

Session - 9

ML - SYSTEM DESIGN - 1

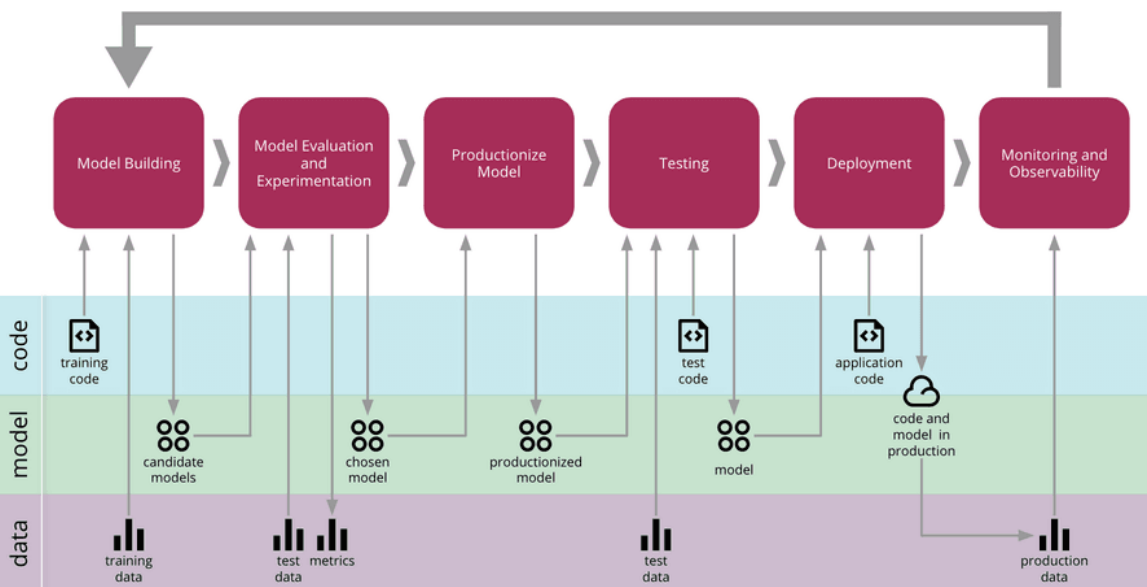
June 23,
2025

WHEN YOU'VE PREPARED JUST DSA FOR INTERVIEW



Agenda

- * What does ML system design includes.



60 → EDA / Data cleaning
 10% → model dev
 30 → maintenance

Things to observe post deployment:




1. Online Accuracy Metrics
2. Latency shouldn't increase beyond the threshold.
3. Server Load (typically this should be less than 50%)
4. Memory Consumption (should be below threshold)
5. Throughput - (Requests you can serve per second, this shouldn't decrease, beyond a certain threshold)

Which of these is NOT a system health/performance metric typically monitored after deploying an ML model?

1 user has participated

A	Requests per second	0%
<input checked="" type="checkbox"/>	B Accuracy of the model's predictions	100%
C	CPU/memory/data utilization	0%
D	Latency per request	0%

[End Quiz Now](#)

<div> <div>  <div> <div></div> <div>Rajeev Ranjan K...</div> <div>1/1 87.86</div> </div> </div> <div>  <div> <div></div> <div>Phani Kishore R...</div> <div>1/1 88.87</div> </div> </div> <div>  <div> <div></div> <div>Gajanan Todeti</div> <div>1/1 89.93</div> </div> </div> </div>		
4	Ashwini Uday	1/1 87.66
5	Santhosh K	1/1 87.23
6	Rajaraman	1/1 79.20
7	Sonali Ghosh	1/1 76.90
8	Lokesh Sonawane	1/1 71.16

What is the purpose of the 'Model Evaluation and Experimentation' phase in the ML model lifecycle?

0 users have participated

A	Deploying the model to a production environment	0%
B	Preparing data and initial code formulation	0%
✓ C	Feature selection, hyperparameter tuning, and algorithm comparison	0%
D	Monitoring and observing the model in a live environment	0%

End Quiz Now

Based on all quizzes from the session

Rajeev Ranjan K...	Phani Kishore R...	Santhosh K...
2/2	2/2	2/2
187.10	190.00	182.43

1/2	Ashwini Uday	2/2	177.76
1/2	Gajanan Todeti	2/2	177.73
1/2	Sonali Ghosh	2/2	163.17
1/2	Avinadam Maji	1/2	85.00
1/2	Chirag Bhatt	1/2	82.73
1/2	Rajaraman	1/2	79.20
1/2	Lokesh Sonawane	1/2	71.16

Drift in Data/Model

→ Something has changed → ① Target
② Features.

① Concept drift: Properties of target has changed.

①. Fraud detection

Soln: Train the model with more recent data

② Data drift: Properties of input Feature has changed.

Soln: Engineer new Features.

Sudden Drift:

A new concept occurs within a short time.



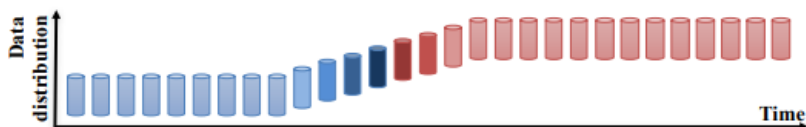
Gradual Drift:

A new concept gradually replaces an old one over a period of time.



Incremental Drift:

An old concept incrementally changes to a new concept over a period of time.



Reoccurring Concepts:

An old concept may reoccur after some time.



How to Notice Drift

① Visual Inspection: Input data over time

① distribution changes

→ mean (over months)
→ median (over months)

② Plot data & see if there's big change

② Statistical Test For same column over diff periods

2

Model performance monitor - Key of PLive metrics
Over past 6 months

What is drift in the context of machine learning?

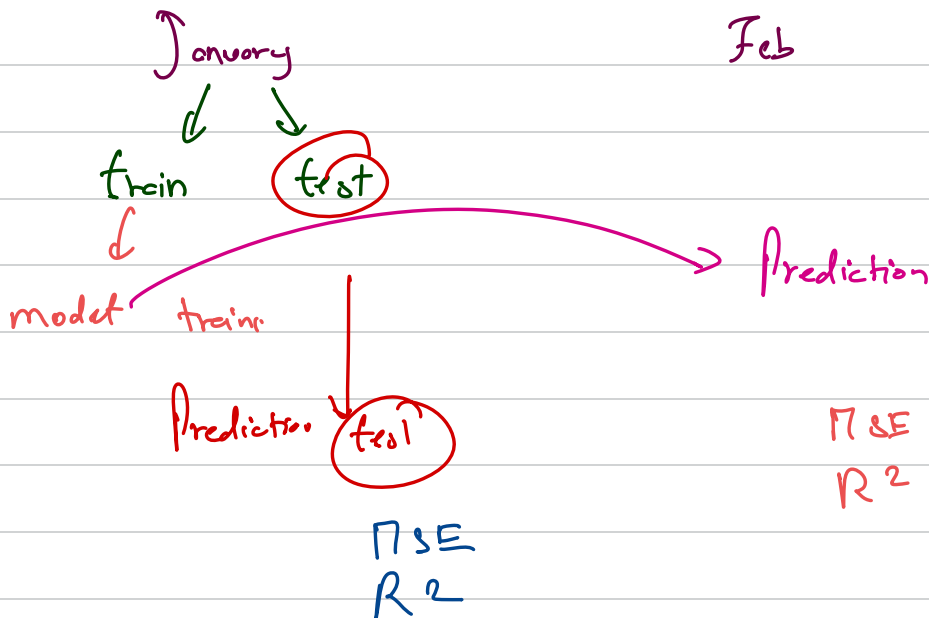
1 user has participated

- A The movement of an autonomous vehicle 0%
- B The physical displacement of a computing server 0%
- C The change in the statistical properties of model data over time 100%
- D The time it takes for a model to make a prediction 0%

End Quiz Now

Based on all quizzes from the session

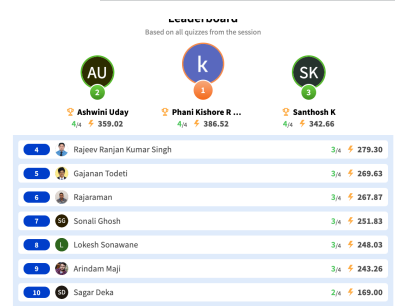
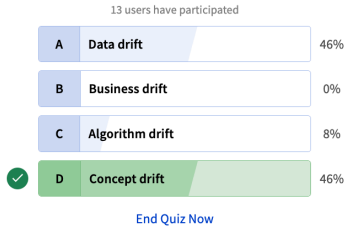
2	Rajeev Ranjan K...	3/3	279.30
1	Phani Kishore R ...	3/3	288.43
3	Gajanan Todeti	3/3	269.63
4	Ashwini Uday	3/3	269.13
5	Santhosh K	3/3	247.47
6	Arindam Maji	2/3	178.37
7	Rajaraman	2/3	173.87
8	Sonali Ghosh	2/3	163.17
9	Lokesh Sonawane	2/3	162.70
10	Sayyid Thajudheen Thangal K C	2/3	130.99



After observing change in model metrics:

1. First check if my features have changed, their distribution.
2. If they've changed, you know the perp, otherwise it's concept drift.

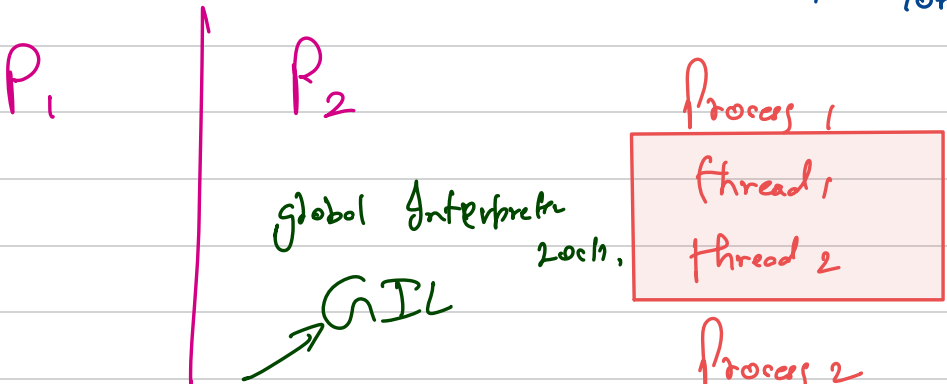
What type of drift occurs when there's a change in the relationship between input data 'x' and output data 'y'?



MULTI PROCESSING

Python applications → Process → Reserves RAM
" compute

multiple workers (thread)



Inside a process only one thread is active at any instance

In C/C++/Java \rightarrow a process

multiple threads,
can execute
concurrently

multi-processing

multi-threads

multi-processes

Website - 1K image

