increases.

A bubble plot is often a better alternative in such cases.

They are like scatter plots but:

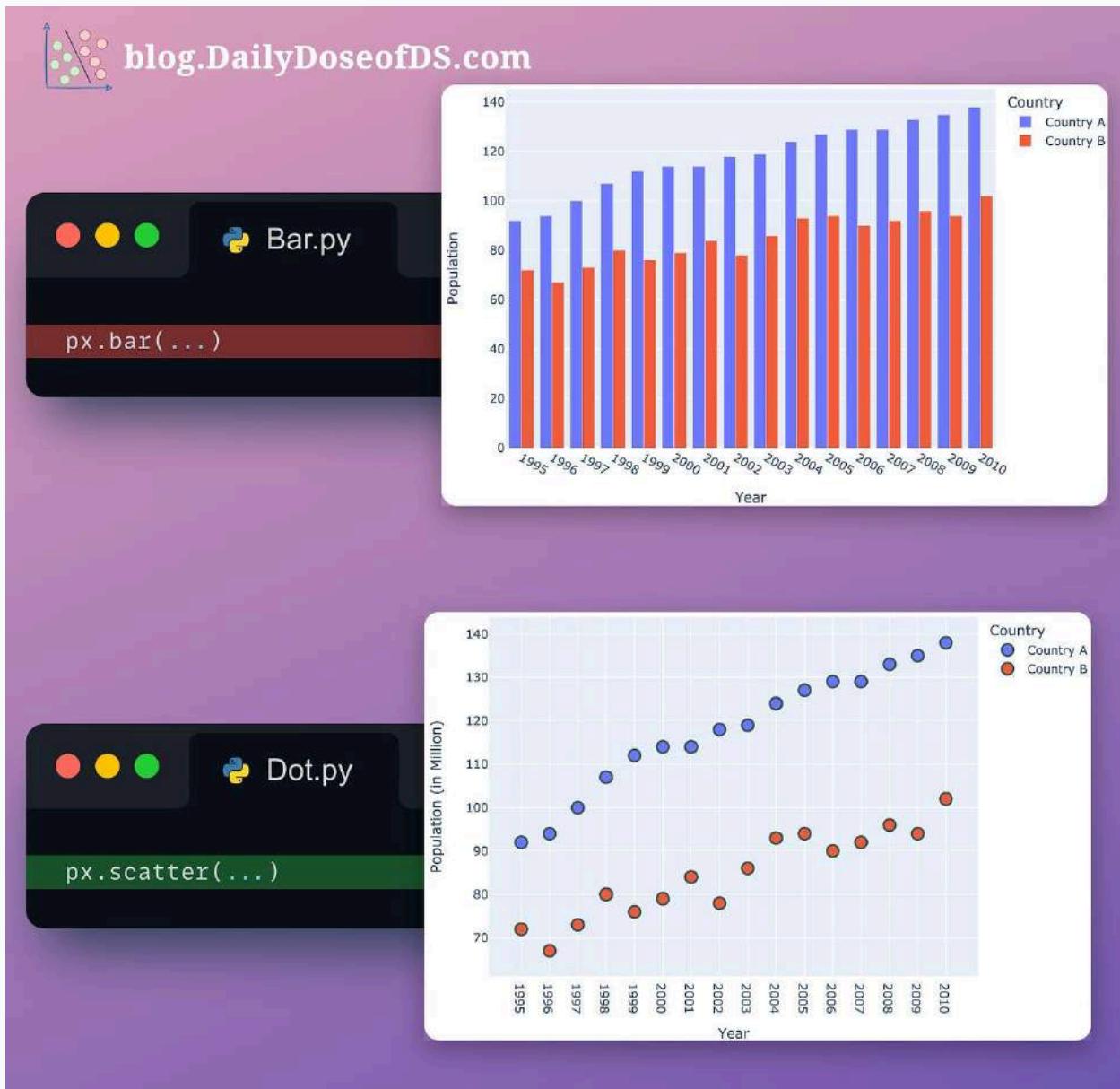
- with one categorical axis
- and one continuous axis



As depicted above:

- It is difficult to interpret the bar plot because it has too many bars packed into a small space,
- But size-encoded bubbles make it pretty easy to visualize the change over time.

Another alternative to bar plots in such situations is dot plots.



Both dot plots and bubble charts are based on the idea that, at times, when we have a bar plot with many bars, we're often not paying attention to the individual bar lengths.

Instead, we mostly consider the individual endpoints that denote the total value.

These plots precisely help us depict that while also eliminating the long bars of little to no use.

Interactive Controls

While using Jupyter, one often finds themselves in situations where they repeatedly modify a cell and re-rerun it. This makes data exploration irreproducible and time-consuming. What's more, the notebook also gets messy and cluttered.

To address this, one of the things I actively leverage in my Jupyter notebooks is interactive controls using `IPywidgets`.

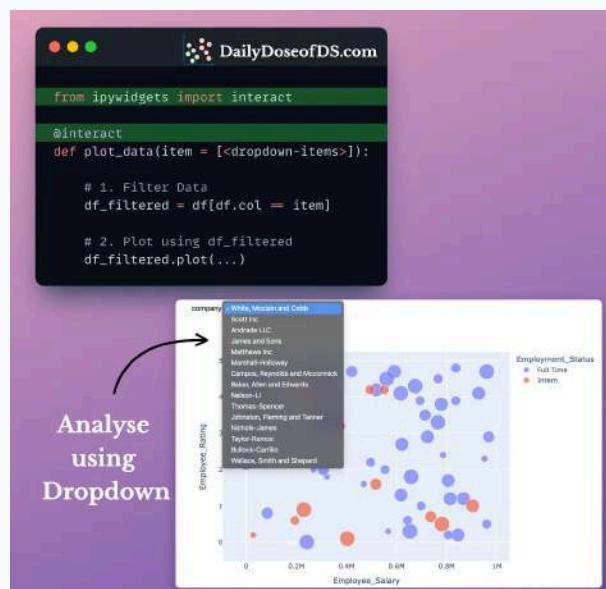
A single decorator (`@interact`) allows us to add:

- sliders
- dropdowns
- text fields, and more.

As a result, one can:

- explore the data interactively
- speed-up data exploration
- avoid repetitive cell modifications and executions
- organize the data analysis.

A Jupyter Notebook is available here to get started: <https://bit.ly/3XZtsea>.



Mosaic Plots

We mostly use `plt.subplots()` method to create subplots using Matplotlib. But this, at times, gets pretty tedious and cumbersome. For instance, it offers limited flexibility to create a custom layout, it is prone to indexing errors, and more.

Instead, use the `plt.subplot_mosaic()` method. Here, you can create a plot of any desired layout by defining the plot structure as a string.

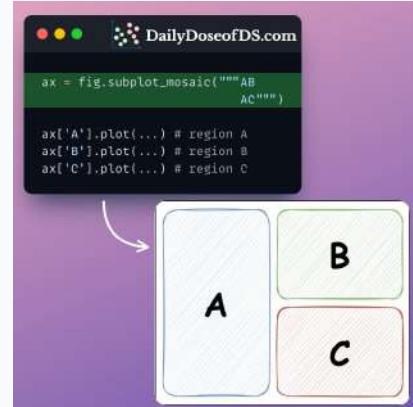


For instance, the string layout:

- AB
- AC

...will create three subplots, wherein:

- subplot "A" spans 1st column
- subplot "B" spans top half of the 2nd column
- subplot "C" spans the bottom half of 2nd column



Next, create a subplot in a specific region by indexing the axes dictionary with its subplot key ("A", "B", or "C").



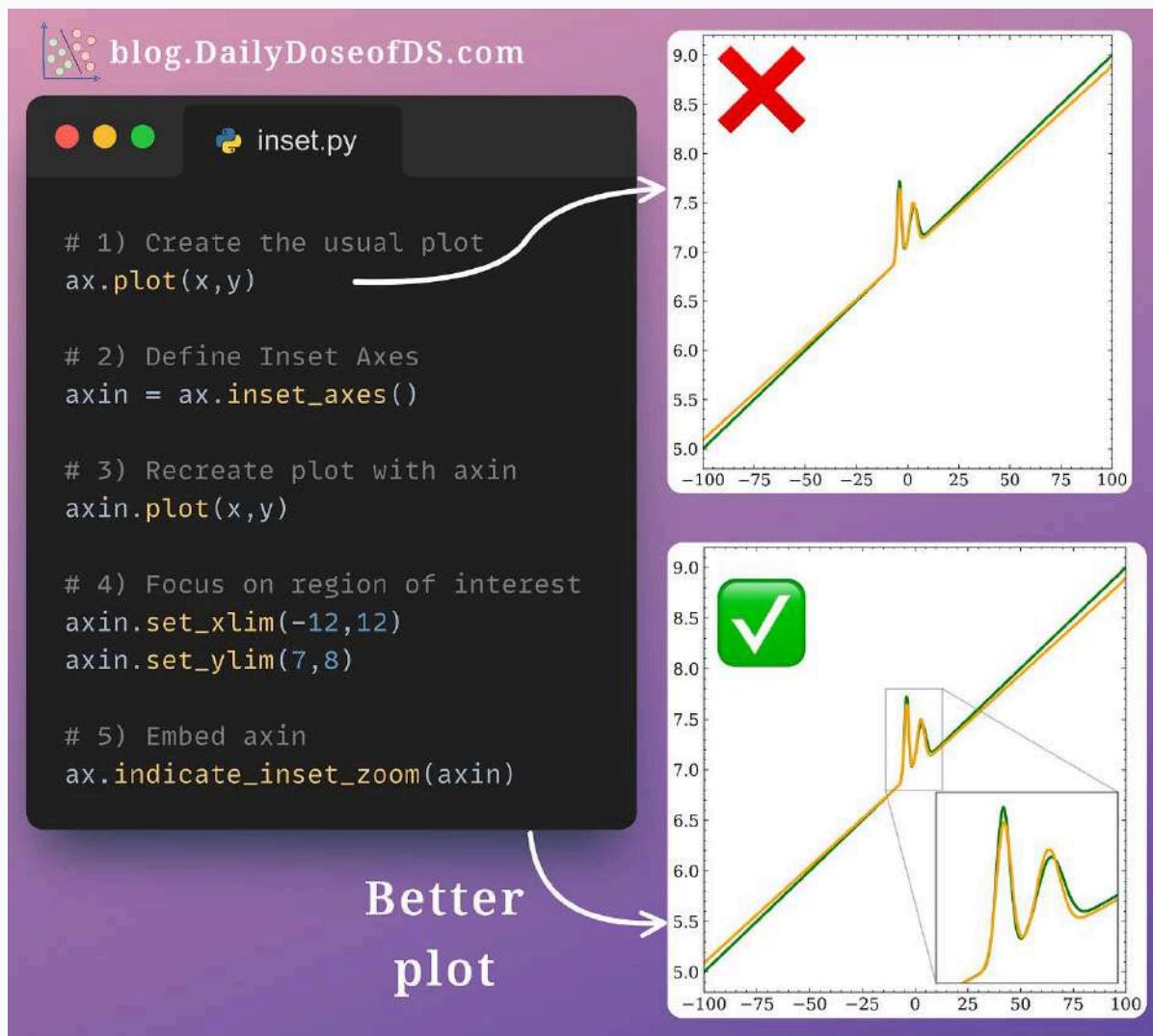
Enrich Matplotlib Plots with Inset Axis and Annotations

While creating data visualizations, there are often certain parts that are particularly important.

Yet, they may not be immediately obvious to the viewer.

A good data storyteller always ensures that the plot guides the viewer's attention to these key areas.

One great way is to zoom in on specific regions of interest in a plot, as depicted below.



In contrast to the usual plot, the other plot guides the viewer's attention to a specific area of interest.

Such efforts always ensure that the plot communicates what we intend it to depict — even if the plot's creator is not present at that time.

In matplotlib, we can do so using `indicate_inset_zoom()`. It adds an indicator box, which can be zoomed-in for better clarity.

The embedded plot is treated like any other matplotlib plot. Thus, we can add axis labels to it, if needed.

Another great way to provide extra info is by adding text annotations to a plot, as depicted below:



Such efforts always ensure that the plot indeed communicates what we intend it to depict — even if the plot's creator is not present at that time.

In matplotlib, you can use `annotate()`, as depicted below:



It adds explanatory texts to your plot, which lets you guide a viewer's attention to specific areas and aid their understanding.

Jupyter notebook for plot annotations is available here: <https://bit.ly/4bzQ5Ja>.

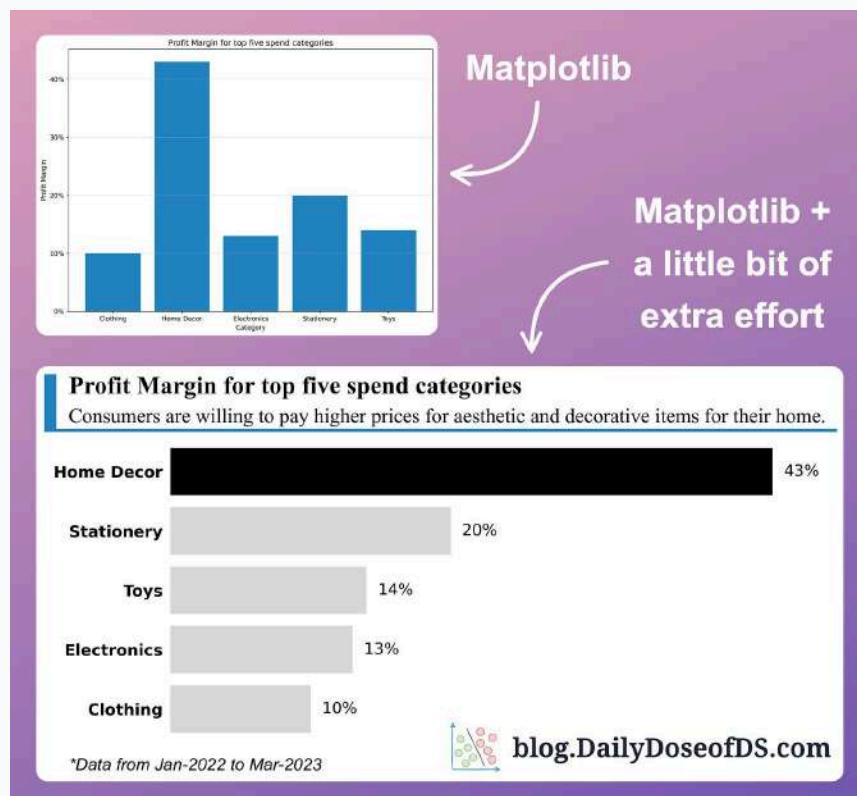
Professionalize Matplotlib Plots

Context

I have been using matplotlib for many years now. Based on that experience, I believe that one of the best yet underrated potentials of matplotlib is the amount of customizability it offers.

But being unaware of that, most matplotlib users use it as a naive plotting utility with almost zero customization. And as the default plots never appear “appealing”, they resort to other libraries, Plotly, for instance, to create elegant plots.

Yet, I believe that in 90-95% of cases, you would NEVER need to look beyond matplotlib. It can do much more than what most users think. For instance, consider the two plots in this image.



Yes! Both plots were created using matplotlib.

But some custom formatting makes the second plot much more elegant, informative, appealing, and easy to follow.

- The title and subtitle significantly aid the story.

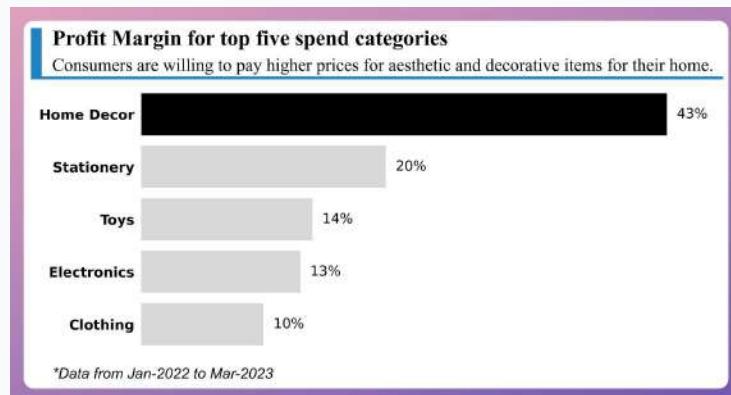
- Also, the footnote offers extra important information, which is nowhere to be seen in the basic plot.
- Lastly, the bold bar immediately draws the viewer's attention and conveys the purchase category's importance.

Thus, in my opinion, the overwhelming potential for customization makes matplotlib far more capable than what most users think.

Conclusion

One of the things I always ensure towards being a good storyteller in my data science projects is that my plot must demand minimal effort from the viewer. Thus, I never shy away from putting in that extra effort.

This has been especially true for professional environments. At times, it is also good to ensure that our visualizations convey the right story, even if they are viewed in our absence. The below plot is a classic example of that.

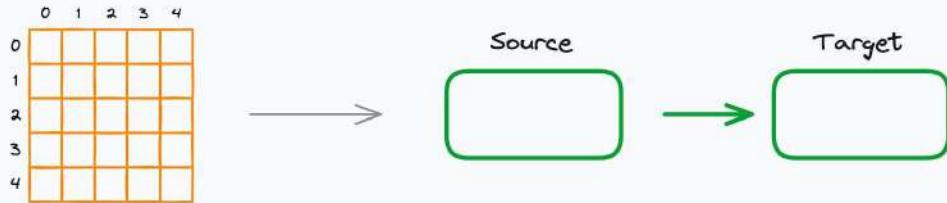


In this entire chapter, I never discussed what that plot is about — somewhat indicating my absence. Yet, by staring at this plot for a few seconds, you can quickly figure out what I intended to highlight here, can't you?

You can download the code notebook for this post here: <https://bit.ly/3zv8pWG>.

Sankey Diagrams

Many tabular data analysis tasks can be interpreted as a flow between several source and target entities.



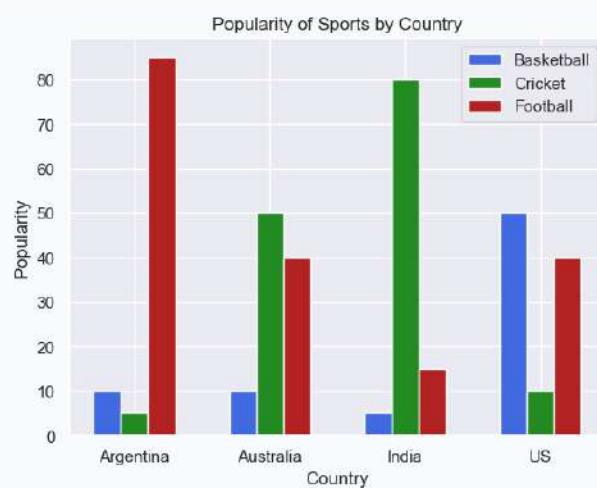
For instance, consider we have a sports popularity dataset, which lists the popularity index of a sport in a country.

In this dataset:

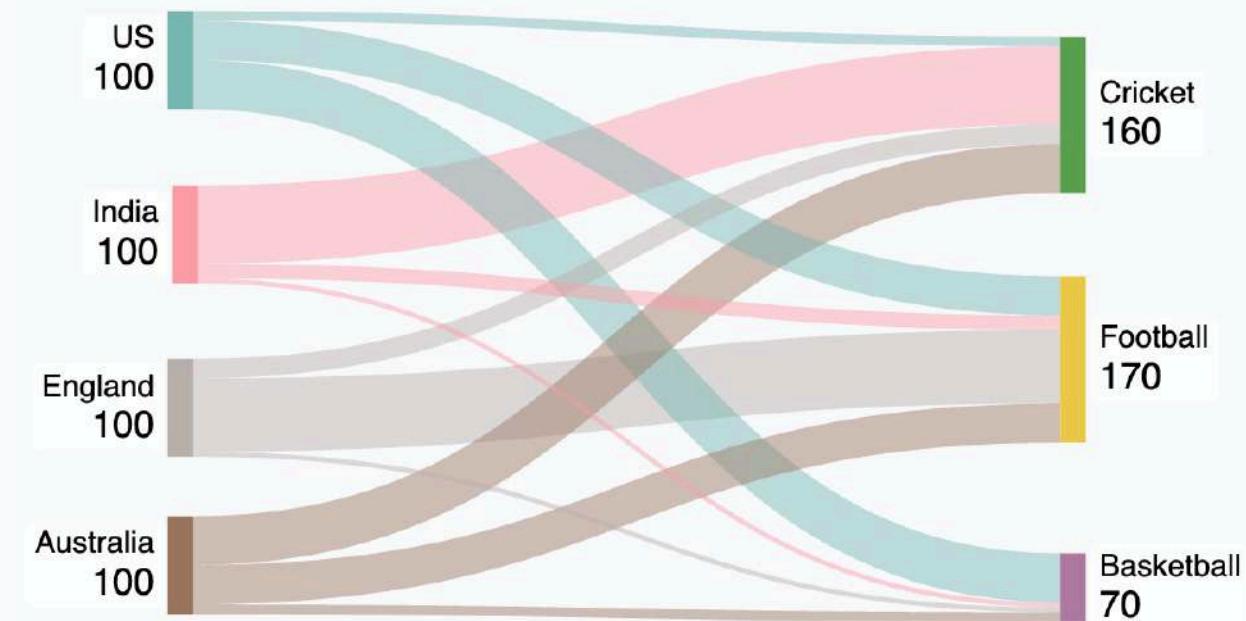
- Countries are entities.
- Sports are entities.
- Information flowing between them is the popularity value.

	Country	Sport	Popularity
0	US	Cricket	10
1	US	Football	40
2	US	Basketball	50
3	India	Cricket	80
4	India	Football	15
5	India	Basketball	5
6	Australia	Cricket	50
7	Australia	Football	40
8	Australia	Basketball	10

While typical plots, like grouped bar plots, could be used here to understand the popularity distribution of country-wise sports, as shown below:



...in my opinion, however, Sankey diagrams stand out as a pretty cool and elegant alternative to represent such flow datasets:



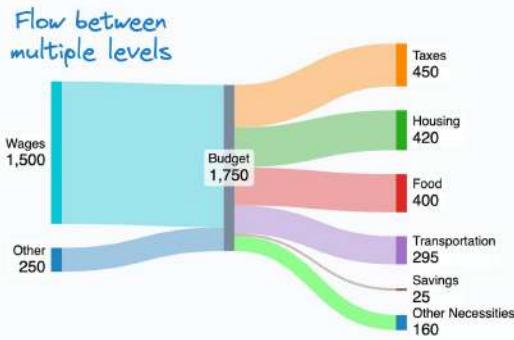
Their links are represented with arcs whose width is proportional to the value of the flow. This immensely simplifies the data analysis process. For instance, from the Sankey diagram above, one can quickly infer that:

- The most popular sport in India is Cricket.
- Basketball and Football are almost equally popular in the US.
- Basketball is hardly popular in India and England.
- England's most popular sport is Football.
- Cricket and Football are almost equally popular in Australia.
- Overall, Football is the most popular sport in this dataset.
- and many many more.

Imagine doing that by looking at the tabular data or a grouped bar chart.

- This process will be time-consuming.
- You may miss out on a few insights.
- The grouped bar chart can appear pretty cluttered and messy at times.

Of course, Sankey diagrams can have multiple levels as well, as shown below:



To determine when to use Sankey diagrams, see if the data involves any kind of flow of resources, energy, or information flow between multiple stages or entities. If yes, Sankey diagrams could be pretty valuable.

There are multiple ways to create Sankey diagrams:

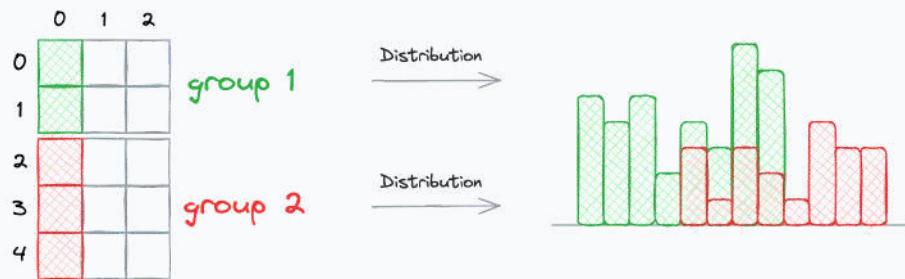
- To generate them programmatically, you may use the ipysankeywidget library: <https://github.com/ricklupton/ipysankeywidget>.
- If you prefer GUI, SankeyMATIC is a pretty cool, and easy-to-use tool to create Sankey diagrams which I often use: <https://bit.ly/3WdMANx>.

The screenshot shows the SankeyMATIC interface. On the left, there's a sidebar with 'Inputs:' containing node definitions like 'Australia [40] Football' and 'Australia [10] Basketball'. Below that is a code editor with comments explaining node creation and styling. The main workspace on the right displays a Sankey diagram with nodes for India, Australia, England, and US, connected by flows for Cricket, Football, and Basketball. The diagram is styled with various colors (red, purple, blue, green) and numerical values indicating flow magnitude.

Ridgeline Plots

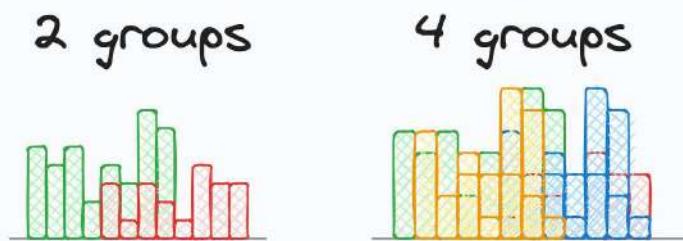
Understanding the distributional differences of distinct groups in a variable is quite useful in uncovering insights around:

- behavioral disparities,
- feature engineering,
- predictive modeling, and more.



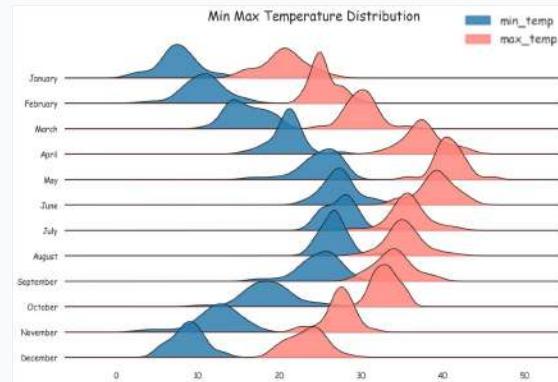
But in such situations, many data scientists tend to create group-level distribution plots (histograms or density plots) on a single axis and compare them.

While this is (somewhat) okay when there are limited groups, in the presence of many groups, it can create cluttered plots, which may not reveal many insights about distributional differences:

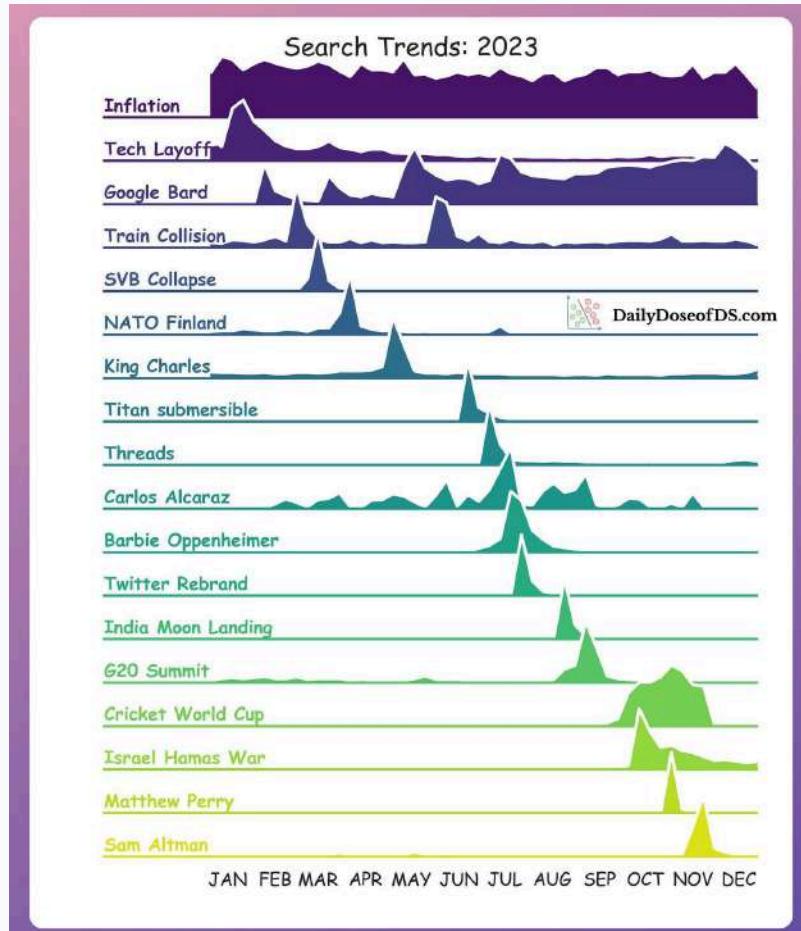


Ridgeline plots (shown below) are a pretty compact and elegant way to visualize the distribution of different variables (or categories of a variable).

More specifically, the vertical stacking on a common axis provides an easy comparison between groups and reveals many insights into the shape and variation of the distributions, which otherwise would be difficult to understand.



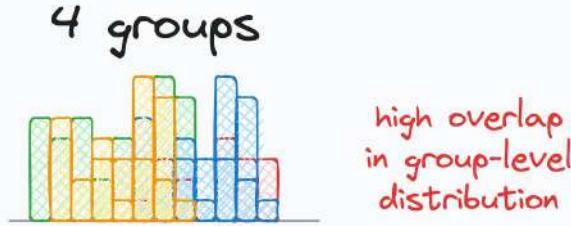
This allows us to compare the distributions of multiple groups side by side and understand how they differ. The image below is another classic example of Ridgeline plots. It depicts the search interest across various events that happened in 2023, and it's so easy to visualize.



While Seaborn provides a way to create Ridgeline plots, I have often found the [Joyplot library](#) to be pretty useful and easy to use.

When to consider Ridgeline plot?

Typically, creating a Ridgeline plot makes sense when the variable has anything above 3-4 groups. This is to avoid the overlap that might appear when visualizing them in a single plot:



Also, Ridgelines plots are relatively more useful when there is a clear pattern and/or ranking on the continuous variable plotted between groups like:

- monotonically increasing,
- monotonically decreasing,
- increasing then decreasing,
- decreasing then increasing, etc.

That is why the order in which you vertically stack the distribution of groups becomes quite important. For instance, consider the above “Search trends” plot again but with a random arrangement of groups:



I don't think I have to ask you which one is easier to visualize and understand the flow of events in 2023. So these were some points that will help you determine whether a Ridgeline plot will be a good fit for visualizing your data.

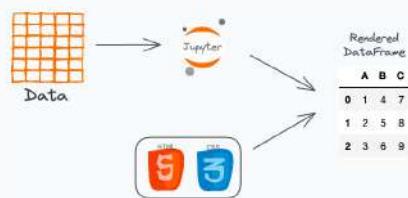
I created this notebook for you to get started with Ridgeline plots using JoyPy:
<https://bit.ly/3xMUpXj>.

Sparkline Plots

When doing any data analysis task in Jupyter, we mostly create standalone charts and visualizations. Of course, there's nothing wrong, but in this chapter, let me introduce you to another pretty cool and elegant way I often use to create visualizations in Jupyter Notebook.

Background

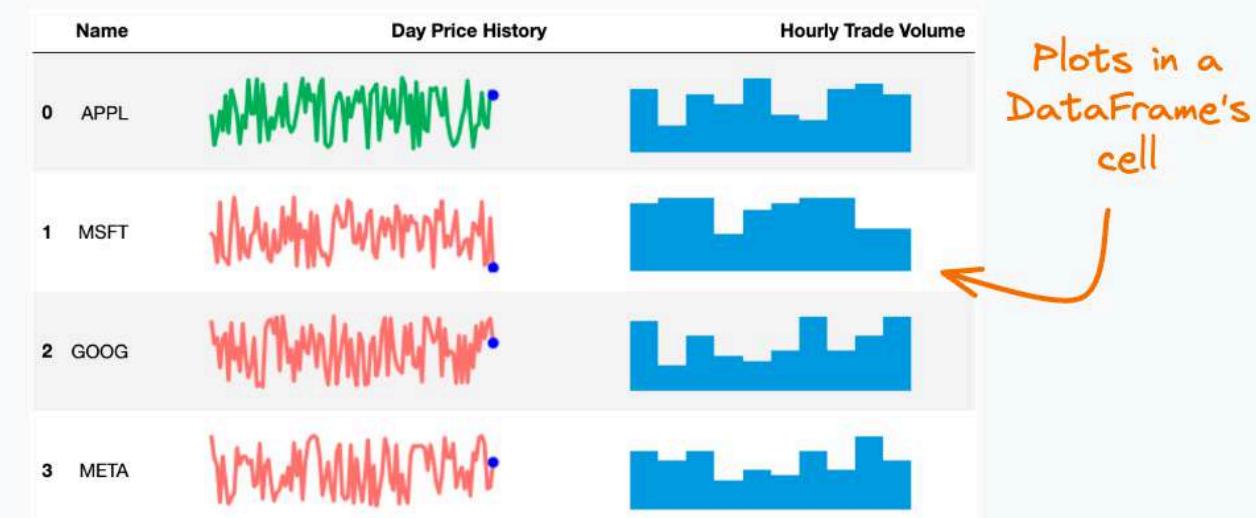
Whenever we display a DataFrame in Jupyter, it is rendered using HTML and CSS.



This means we can format its output just like any other web page. Moreover, when we create any plot in Jupyter, it is rendered as an image.

However, the same image can also be rendered within an HTML image tag in Jupyter.

And magic happens when we combine the above two ideas — rendering plots inside a DataFrame, as depicted below:



Formally, these compact plots are called Sparklines, and they provide a pretty elegant way to visualize data without taking up too much space.

Unlike traditional plots, we typically do not create axis ticks and labels in Sparklines.

One of the coolest ways I prefer using Sparklines is by adding them to a DataFrame's cell, as depicted in the image above.

A Jupyter Notebook on creating Sparkline plots is available here:
<https://bit.ly/4eVhVm>.

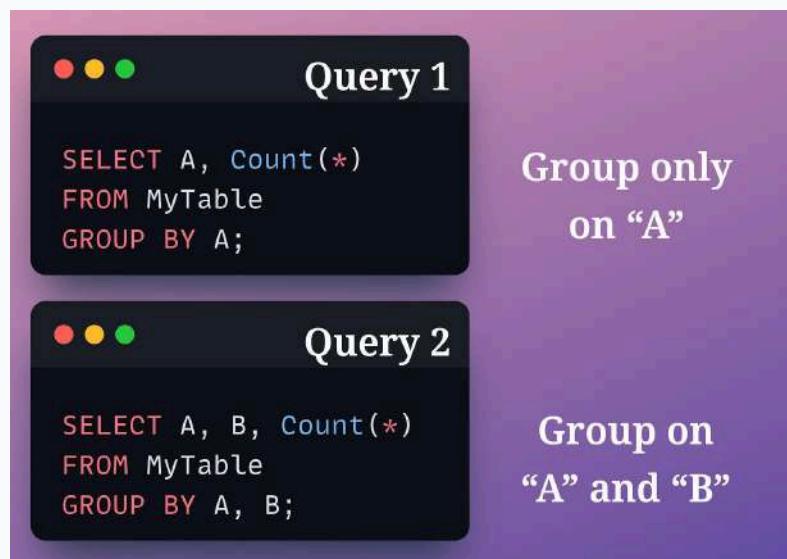
In this notebook, I have provided the step-by-step procedure to embed Sparklines in a DataFrame.

SQL

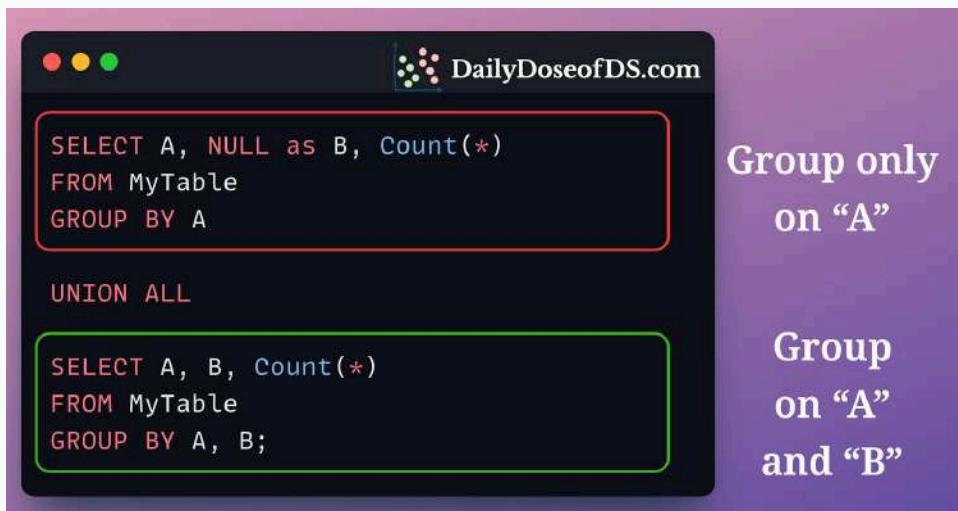
Grouping Sets, Rollup and Cube in SQL

A typical GroupBy query aggregates on just one set of columns. For instance:

- Grouping data on column “A” will require one query.
- Grouping data on columns “A” and “B” will require a separate query.



Next, if the two outputs must be gathered in a single table, we use UNION or UNION ALL (as needed).



But this is not efficient as it scans the same table twice. Instead, there are three ways to run multiple aggregations on the same table by scanning the table just once. This makes our query much more efficient. These are Grouping Sets, Rollup, and Cube, and their usage has been depicted in the image below:

Grouping Sets , Rollup and CUBE in SQL

Spends

Fruit	City	Amt
🍎🍎	NY	10
🍊🍊	CA	15
🍎🍎	CA	20
🍏🍏	NY	18



blog.DailyDoseofDS.com

1) GROUPING SETS

```
Select Fruit, City, SUM(Amt)
FROM Spends
Group By
GROUPING SETS ((Fruit), (City));
```

Fruit	City	Sum(Amt)
-	NY	28
	CA	35
Fruit group	-	30
	-	18
	-	15

2) ROLLUP

```
Select Fruit, City, SUM(Amt)
FROM Spends
Group By
ROLLUP (Fruit, City);
```

Fruit	City	Sum(Amt)
Fruit group	-	30
	-	15
	-	18
Fruit and city group	NY	10
	CA	15
	CA	20
	NY	18
City group	NY	28
	CA	35
	-	63
Total	-	63

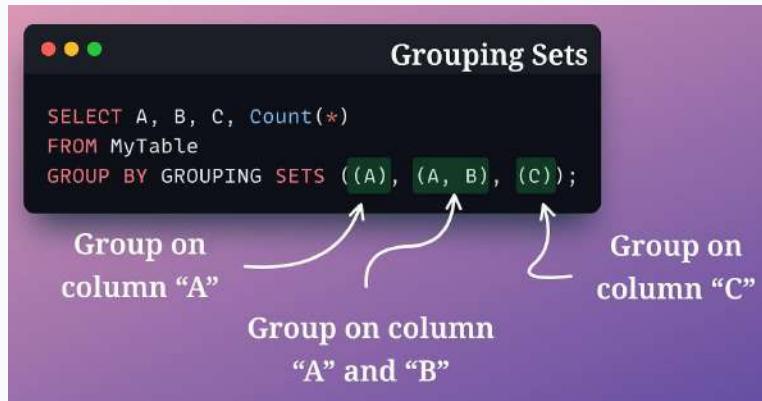
3) CUBE

```
Select Fruit, City, SUM(Amt)
FROM Spends
Group By
CUBE (Fruit, City);
```

Fruit	City	Sum(Amt)
Fruit group	-	30
	-	15
	-	18
Fruit and city group	NY	10
	CA	15
	CA	20
	NY	18
City group	NY	28
	CA	35
	-	63
Total	-	63

#1) Grouping sets

The GROUPING SETS clause allows us to define multiple groupings in a single query. Each grouping set defines a combination of columns by which the data is grouped.



Given the above query, here are the group aggregations that will be created:

- (A): Aggregated by A, counting all rows across all other columns.
- (A, B): Aggregated by (A, B), counting all rows across all other columns.
- (C): Aggregated by C, counting all rows across all other columns.

A demonstration is shown below:

1) GROUPING SETS

Fruit	City	Amt
🍎🍎	NY	10
🍊🍊	CA	15
🍎🍎	CA	20
🍏🍏	NY	18

Select Fruit, City, SUM(Amt)
FROM Spends
Group By
GROUPING SETS ((City), (Fruit));

Fruit	City	Sum(Amt)
-	NY	28
-	CA	35
🍎🍎	-	30
🍏🍏	-	18
🍊🍊	-	15

City group

Fruit group

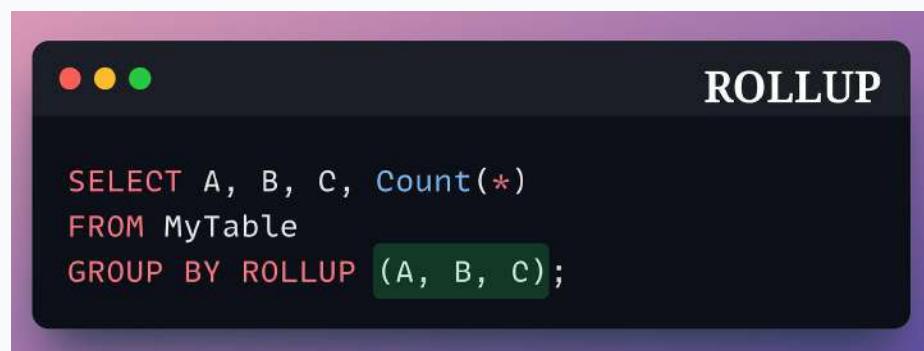
The above query:

- Generates an aggregation on City column.
- Generates another aggregation on Fruit column.

#2) ROLLUP

ROLLUP creates a result set that includes subtotals and a grand total in addition to the regular grouped results.

It does this by grouping the data at multiple levels of aggregation.



```
SELECT A, B, C, Count(*)
FROM MyTable
GROUP BY ROLLUP (A, B, C);
```

Given the above query, here are the group aggregations that will be created:

- (A): Subtotal for each A, aggregated across all B and C.
- (A, B): Subtotal for each A and B combination, aggregated across all C.
- (A, B, C): Regular group by all three columns.
- (): Grand total, aggregated across all A, B, and C.

A demonstration is shown below:

2) ROLLUP

Fruit	City	Amt
🍎🍎	NY	10
🍊🍊	CA	15
🍎🍎	CA	20
🍏🍏	NY	18

```
Select Fruit, City, SUM(Amt)
FROM Spends
Group By
ROLLUP (Fruit, City);
```

Fruit	City	Sum(Amt)
🍎🍎	-	30
🍊🍊	-	15
🍏🍏	-	18
🍎🍎	NY	10
🍊🍊	CA	15
🍎🍎	CA	20
🍏🍏	NY	18
Total		63

The above query:

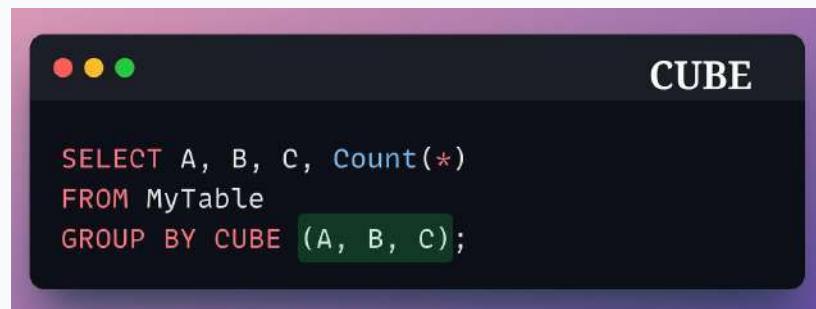
- Generates an aggregation on Fruit column.
- Generates another aggregation on (Fruit, City) column.
- Generates a grand total.

Unlike Grouping Sets, the order is important in ROLLUP.

More specifically, ROLLUP (A, B) will not be the same as ROLLUP (B, A).

#3) CUBE

Finally, CUBE creates a result set that includes all possible combinations of aggregations for the specified columns.



```
CUBE

SELECT A, B, C, Count(*)
FROM MyTable
GROUP BY CUBE (A, B, C);
```

Given the above query, here are the group aggregations that will be created:

- (A, B, C): Regular group by all three columns.
- (A, B): Subtotal for each A and B combination, aggregated across all C.
- (A, C): Subtotal for each A and C combination, aggregated across all B.
- (B, C): Subtotal for each B and C combination, aggregated across all A.
- (A): Subtotal for each A, aggregated across all B and C.
- (B): Subtotal for each B, aggregated across all A and C.
- (C): Subtotal for each C, aggregated across all A and B.
- (): Grand total, aggregated across all A, B, and C.

A demonstration is shown below:

3) CUBE

Fruit	City	Amt
🍎🍎	NY	10
🍊🍊	CA	15
🍎🍎	CA	20
🍏🍏	NY	18

Select Fruit, City, SUM(Amt)
FROM Spends
Group By
CUBE(Fruit, City);

Fruit	City	SUM(Amt)
🍎🍎	-	30
🍊🍊	-	15
🍏🍏	-	18

Fruit	City	SUM(Amt)
🍎🍎	NY	10
🍊🍊	CA	15
🍎🍎	CA	20
🍏🍏	NY	18

Fruit	City	SUM(Amt)
-	NY	28
-	CA	35

Total	-	-	63
-------	---	---	----

The above query:

- Generates an aggregation on Fruit column.
- Generates an aggregation on City column.
- Generates another aggregation on (Fruit, City) column.
- Generates a grand total.

Try it out by downloading this Jupyter Notebook: <https://bit.ly/4bx4SUX>.

Semi, Anti, and Natural Joins

We all have heard of LEFT JOIN, RIGHT JOIN, INNER JOIN, and OUTER JOIN, haven't we? These four are the most prevalent types of SQL joins. But there are more.

And in this chapter, I want to introduce you to three of them which I find pretty handy at times. These are – Semi Join, Anti Join and Natural Join.

Semi Join, Anti Join and Natural Join in SQL

 blog.DailyDoseofDS.com

SEMI JOIN

User_ID	User_Name
123	Bob
124	Alice
125	Charlie

User_ID	Order_ID
123	333
125	222
126	111

```
Select *
FROM USERS u
SEMI JOIN ORDERS o
on u.User_ID = o.User_ID;
```

User_ID	User_Name
123	Bob
125	Charlie

ANTI JOIN

User_ID	User_Name
123	Bob
124	Alice
125	Charlie

```
Select *
FROM USERS u
ANTI JOIN ORDERS o
on u.User_ID = o.User_ID;
```

User_ID	User_Name
124	Alice

NATURAL JOIN

User_ID	User_Name
123	Bob
124	Alice
125	Charlie

User_ID	Order_ID
123	333
125	222
126	111

```
Select *
FROM USERS u
NATURAL JOIN ORDERS o;
```

User_ID	User_Name
123	Bob
125	Charlie

No explicit join condition

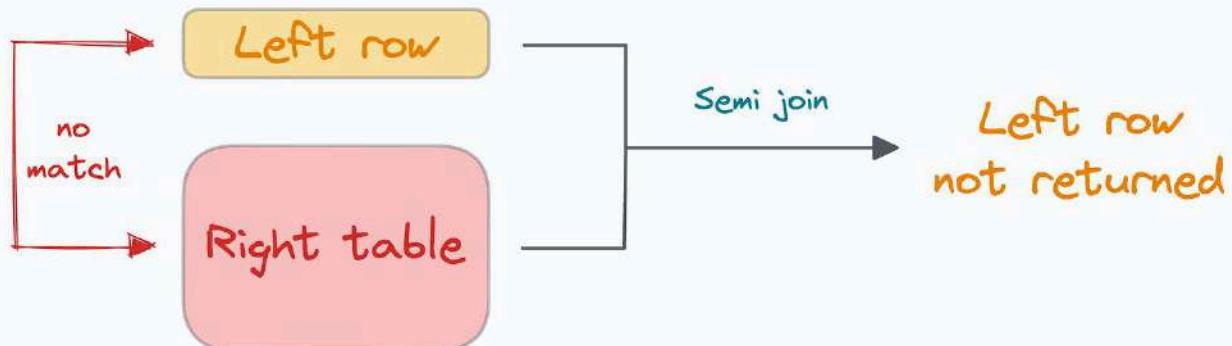
#1) Semi-Join

Semi-join appears quite similar to a left join, but there are three notable differences:

1. If the join condition between two rows is TRUE, columns from only the left table are returned. Compare this to the left join, which returns columns from both tables.



2. If a row in the left table has no match, then that row is not returned. In left join, however, all rows from the left table are returned irrespective of whether they have a match or not.

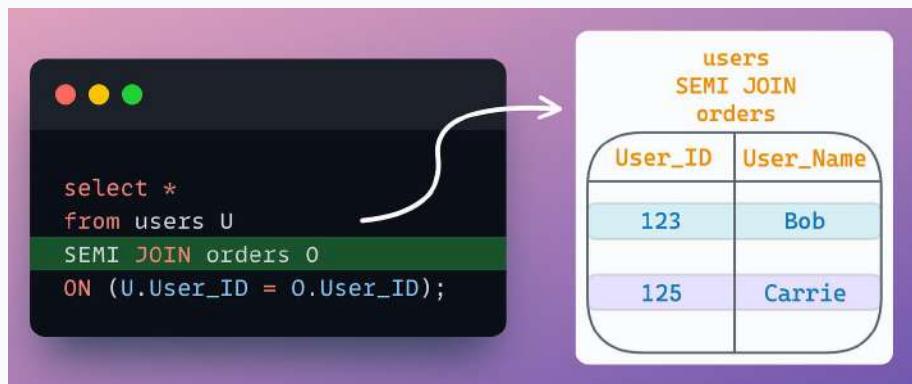


3. If a row in the left table has multiple matches, only one entry is returned. In left join, however, multiple matches are returned an equivalent number of times.

For instance, consider we have the following two tables:

users		orders	
User_ID	User_Name	User_ID	Order_ID
123	Bob	123	333
124	Alice	125	222
125	Carrie	126	111

Executing a semi-join, we get the following results:



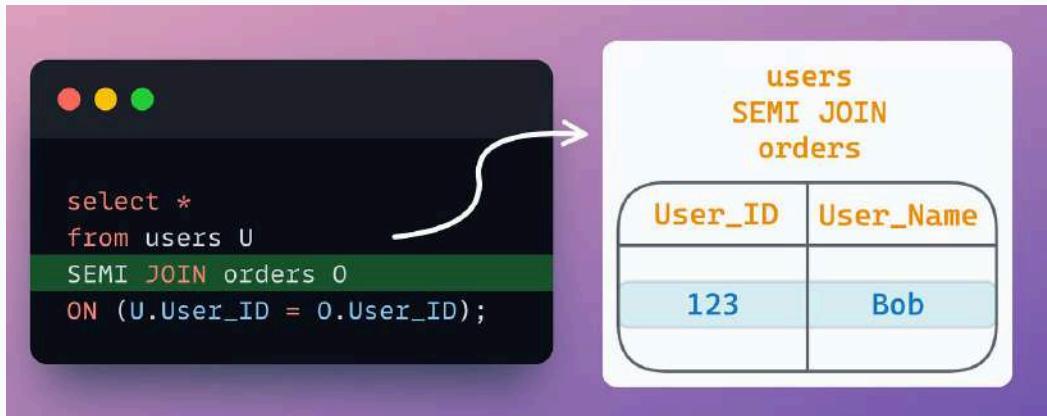
As depicted above, unlike left-join:

- It only returns columns from the left table.
- It only returns the matched rows from the left table.

If a record has multiple matches, like in this case:

users		orders	
User_ID	User_Name	User_ID	Order_ID
123	Bob	123	333
124	Alice	123	222
125	Carrie	126	111

...then we notice that semi-join only returns one record from the left table:



I find semi-join to be particularly useful when I only care about the existence of records in another table. Left join returns duplicates which are not of interest at that point.

#2) Anti-Join

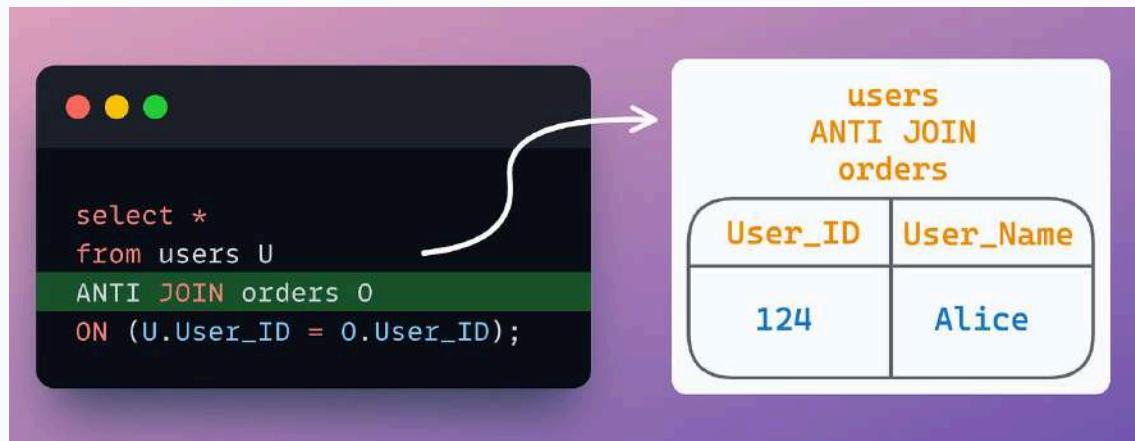
The rows discarded by the semi-join from the left table are the results of an anti-join. So, in a way, we can say that:

$$[\text{SEMI JOIN}] + [\text{ANTI JOIN}] = [\text{LEFT TABLE}]$$

Consider the above orders and users table again (one in which there were no multiple matches):

users		orders	
User_ID	User_Name	User_ID	Order_ID
123	Bob	123	333
124	Alice	125	222
125	Carrie	126	111

Executing an anti-join, we get the following results:



It is clear from the semi-join and anti-join results that:

$$[\text{SEMI JOIN}] + [\text{ANTI JOIN}] = [\text{LEFT TABLE}]$$

User_ID	User_Name
123	Bob
125	Carrie

User_ID	User_Name
124	Alice

User_ID	User_Name
123	Bob
124	Alice
125	Carrie

Of course, the order can be different. When I say “[SEMI JOIN] + [ANTI JOIN] = [LEFT TABLE]”, I mean the collection of all records.

I find anti-join to be particularly useful when I wish to know which records do not exist in another table.

#3) Natural Join

This one is similar to INNER JOIN, but there's no need to specify a join condition explicitly.

Instead, it automatically considers a join condition on ALL the matching column names.

Consider the users and orders table yet again:

users		orders	
User_ID	User_Name	User_ID	Order_ID
123	Bob	123	333
124	Alice	125	222
125	Carrie	126	111

Here, the User_ID column is present in both tables.

Executing a natural join, we get the following results:

The image shows a Jupyter Notebook cell with a black background and white text. The code is:

```
select *
from users U
NATURAL JOIN orders O;
```

Below the cell, the text "No explicit join condition" is displayed in white. To the right, the results of the query are shown in a table:

users
NATURAL JOIN
orders

User_ID	User_Name	Order_ID
123	Bob	333
125	Carrie	222

As depicted above, the results are similar to what we would get with INNER JOIN.

However, we did not have to explicitly specify a JOIN condition, which, admittedly, could be good or bad.

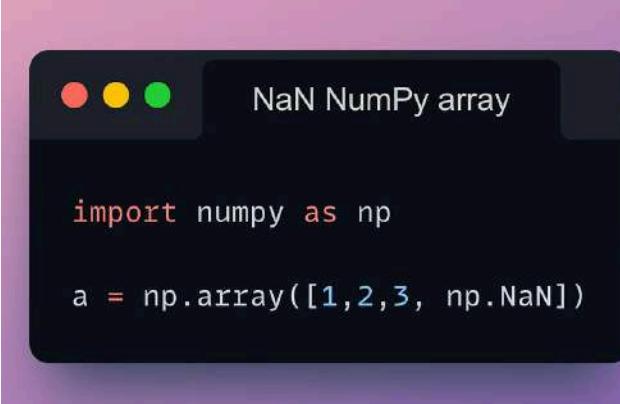
- It is good because it helps us write concise queries.
- It is bad because we are not explicit about the columns being joined.

These were three more types of SQL Joins, which I use at times to write concise and elegant SQL queries.

If you wish to experiment with what we discussed, download this Jupyter Notebook: <https://bit.ly/4cRGMWg>.

Use SQL "NOT IN" With Caution

In my experience, many unexpected errors in libraries/tools/languages can be attributed to the presence of missing values. For instance, consider a 1D NumPy array with NaN values.



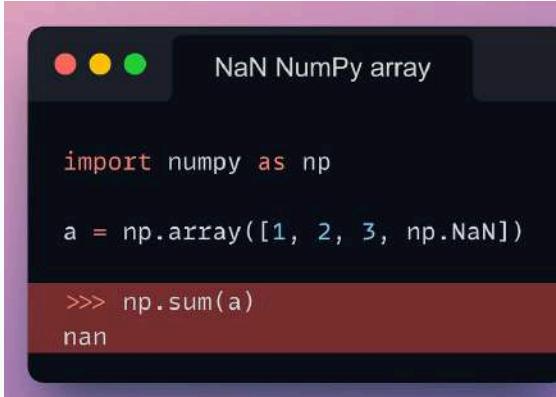
```
NaN NumPy array

import numpy as np

a = np.array([1,2,3, np.NaN])
```

**NumPy array
with missing
values**

When we aggregate this NumPy array to calculate, say, its sum, we get the following output:



```
NaN NumPy array

import numpy as np

a = np.array([1, 2, 3, np.NaN])

>>> np.sum(a)
nan
```

**NumPy returns
a NaN sum**

Strange, right? Although an output of 6 may have made more sense here, NumPy produces a NaN value instead.

Anyway, this is not the topic of this chapter, but I hope you get the point

A similar silent mistake can be found in SQL as well, specifically in the usage of “NOT IN” clause, which many SQL users are not aware of.

Let's understand!

Consider we have the following tables (students and names), in our database:

students				names	
	roll_no	first_name	last_name	grade	first_name
0	101	John	Doe	A	0 John
1	102	Emma	Smith	B	1 Peter
2	103	Michael	Johnson	C	2 Ethan
3	104	Sophia	Brown	A	
4	105	Ethan	Williams	B	

The task is to select records from students table where `first_name` is not in the names table. One way to do this is by using the `NOT IN` clause:

Filter using NOT IN clause

```
select *
from students
where first_name NOT IN (
    select first_name
    from names
);
```

	roll_no	first_name	last_name	grade
0	102	Emma	Smith	B
1	103	Michael	Johnson	C
2	104	Sophia	Brown	A

This provides the expected results as well.

Now, say our names table had a `NONE` value:

students				names	
	roll_no	first_name	last_name	grade	first_name
0	101	John	Doe	A	0 John
1	102	Emma	Smith	B	1 Peter
2	103	Michael	Johnson	C	2 Ethan
3	104	Sophia	Brown	A	3 None
4	105	Ethan	Williams	B	NULL value

If we run the above query again, we get no records this time:

Filter using NOT IN clause

```
select *
from students
where first_name NOT IN (
    select first_name
    from names
);
```

No records

roll_no	first_name	last_name	grade

On a side note, had we used the IN clause to select rows from the `students` table where `first_name` was in the `names` table, we notice that it works as expected:

Tables

	students				names	
	roll_no	first_name	last_name	grade	first_name	
0	101	John	Doe	A	0	John
1	102	Emma	Smith	B	1	Peter
2	103	Michael	Johnson	C	2	Ethan
3	104	Sophia	Brown	A	3	None
4	105	Ethan	Williams	B		NULL value

Filter using IN clause

```
select *
from students
where first_name IN (
    select first_name
    from names
);
```

Correct output

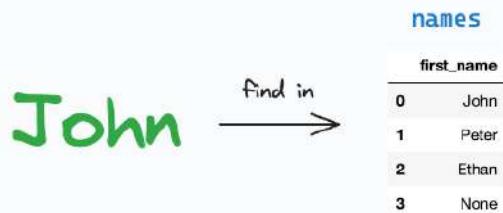
roll_no	first_name	last_name	grade
0	101	John	Doe
1	105	Ethan	Williams

That's strange, isn't it?

Why do we get no records with NOT IN?

The reason we get no records when we use NOT IN but the IN clause works as expected has to do with how these two clauses operate internally.

For simplicity, consider we are currently checking the first record (where `first_name` is "John") from the `students` table:



The **WHERE** clause needs a boolean value to determine whether a record must be filtered or not. When we use “IN”, this boolean value is evaluated using the OR operator as follows:

IN Clause

```
("John" = "John") OR
("John" = "Peter") OR
("John" = "Ethan") OR
("John" = None)
```

IN is based on the OR operator

If any condition is TRUE, the row gets filtered. However, when we use “NOT IN”, the boolean value is evaluated using the AND operator as follows:

NOT IN Clause

```
("John" ≠ "John") AND
("John" ≠ "Peter") AND
("John" ≠ "Ethan") AND
("John" ≠ None)
```

NOT IN is based on the AND operator

For the above expression to be TRUE, all individual conditions **MUST** be TRUE.

But the (JOHN != None) condition produces a conflict because, typically, this condition results in an UNKNOWN value. Thus, the entire expression evaluates to UNKNOWN — producing no records. This happens for every record in the students table. As a result, the query results in no records:

Filter using NOT IN clause

```
select *
from students
where first_name NOT IN (
    select first_name
    from names
);
```

roll_no	first_name	last_name	grade
No records			

What to do instead?

There are many ways to avoid this. As our task is to select records from `students` table where `first_name` is not in the `names` table, we can:

- Filter out the NULL values in the sub-query:

```
select *
from students
where first_name NOT IN (
    select first_name
    from names
    where first_name is NOT NULL
);
```

Correct output

	roll_no	first_name	last_name	grade
0	102	Emma	Smith	B
1	103	Michael	Johnson	C
2	104	Sophia	Brown	A

- Use Anti Joins: Anti join returns only those rows from the left table where no match is found in the right table. This is precisely what we need in our case and it is implemented below:

Use ANTI Join

```
select *
from students s
ANTI JOIN names n
ON (s.first_name = n.first_name);
```

Correct output

	roll_no	first_name	last_name	grade
0	102	Emma	Smith	B
1	103	Michael	Johnson	C
2	104	Sophia	Brown	A

At times, such mistakes can take some serious time to debug if you are not aware of them beforehand.

If you wish to experiment with what we discussed about the NOT IN clause, download this Jupyter Notebook: <https://bit.ly/3xVOIXF>.

Python OOP

Getters and Setters

Dot notation provides a concise and elegant way to access and modify an object's attributes.

A screenshot of a Python code editor window titled "myclass.py". The code defines a class "Square" with an __init__ method that initializes a "side" attribute. A variable "s" is created as an instance of "Square" with a side length of 10. The line "s.side" is highlighted in green, and a white arrow points from the word "Dot notation" to this line. The code is as follows:

```
class Square:
    def __init__(self, length):
        self.side = length

>>> s = Square(10)
```

Dot notation

```
>>> s.side
```

Yet, with dot notation, we can not validate the updates made to an attribute. This means we can assign invalid values to an instance's attributes, as shown below:

A screenshot of a Python code editor window titled "myclass.py". The code creates an instance "s" of the "Square" class with a side length of 10. The line "s.side = -2" is highlighted in green, and a white arrow points from the text "No error" to this line. The code is as follows:

```
>>> s = Square(10)
```

No error

```
>>> s.side = -2
```

```
>>> s.side
```

One common way to avoid this is by defining a setter (`set_side()`), which validates the assignment step.

```

class Square:
    def __init__(self, length):
        self._side = length

    def set_side(self, length):
        if length < 0:
            raise ValueError("Side cannot be negative")
        else:
            self._side = length

>>> s.set_side(-2)
ValueError: Side cannot be negative

```

But explicitly invoking a setter method isn't as elegant as dot notation, is it?

Ideally, we would want to use dot notation and still apply those validation checks. The `@property` decorator in Python can help. Here's how we can use it here.

First, define a getter as follows:

Declare a method with the attribute's name.

There's no need to specify any parameters for this method.

Decorate it with the `@property` decorator.

```

class Square:
    ...

    @property
    def side(self):
        return self._side

```

Next, define a setter as follows:



```

class Square:
    ...
    @side.setter
    def side(self, length):
        if length<0:
            raise ValueError("Side cannot be negative")
        else:
            self._side = length

```

- Declare a method with the attribute's name.
- Specify the parameter you want to update the attribute with.
- Write the conditions as you usually would in any other setter method.
- Decorate it with the `@attribute-name.setter` decorator.

Now, you can use the dot notation while still having validation checks in place.

This approach offers both the validation and control of explicit setters and getters and the elegance of dot notations.



```

class Square:
    def __init__(self, length):
        self._side = length

    @property
    def side(self):
        return self._side

    @side.setter
    def side(self, length):
        if length<0:
            raise ValueError("Side cannot be negative")
        else:
            self._side = length

```

Raises errors during assignment

```

>>> s = Square(10)
>>> s.side # Getter
10

>>> s.side = -2 # Setter (with dot)
ValueError: Side cannot be negative

```

Descriptors in Python

In this chapter, let's continue our discussion on the above topic and discuss a limitation of the above approach. Moving on, we shall see how Descriptors in Python provide a much more elegant way of setting and getting values.

Limitations of @property decorator

Consider the above class implementation again:

```

class Square:
    def __init__(self, length):
        self._side = length

    @property
    def side(self):
        return self._side

    @side.setter
    def side(self, length):
        if length<0:
            raise ValueError("Side cannot be negative")
        else:
            self._side = length

```

The biggest issue here is that we must define a getter and setter for every instance-level attribute.

So what if our class has, say, 3 such attributes, and all must be positive?

Of course, we will have 3 getters and 3 setters, which makes the overall implementation long, messy, and redundant.

```

class DummyClass:

    def __init__(self, num1, num2, num3):
        self._number1 = num1
        self._number2 = num2
        self._number3 = num3

    @property
    def number1(self):
        return self._number1

    @number1.setter
    def number1(self, value):
        if value<0:
            raise ValueError("Enter positive value")
        self._number1 = value

    @property
    def number2(self):
        ...

    @number2.setter
    def number2(self, value):
        ...

    @property
    def number3(self):
        ...

    @number3.setter
    def number3(self, value):
        ...

```

**Getter and Setter
for the attribute
"number1"**

**Getter and Setter
for the attribute
"number2"**

**Getter and Setter
for the attribute
"number3"**

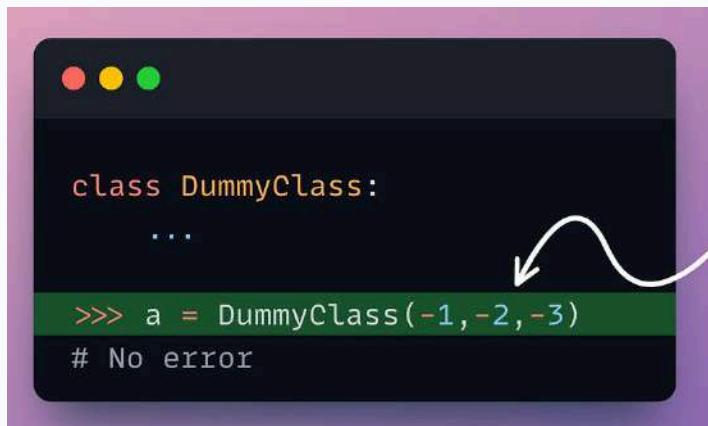
There's redundancy because every setter method will have almost the same lines of code (the if statements for validation).

Also, if you think about it, the getter methods are somewhat redundant and unnecessary too, as they just return an attribute.

If that is clear, there's one more issue with the above implementation.

Recall what I mentioned earlier: "Our class will have 3 such instance-level attributes, and all must be positive?"

See what happens when we create an object with an invalid input:



The screenshot shows a terminal window with a dark background. At the top, there are three colored window control buttons (red, yellow, green). Below them, the terminal displays Python code and its output. The code defines a class `DummyClass` with a constructor that takes three parameters. A call to the constructor with negative values is shown, followed by a comment `# No error`. A white arrow points from the text "Python raises no error" to the word `a` in the code line `>>> a = DummyClass(-1,-2,-3)`.

Python raises
no error

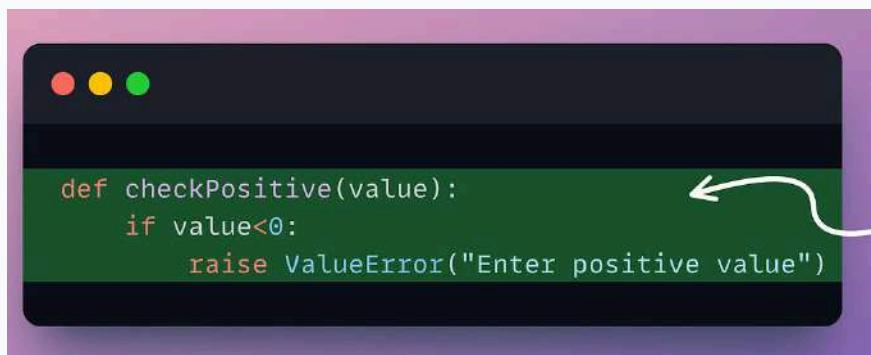
```
class DummyClass:  
    ...  
>>> a = DummyClass(-1,-2,-3)  
# No error
```

As depicted above, Python does not raise any error, when ideally, it should.

An imperfect solution

One common way programmers try to eliminate redundancy is by defining explicit validation functions.

For instance, we can define a function that just validates the value received, as demonstrated below:



The screenshot shows a terminal window with a dark background. At the top, there are three colored window control buttons (red, yellow, green). Below them, the terminal displays Python code. It defines a function `checkPositive` that takes a single parameter `value`. Inside the function, it checks if `value` is less than 0. If so, it raises a `ValueError` with the message "Enter positive value". A white arrow points from the text "Input validation method" to the word `value` in the code line `if value<0:`.

Input
validation
method

```
def checkPositive(value):  
    if value<0:  
        raise ValueError("Enter positive value")
```

Next, we can invoke this method wherever needed:

```

class DummyClass:

    def __init__(self, num1, num2, num3):
        checkPositive(num1)
        checkPositive(num2)
        checkPositive(num3)
        self._number1 = num1
        self._number2 = num2
        self._number3 = num3

    @property
    def number1(self):
        return self._number1

    @number1.setter
    def number1(self, value):
        checkPositive(value)
        self._number1 = value

    @property
    def number2(self):
        return self._number2

    @number2.setter
    def number2(self, value):
        checkPositive(value)
        self._number2 = value

    @property
    def number3(self):
        return self._number3

    @number3.setter
    def number3(self, value):
        checkPositive(value)
        self._number3 = value

```

Check value in every method by invoking the function

But this does not solve the problem either:

We still have explicit and redundant function calls in each setter method.

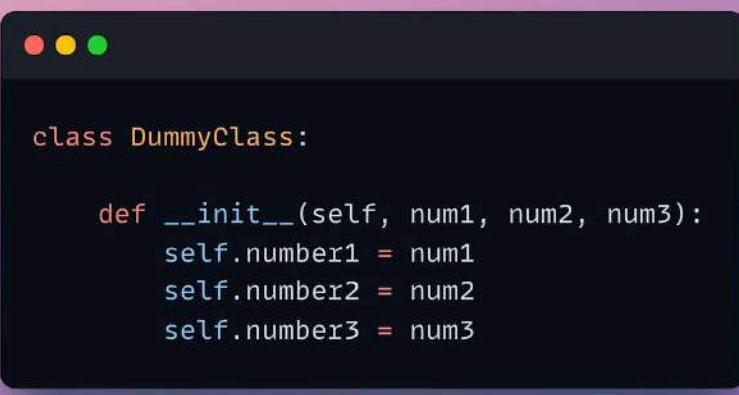
All getter methods still do the same thing and have high redundancy.

And most importantly, the `__init__` method is now messed up with multiple function calls.

Descriptors

Simply put, Descriptors are objects with methods (like `__get__`, `__set__`, etc.) that are used to manage access to the attributes of another class. So, every descriptor object is assigned to only one attribute of another class.

And just to be clear, this “another class” is the class we are primarily interested in — the `DummyClass` we saw earlier, for instance.



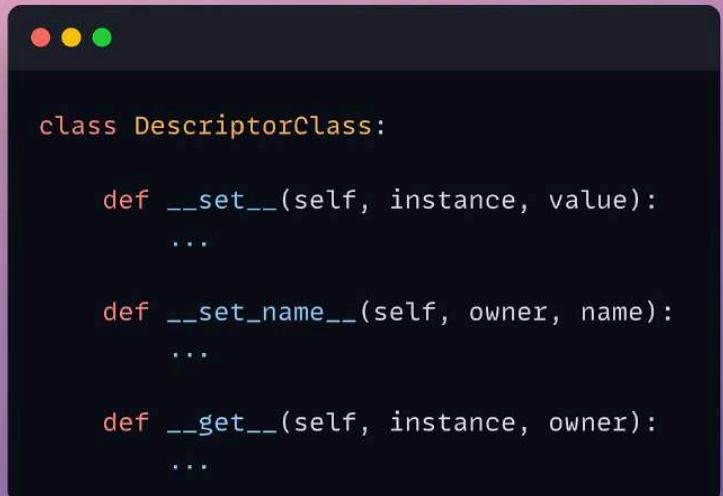
```
class DummyClass:

    def __init__(self, num1, num2, num3):
        self.number1 = num1
        self.number2 = num2
        self.number3 = num3
```

Every attribute will be assigned to one Descriptor object

- The attribute `number1` → gets its own descriptor.
- The attribute `number2` → gets its own descriptor.
- The attribute `number3` → gets its own descriptor.

A typical Descriptor class is implemented with three methods, as shown below:



```
class DescriptorClass:

    def __set__(self, instance, value):
        ...

    def __set_name__(self, owner, name):
        ...

    def __get__(self, instance, owner):
        ...
```

Descriptor class definition

The `__set__` method is called when the attribute is assigned a new value. We can define the custom checks here.

The `__set_name__` method is called when the descriptor object is assigned to a class attribute. It allows the descriptor to keep track of the name of the attribute it's assigned to within the class.

The `__get__` method is called when the attribute is accessed.

Also:



A screenshot of a Mac OS X terminal window. The window title is "Descriptor class definition". The terminal contains the following Python code:

```
class DescriptorClass:

    def __set__(self, instance, value):
        ...

    def __set_name__(self, owner, name):
        ...

    def __get__(self, instance, owner):
        ...
```

The `instance` parameter refers to the object of the desired class — `DummyClass()`.

The `owner` parameter is the desired class itself — `DummyClass`.

The `value` parameter is the value being assigned to an attribute of the desired class.

The `name` parameter is the name of the attribute.

If it's unclear, let me give you a simple demonstration.

Consider this Descriptor class:

```
class DescriptorClass:

    def __set__(self, instance, value):
        if value < 0:
            raise ValueError("Enter positive value")
        instance.__dict__[self.attribute_name] = value

    def __set_name__(self, owner, name):
        self.attribute_name = name

    def __get__(self, instance, owner):
        return instance.__dict__[self.attribute_name]
```

Setter checks

**Assign name
to attribute**

Getter method

I'll explain this implementation shortly, but before that, let's consider its usage, which is demonstrated below:

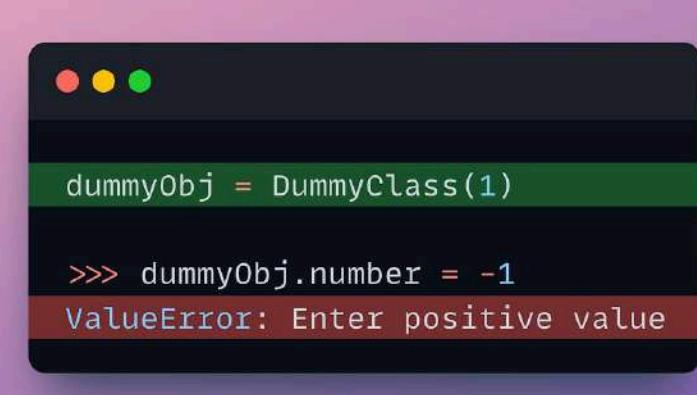
```
1 class DummyClass:
2
3     number = DescriptorClass()
4
5     def __init__(self, num):
6         self.number = num
```

**Create
descriptor
object**

Now, let's go back to the DescriptorClass implementation:

- `__set_name__(self, owner, name)`: This method is called when the descriptor is assigned to a class attribute (line 3). It saves the name of the attribute in the descriptor for later use.
- `__set__(self, instance, value)`: When a value is assigned to the attribute (line 6), this method is called. It raises an error if the value is negative. Otherwise, it stores the value in the instance's dictionary under the attribute name we defined earlier.
- `__get__(self, instance, owner)`: When the attribute is accessed, this method is called. It returns the value from the instance's dictionary.

Now, see how this solution smartly solves all the problems we discussed earlier. Let's create an object of the DummyClass:

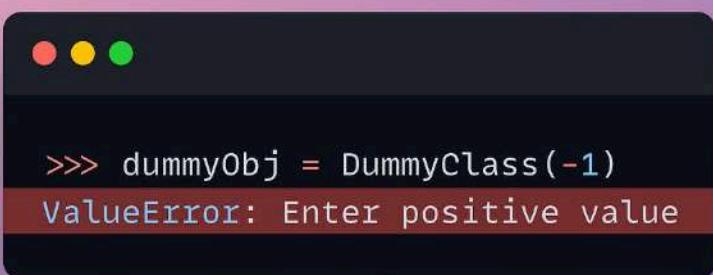


```
dummyObj = DummyClass(1)

>>> dummyObj.number = -1
ValueError: Enter positive value
```

**Error raised
on invalid
assignment**

As depicted above, assigning an invalid value to the attribute raises an error. Next, let's see what happens when the attribute specified during the initialization is invalid:



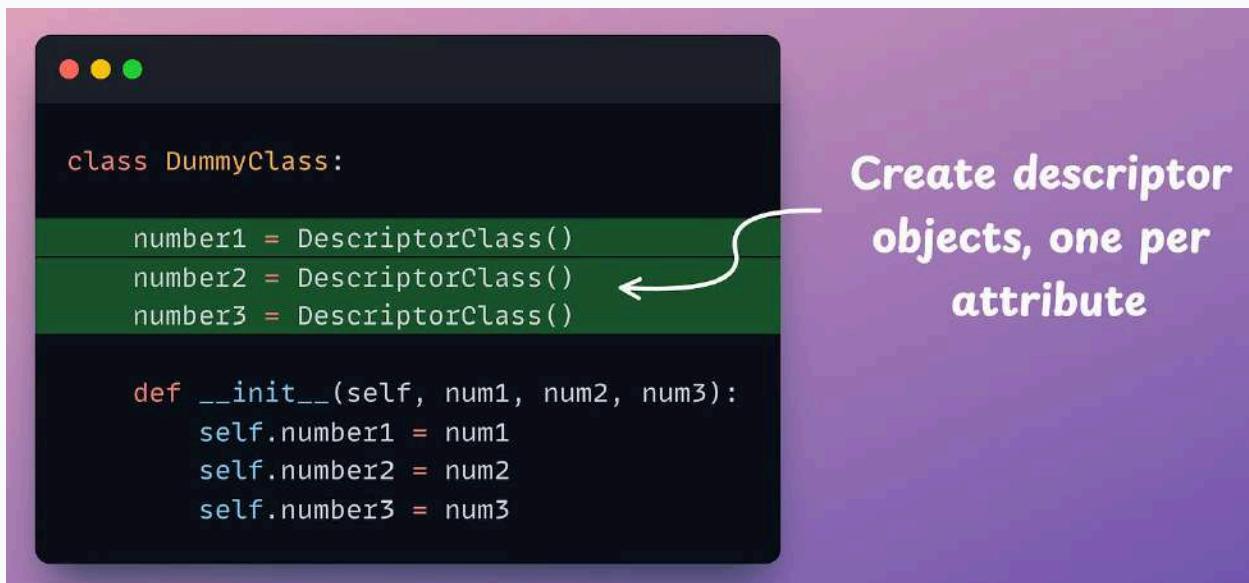
```
>>> dummyObj = DummyClass(-1)
ValueError: Enter positive value
```

**Error raised
when creating
object**

Great! It validates the initialization too.

Here, recall that we never defined any explicit checks in the `__init__` method, which is super cool.

Moving on, let's define multiple attributes in the `DummyClass` now:



```

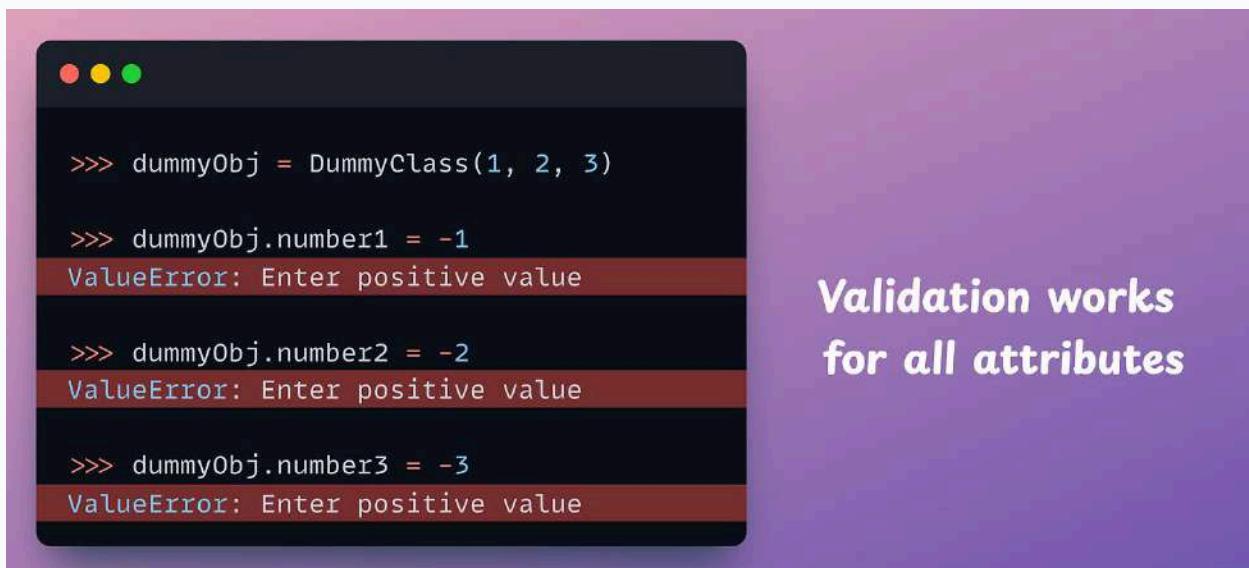
class DummyClass:

    number1 = DescriptorClass()
    number2 = DescriptorClass()
    number3 = DescriptorClass()

    def __init__(self, num1, num2, num3):
        self.number1 = num1
        self.number2 = num2
        self.number3 = num3

```

Creating an object and setting an invalid value for any of the attributes raises an error:



```

>>> dummyObj = DummyClass(1, 2, 3)

>>> dummyObj.number1 = -1
ValueError: Enter positive value

>>> dummyObj.number2 = -2
ValueError: Enter positive value

>>> dummyObj.number3 = -3
ValueError: Enter positive value

```

Works seamlessly!

Recall that we never defined multiple getters and setters for each attribute individually, like we did with the `@property` decorator earlier.

This is great, isn't it?

I find descriptors to be massively helpful in reducing work and code redundancy while also making the entire implementation much more elegant.

If you want to try them out, I prepared this notebook for you to get started:

<https://bit.ly/4cqsyM6>.

20 Most Common Magic Methods

Here are 20 most common magic methods used in Python OOP:

Magic Method	Syntax	Usage/Description
<code>__new__</code>	<code>__new__(cls, *args, **kwargs):</code>	Invoked before <code>__init__</code> to allocate memory to object
<code>__init__</code>	<code>__init__(self, *args, **kwargs):</code>	Invoked after <code>__new__</code> to initialise the object
<code>__str__</code>	<code>__str__(self):</code>	Invoked when <code>str(obj)</code> or <code>print(obj)</code> is used
<code>__int__</code>	<code>__int__(self):</code>	Invoked when <code>int(obj)</code> is used
<code>__len__</code>	<code>__len__(self):</code>	Invoked when <code>len(obj)</code> is used
<code>__call__</code>	<code>__call__(self, *args, **kwargs):</code>	Invoked when class object is called as a function: <code>obj()</code>
<code>__getitem__</code>	<code>__getitem__(self, key):</code>	Invoked when object is indexed: <code>obj[key]</code>
<code>__setitem__</code>	<code>__setitem__(self, key, value):</code>	Invoked when object is indexed and value is set: <code>obj[key]=value</code>
<code>__delitem__</code>	<code>__delitem__(self, key):</code>	Invoked when object's index is deleted: <code>del obj[key]</code>
<code>__contains__</code>	<code>__contains__(self, item):</code>	Invoked when the <code>in</code> operator is used: <code>item in obj</code>
<code>__bool__</code>	<code>__bool__(self):</code>	Invoked when object is used in boolean context: <code>if obj</code> or <code>bool(obj)</code>
<code>__iter__</code>	<code>__iter__(self):</code>	Invoked when object is iterated: <code>for x in obj</code>
<code>__eq__</code>	<code>__eq__(self, other):</code>	Invoked when <code>==</code> operator is used to compare two objects: <code>obj1 == obj2</code>
<code>__ne__</code>	<code>__ne__(self, other):</code>	Invoked when <code>!=</code> operator is used to compare two objects: <code>obj1 != obj2</code>
<code>__gt__</code>	<code>__gt__(self, other):</code>	Invoked when <code>></code> operator is used to compare two objects: <code>obj1 > obj2</code>
<code>__add__</code>	<code>__add__(self, other):</code>	Invoked when two objects are added: <code>obj1 + obj2</code>
<code>__mul__</code>	<code>__mul__(self, other):</code>	Invoked when two objects are multiplied: <code>obj1 * obj2</code>
<code>__abs__</code>	<code>__abs__(self):</code>	Invoked to compute absolute value of object: <code>abs(obj)</code>
<code>__neg__</code>	<code>__neg__(self):</code>	Invoked when unary operator <code>-</code> is used on an object: <code>-obj</code>
<code>__invert__</code>	<code>__invert__(self):</code>	Invoked when <code>~(tilde)</code> operator is used to invert an object: <code>~obj</code>

In my experience, these are possibly the only 20 magic methods you would ever need in most Python projects utilizing OOP.

Syntactically, they are prefixed and suffixed with double underscores, such as `__len__`, `__str__`, and many more. That is why they are also called “Dunder methods” — short for Double UNDERscore.

Memory Efficient Class Objects using Slots

When we define a class in Python, it is possible to dynamically add new attributes to its objects during run-time.

For instance, consider the following Python class here.

```

class Student:

    def __init__(self, name, age):
        self.name = name
        self.age = age

    def __repr__(self):
        ...

StudentA = Student("Mark", 19)
StudentB = Student("Joey", 18)

```

Python class with two instance-level attributes

Here, it is perfectly legal to add new attributes to any object at run-time, as shown below:

```

StudentA = Student("Mark", 19)

>>> StudentA.friend = "Julie"

```

Add new attributes

However, this is not always recommended because:

- It may lead to bugs if the code assumes all class instances will always have the same attributes.
- It makes it difficult to debug code when objects keep on accumulating new attributes.
- It leads to a conflicting schema, etc.

```
>>> print(StudentA)
Student(name="Mark", age=19, friend="Julie")

>>> print(StudentB)
Student(name="Joey", age=18)
```

Conflicting schema

Can we restrict this dynamicity?

Of course we can!

Defining a slotted class helps us achieve this. Simply put, it allows us to fix the instance-level attributes a class object can ever possess.

A slotted class is declared as follows:

```
class Student:

    __slots__ = ["name", "age"]

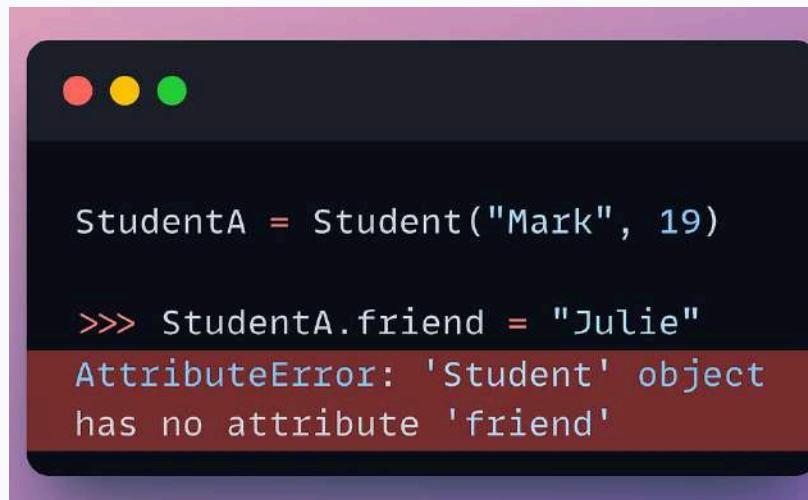
    def __init__(self, name, age):
        self.name = name
        self.age = age
```

Declare __slots__ attribute

Define a `__slots__` attribute with all attributes an object may possess.

Done!

Now, if we try to add a new attribute at run-time, it raises an error:



The screenshot shows a Jupyter Notebook cell with the following code and output:

```
StudentA = Student("Mark", 19)

>>> StudentA.friend = "Julie"
AttributeError: 'Student' object
has no attribute 'friend'
```

A white curly brace on the right side of the slide points to the word "AttributeError" in the output.

Adding new attribute raises an error

Declaring a slotted class has many advantages as well. For instance, it can help us avoid typical code typos. Say we want to change the name of `studentA`.

Here, instead of referring to the attribute name (with a lowercase n), we mistakenly wrote `Name` (with an uppercase n). A normal unslotted class will not raise an error, as shown below:



The screenshot shows a Jupyter Notebook cell with the following code and output:

```
class Student:

    def __init__(self, name, age):
        self.name = name
        self.age = age

StudentA = Student("Mark", 19)

>>> StudentA.Name = "Peter"
```

A white curly brace on the right side of the slide points to the word "Name" in the last line of the code.

Unslotted class

No error

But a slotted class will catch this typo:

The screenshot shows a Python terminal window. The code defines a class `Student` with slots for `"name"` and `"age"`. It includes an `__init__` method to initialize these attributes. A call to `StudentA = Student("Mark", 19)` is followed by an attempt to set the attribute `Name` to "Peter". The terminal returns an `AttributeError`: `'Student' object has no attribute 'Name'`. A curly brace on the right side of the terminal window is labeled **Slotted class**, and another one below it is labeled **Typo caught**.

```

class Student:
    __slots__ = ["name", "age"]

    def __init__(self, name, age):
        self.name = name
        self.age = age

StudentA = Student("Mark", 19)

>>> StudentA.Name = "Peter"
AttributeError: 'Student' object
has no attribute 'Name'

```

Declaring a slotted class provides memory advantages as well. As shown below, an object of the slotted class consumes 2.5x less memory than an object of a normal class:

The screenshot shows a Python terminal window. It creates two objects: `studentA` of type `Student` and `studentB` of type `StudentWithoutSlots`, both initialized with the same parameters. It then uses the `pympler` library to measure their memory sizes. The output shows that `student1` (slotted) is 136 bytes, while `student2` (non-slotted) is 352 bytes. A curly brace on the right side of the terminal window is labeled **Slotted class object consumes 2.5x less memory**.

```

studentA = Student("Mark", 24)
studentB = StudentWithoutSlots("Mark", 24)

import pympler

>>> pympler.asizeof.asizeof(student1)
136

>>> pympler.asizeof.asizeof(student2)
352

```

Why, you may wonder?

Typically, Python creates a dictionary for every object that maps attributes to their values. As a dictionary is mutable, this is precisely what allows us to add/delete attributes dynamically. But this introduces memory overheads as Python always tries to keep room for new attributes that may get added at some point.

In case of slotted class, however, Python can do away with this dictionary:

```
class StudentWithoutSlots:
    def __init__(self, name, age):
        self.name = name
        self.age = age

StudentB = StudentWithoutSlots("Mark", 19)

print(StudentB.__dict__)
{'name': 'Mark', 'age': 19} Python Dictionary
```

```
class Student:
    __slots__ = ["name", "age", "dummy"]

    def __init__(self, name, age):
        self.name = name
        self.age = age

StudentA = Student("Mark", 19)

print(StudentA.__dict__)
```

No Python Dictionary

```
AttributeError                                     Traceback (most recent call last)
Input In [11], in <cell line: 11>()
    7     self.age = age
    9 StudentA = Student("Mark", 19)
--> 11 print(StudentA.__dict__)

AttributeError: 'Student' object has no attribute '__dict__'
```

As a result, the object consumes less memory. As a departing note, if you don't want to dynamically add new attributes to an object, it is better to create a slotted class.

In fact, even if you know the attributes that a class object will ever possess, but some of these attributes are not available during the object's initialization, you can still declare them in the `__slots__` class attribute and assign a value to them later in the program whenever they are available.

Why Don't We Invoke `model.forward()` in PyTorch?

In PyTorch, the forward pass is implemented in the `forward()` method, as demonstrated below:



```

import torch.nn as nn

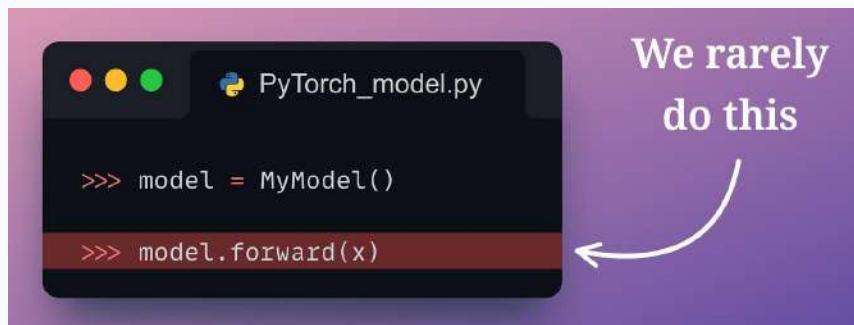
class MyModel(nn.Module):
    def __init__(self):
        super(MyModel, self).__init__()
        self.fc = nn.Linear(10, 5)

    def forward(self, x):
        return self.fc(x)
  
```

Method for forward pass

A white curly arrow points from the text "Method for forward pass" to the `def forward(self, x):` line in the code.

Here, have you ever wondered that when we want to run the forward pass, we rarely invoke this `forward()` method:



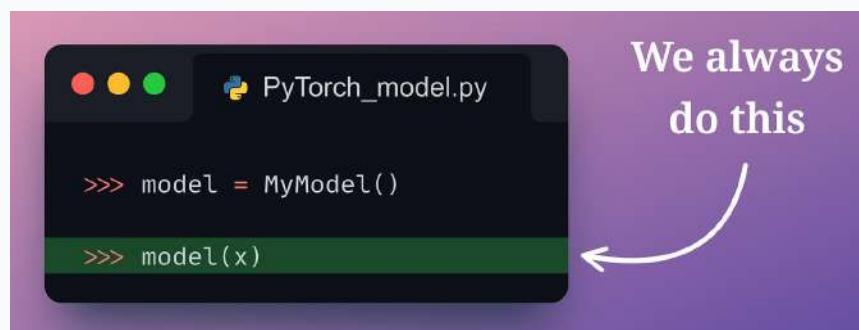
```

>>> model = MyModel()
>>> model.forward(x)
  
```

We rarely do this

A white curly arrow points from the text "We rarely do this" to the `model.forward(x)` line in the code.

Instead, we always invoke the model (a class object) as demonstrated below, as it was a function:



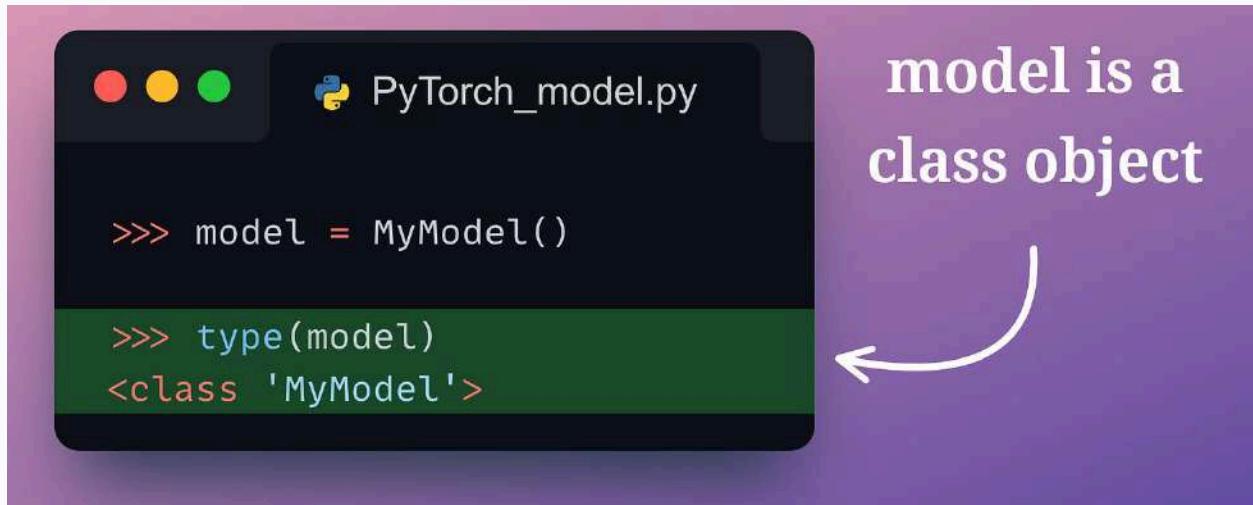
```

>>> model = MyModel()
>>> model(x)
  
```

We always do this

A white curly arrow points from the text "We always do this" to the `model(x)` line in the code.

We can also verify that model is a class object:



```
>>> model = MyModel()
>>> type(model)
<class 'MyModel'>
```

model is a class object

How can a class object be invoked like a function and what are we missing here?

Let's understand!

A simple example

Consider we want to evaluate the following quadratic:

$$f(x) = ax^2 + bx + c$$

One way is to define a method that accepts the input and returns the value of the quadratic, as shown below:



```
Quadratic class

class Quadratic:
    def __init__(self, a, b, c):
        self.a = a
        self.b = b
        self.c = c

    def evaluate(self, x):
        return self.a*(x**2) +
               self.b*x + self.c
```

```
>>> my_quad = Quadratic(1, 2, 3)
>>> my_quad.evaluate(1)
6
```

Invoke method

Of course, there is nothing wrong with this approach.

But there is one smart and elegant way of doing this in Python. Instead of explicitly invoking a method, we can define the `__call__()` magic method.

This magic method allows you to define the behavior of the class object when it is invoked like a function (like this: `object()`).

Let's rename the `evaluate()` method to `__call__()`.

The diagram illustrates the transition from a standard class definition to one that supports direct invocation. On the left, a code editor shows a Python class named `Quadratic` with an `__init__` method and an `evaluate` method. A callout bubble points to the `evaluate` method with the text "Change to `__call__()`". On the right, a terminal window shows the creation of a `Quadratic` object and its direct invocation with the value `1`, which returns the value `6`. A callout bubble points to the terminal with the text "Invoke object directly".

```

class Quadratic:
    def __init__(self, a, b, c):
        self.a = a
        self.b = b
        self.c = c

    def evaluate(self, x):
        return self.a*(x**2) + self.b*x + self.c
  
```

```

>>> my_quad = Quadratic(1,2,3)
>>> my_quad(1)
6
  
```

As a result, we can now invoke the class object directly instead of explicitly invoking a method. This can have many advantages. For instance:

- It allows us to implement objects that can be used in a flexible and intuitive way.
- It allows us to use a class object in contexts where a callable object is expected — using a class object as a decorator, for instance.

What is callable?

In Python, a callable is any object that can be called using parentheses and may return a value. For instance, a function is a callable object (one that can be called/invoked).

A terminal window shows a simple function named `my_method` defined with a single `pass` statement. When the `callable` function is called with `my_method` as an argument, it returns `True`. A callout bubble points to the terminal with the text "A callable object".

```

def my_method():
    pass
  
```

```

>>> callable(my_method)
True
  
```

Coming back to PyTorch

This is what happens when we build deep learning models with PyTorch. For instance, consider the PyTorch class again:



```
PyTorch_model.py

import torch.nn as nn

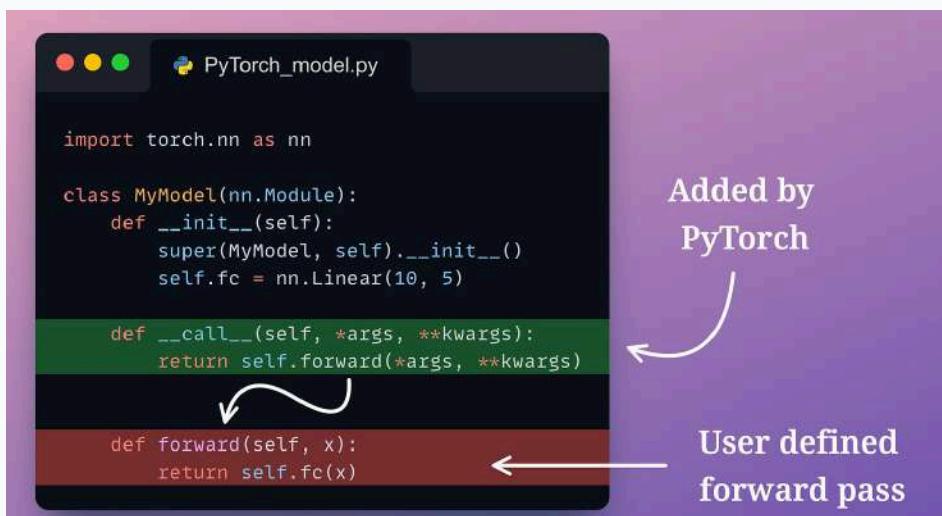
class MyModel(nn.Module):
    def __init__(self):
        super(MyModel, self).__init__()
        self.fc = nn.Linear(10, 5)

    def forward(self, x):
        return self.fc(x)
```

The code defines a PyTorch module named `MyModel`. It includes an `__init__` method that initializes a linear layer `self.fc`. The `forward` method takes an input `x` and returns the result of applying the linear layer to it.

As you may have already guessed, the model object can be invoked because all PyTorch classes implicitly declare the `__call__()` method themselves. Within that `__call__()` method, they invoke the user-defined forward pass.

A simplified version of this is depicted below:



```
PyTorch_model.py

import torch.nn as nn

class MyModel(nn.Module):
    def __init__(self):
        super(MyModel, self).__init__()
        self.fc = nn.Linear(10, 5)

    def __call__(self, *args, **kwargs):
        return self.forward(*args, **kwargs)

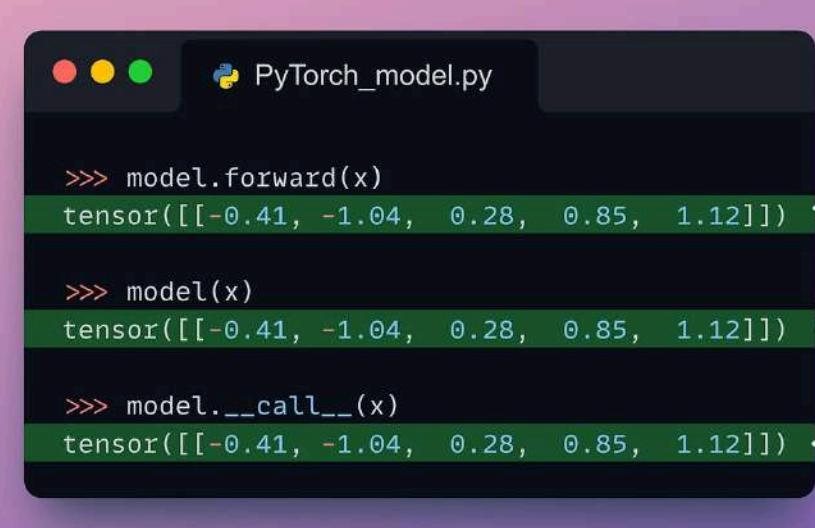
    def forward(self, x):
        return self.fc(x)
```

The diagram illustrates the PyTorch `__call__` method. A green box highlights the `__call__` method, which contains a call to `self.forward(*args, **kwargs)`. A red box highlights the `forward` method, which contains the user-defined logic `return self.fc(x)`. A curved arrow points from the `__call__` method to the `forward` method, labeled "Added by PyTorch". Another curved arrow points from the `forward` method back to the `__call__` method, labeled "User defined forward pass".

- PyTorch itself adds the `__call__()` method.
- The `__call__()` method invokes the user-defined `forward()` method.

This way, Python gets to know that the model object can be invoked like a function — `model()`.

In fact, we can verify that we will get the same output no matter which way we run the forward pass:



The screenshot shows a Jupyter Notebook cell with the following code:

```
>>> model.forward(x)
tensor([[-0.41, -1.04, 0.28, 0.85, 1.12]])

>>> model(x)
tensor([[-0.41, -1.04, 0.28, 0.85, 1.12]])

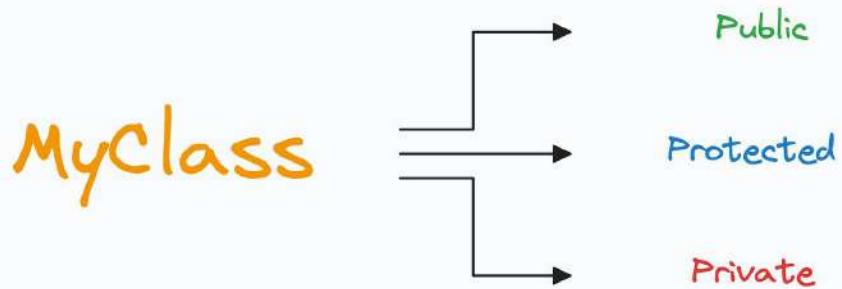
>>> model.__call__(x)
tensor([[-0.41, -1.04, 0.28, 0.85, 1.12]])
```

Three white arrows point from the text "Same output" to the three lines of code above. The text "Same output" is located on the right side of the screen.

All three ways return the same output.

True OOP Encapsulation is Missing From Python

Using access modifiers (public, protected, and private) is fundamental to encapsulation in OOP.



This is not just about Python but applicable to OOP in general. However, when we talk specifically about Python, we notice that it fails to deliver true encapsulation procedures, which one would want to leverage in their OOP code.

How?

Let's understand this in this chapter.

As you may already know, class attributes in OOP can be of three types:

- A public member is accessible everywhere inside and outside the base class, and it is inherited by all child classes.
- A protected member is accessible everywhere inside the base class and child class(es). It is not accessible outside the class.
- A private member is accessible only inside the base class.

But, with Python, there are no such strict enforcements. This is unlike many other programming languages like C++. Thus, protected members behave exactly like public members.

What's more, private members can be (somehow) accessed outside the class as well.

For better clarity, consider the following class implementation:



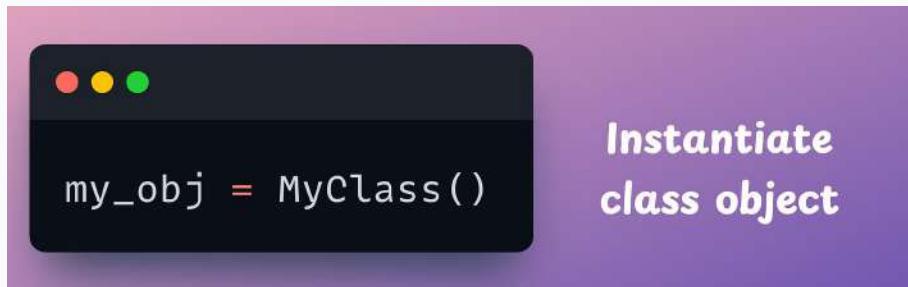
```
class MyClass:

    def __init__(self):
        self.public_attr    = "I'm public"      # 0 underscores
        self._protected_attr = "I'm protected" # 1 underscore
        self.__private_attr = "I'm private"    # 2 underscores
```

Syntactically speaking, in Python:

- A public member is declared with 0 leading underscores.
- A protected member is declared with 1 leading underscore.
- A private member is declared with 2 leading underscores.

Next, we instantiate a class object:



my_obj = MyClass()

**Instantiate
class object**

Moving on, as one would expect, the public attribute is accessible outside the class:



```
class MyClass:
    ...

my_obj = MyClass()

>>> my_obj.public_attr
"I'm public"
```

**Public member is
accessible outside
the class**

However, in a similar way, even the protected attribute is accessible outside the class, which, ideally, should not happen:



The screenshot shows a terminal window with the following code and output:

```

class MyClass:
    ...
myObj = MyClass()
>>> myObj._protected_attr
'I'm protected'

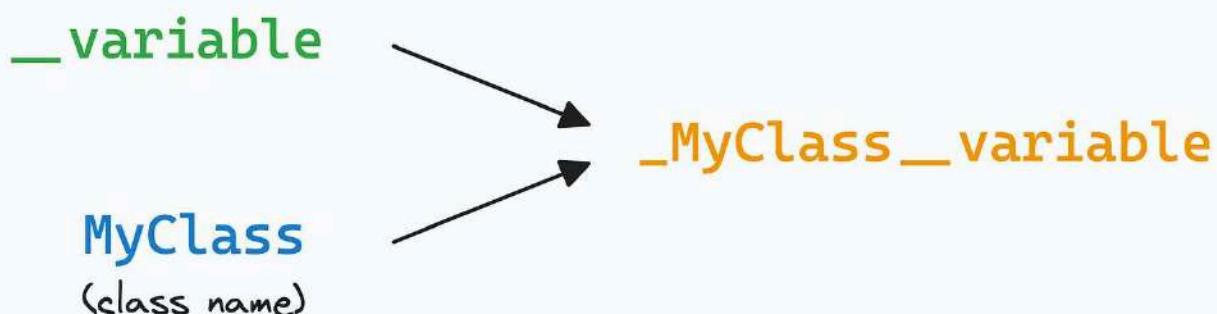
```

To the right of the terminal, the text "Protected member is also accessible outside the class" is displayed in white, with a curved arrow pointing from the text towards the terminal window.

Lastly, when declaring private members, Python performs name mangling.

It is a technique used in programming to avoid naming conflicts between different classes.

Python performs name mangling by attaching underscore-prefixed class name (`_MyClass`) to all members with two leading underscores.



So, while the private member is not directly accessible using its original name, it can still be accessed with the modified name obtained from name mangling.

This is demonstrated below:

The screenshot shows a terminal window with the following code and output:

```

class MyClass:
    ...
my_obj = MyClass()
>>> my_obj.__private_attr
AttributeError: 'MyClass' object has no attribute '__private_attr'
>>> my_obj._MyClass__private_attr
"I'm private"

```

Annotations on the right side of the terminal window explain the behavior:

- Python raises error**: Points to the first line of the error output.
- Private member accessible now**: Points to the second line where the attribute is accessed via its mangled name.

As demonstrated in all the examples above, protected members and private members (with name mangling) can be accessed like public members.

Thus, the point that every Python programmer using OOP must remember here is that Python never enforces encapsulation.

Instead, leveraging encapsulation procedures in Python mainly relies on conventions.

They are used to communicate the accessibility protocols of class members to other programmers.

The screenshot shows a terminal window with the following code and output:

```

class Reviews:
    def __init__(self):
        self.__rating = 5

```



Thus, it is the programmer's responsibility to obey these conventions.

A Common Misconception About `__init__()`

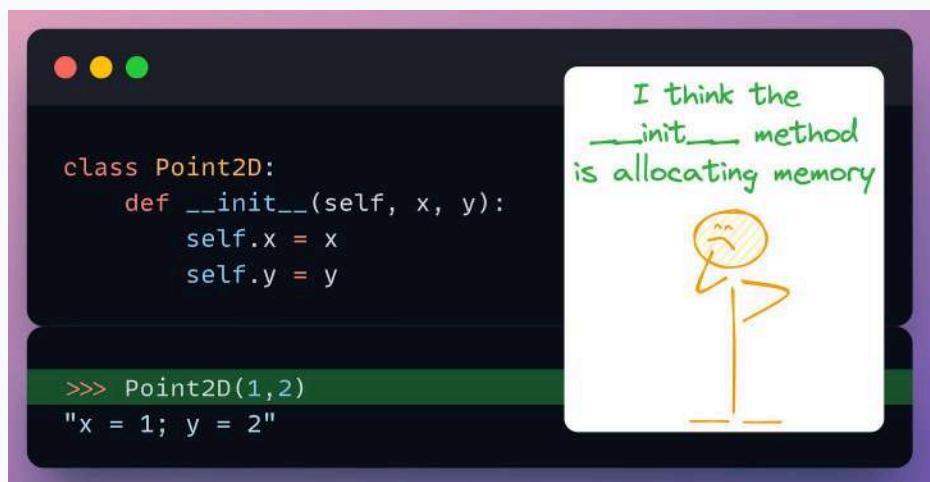
Most Python programmers misinterpret the `__init__()` magic method in Python OOP. They think that it creates a new object, i.e., allocates memory to it. For instance, consider the `Point2D` class below:



Point class with two attributes -> (x, y)

```
class Point2D:
    def __init__(self, x, y):
        self.x = x
        self.y = y
```

Assume that there is also a `__repr__()` method, as it is not shown here. When we create an object (shown below), programmers believe that in the background, it is the `__init__()` method that is allocating memory to their object:



I think the `__init__` method is allocating memory

```
class Point2D:
    def __init__(self, x, y):
        self.x = x
        self.y = y

>>> Point2D(1,2)
"x = 1; y = 2"
```

But that is not true.

When we create ANY object in Python, the `__init__()` method NEVER allocates memory to it.

As the name suggests, `__init__()` only assigns value to an object's attributes, i.e., initialize the attributes.

Instead, it's the `__new__()` magic method that creates a new object and allocates memory to it. To understand better, consider the class implementation below.

```

__new__(cls, x, y)
    if isinstance(x, int) and isinstance(y, int):
        print("Allocating memory")
        return super().__new__(cls)

    else:
        raise ValueError("x and y must be integers")

def __init__(self, x, y):
    self.x = x
    self.y = y

```

Here, we have implemented the `__new__()` method, which checks if the passed arguments are of integer type. Now, if we try to create an object of this class, Python would first validate the checks specified in the `__new__()` method and create a new object only when the specified conditions are true. This is evident from the image below:

Integer values work

```

>>> Point2D(1, 2)
"Allocating memory"
"x = 1; y = 2"

```

Float values raise an error

```

>>> Point2D(1, 2.5)
"ValueError: x and y must be integers"

```

Another popular use case of the `__new__` method is to define singleton classes — classes that can only have one object. For instance, consider the following class implementation:



```
__new__() method
for singleton class

class Point2D:

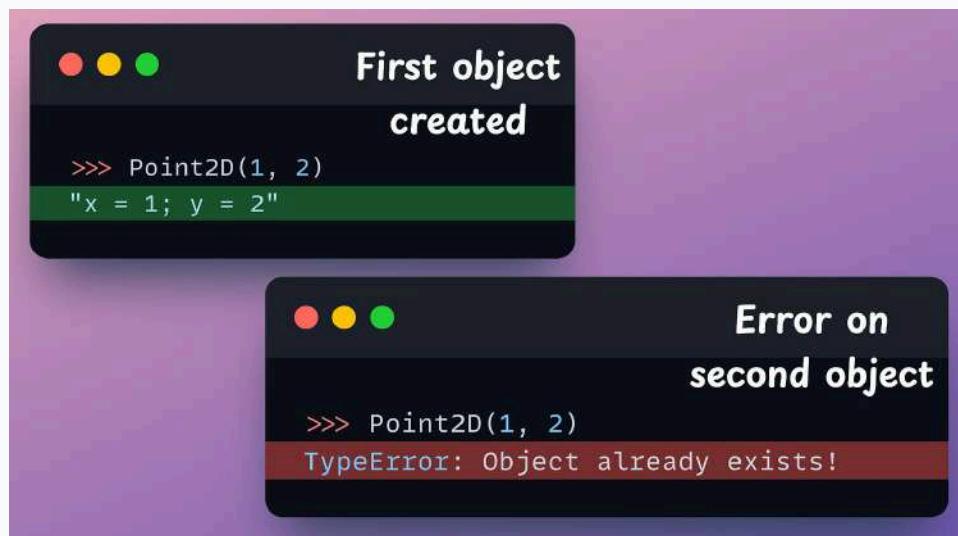
    _class_count = 0

    def __new__(cls, x, y):
        if cls._class_count == 0:
            cls._class_count = 1
            return super().__new__(cls)

        else:
            raise TypeError("Object already exists!")

    def __init__(self, x, y):
        self.x = x
        self.y = y
```

In the above code, the `__new__` method defines a class variable `_class_count`. When a new object is instantiated, if the value of `_class_count=0`, the value of `_class_count` is updated to 1 and the object is returned. After creating the first object, a new object can never be created because the value of `_class_count` will never be 0. This is evident from the image below:



```
First object
created

>>> Point2D(1, 2)
"x = 1; y = 2"

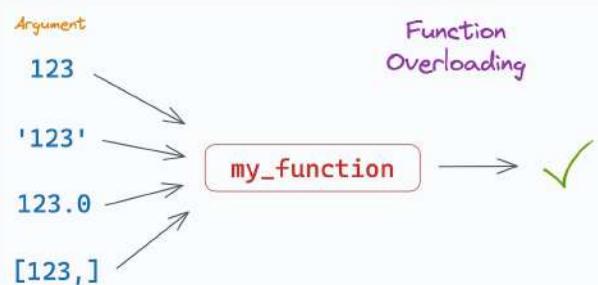
Error on
second object

>>> Point2D(1, 2)
TypeError: Object already exists!
```

Function Overloading in Python

Function overloading is critical to polymorphism, wherein, we may have multiple functions with:

- Same name, and
- Different number (or type) of parameters.



However, Python provides no native support for function overloading. In other words, if we define two (or more) functions with the same name and different parameters, Python will only consider the latest definition corresponding to that function name. This is depicted below.

```
def add(a, b):
    return a + b

def add(a, b, c):
    return a + b + c
```

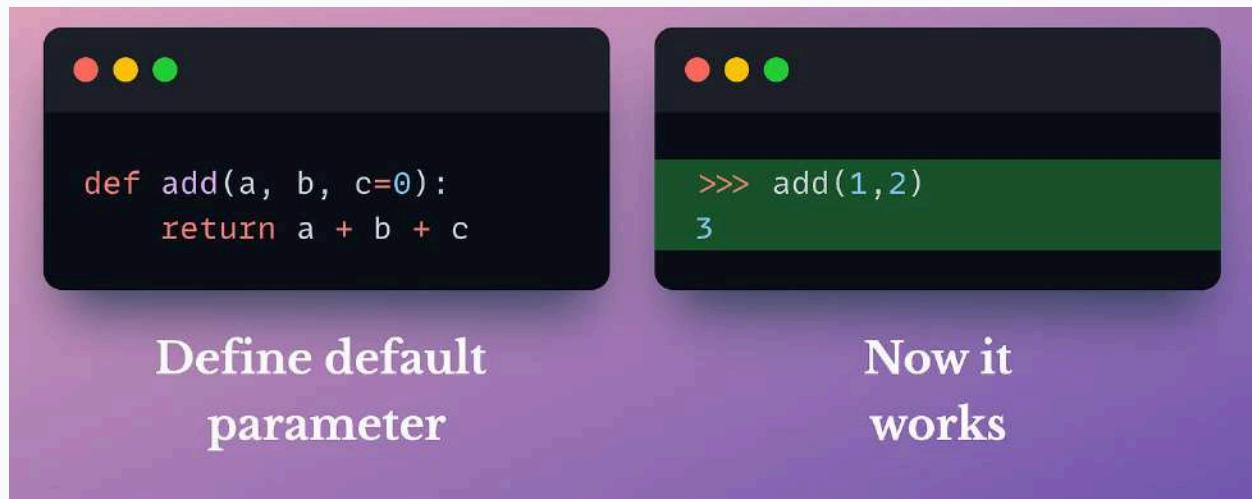
**Functions with
different
parameters**

**Python expects
three arguments**

```
>>> add(1,2)
TypeError: add() missing 1 required
positional argument: 'c'

>>> add(1,2,3)
6
```

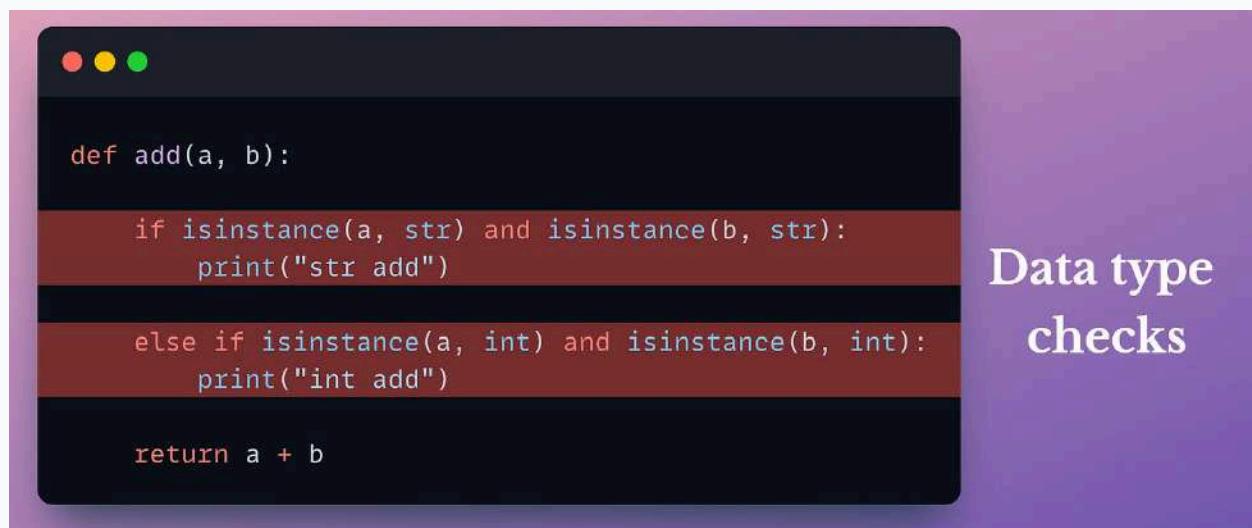
As depicted above, the latest definition of add() had three parameters. That is why passing two arguments raised an error. This restricts us from writing good polymorphic code in Python. Of course, there are ways to prevent this error, as depicted below:



The image shows a Jupyter Notebook interface. On the left, a code cell contains the following Python code:`def add(a, b, c=0):
 return a + b + c`On the right, a command cell shows the output of running the function with two arguments:`>>> add(1,2)
3`

Below the code cell, the text "Define default parameter" is displayed. Below the command cell, the text "Now it works" is displayed.

Also, if the same parameter can take multiple data types, we can write (somewhat) polymorphic code by adding if statements using `isinstance()`.



The image shows a Jupyter Notebook interface. A code cell contains the following Python code:`def add(a, b):
 if isinstance(a, str) and isinstance(b, str):
 print("str add")
 else if isinstance(a, int) and isinstance(b, int):
 print("int add")

 return a + b`

Below the code cell, the text "Data type checks" is displayed.

But this isn't as elegant as defining multiple functions with different data types like we can do in, say, C++, is it?

```

#include <iostream>
using namespace std;

void add(int a, int b){
    cout << "sum = " << (a + b);
}

void add(double a, double b){
    cout << endl << "sum = " << (a + b);
}

int main()
{
    add(10, 2); // int arguments

    add(5.3, 6.2); // float arguments

    return 0;
}

```

Two int parameters

Two float parameters

C++ automatically invokes the correct method corresponding to the argument data type. While exploring this for one of my projects, I found a pretty handy way to enable function overloading in Python

The `@dispatch` decorator from the Multidispatch library allows us to leverage the standard and elegant function overloading in Python, as demonstrated below:

```

from multipledispatch import dispatch

@dispatch(int, int)
def add(a, b):
    print("int add")
    return a + b

@dispatch(str, str)
def add(a, b):
    print("str add")
    return a + b

```

```

>>> add(1, 2) # int arguments
"int add"
3

>>> add('1', '2') # str arguments
"str add"
'12'

```

As depicted above, we have multiple functions with the same name and different parameters. The `@dispatch` decorator allows us to invoke the correct function corresponding to the parameters passed during the function call.

Another cool thing about the `@dispatch` decorator is that it raises an error when the function call does not match any of the function definitions:

The screenshot shows a Jupyter Notebook interface. On the left, there is a code cell containing Python code. The code imports `multipledispatch` and defines two functions, `add`, each decorated with `@dispatch` for different parameter types: `(int, int)` and `(str, str)`. Both functions print a string indicating the type of addition and return the sum. On the right, the output of running `>>> add(1, '2')` is shown, resulting in an `NotImplementedError` because the function signature does not match either definition.

```
from multipledispatch import dispatch

@dispatch(int, int)
def add(a, b):
    print("int add")
    return a + b

@dispatch(str, str)
def add(a, b):
    print("str add")
    return a + b
```

```
>>> add(1, '2')
NotImplementedError: Could not
find signature for add: <int, str>
```

**<int, str> definition
not implemented**

This allows us to identify errors which can often go unnoticed.

In a dummy experimentation, I did notice a slight increment in the run time. But the order of increase was in nanoseconds or so, which can be safely ignored.