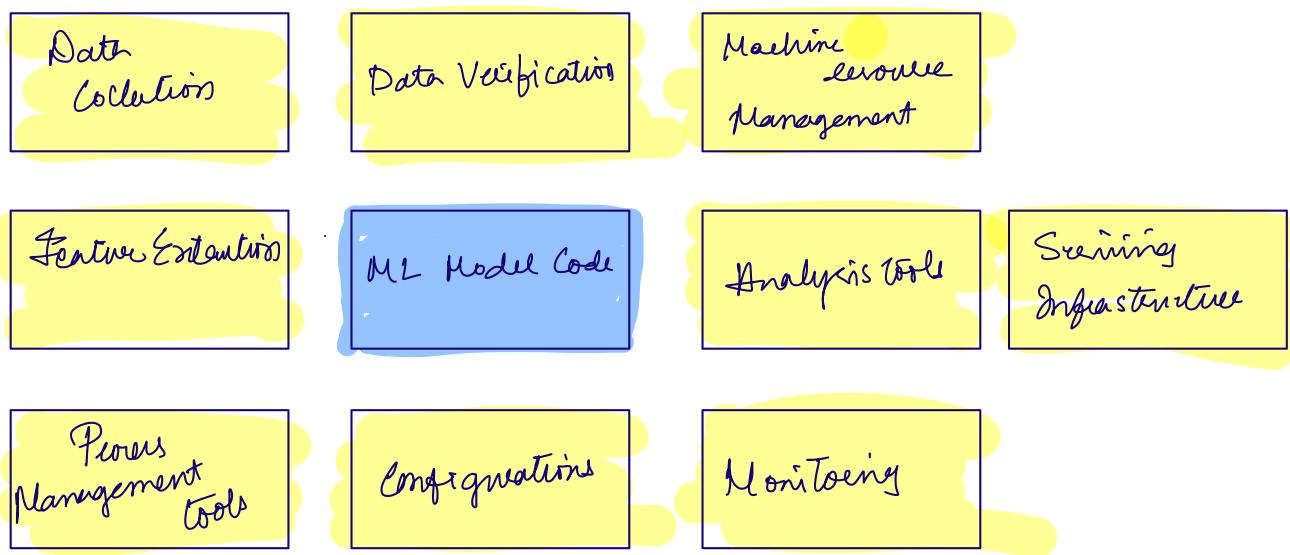


Real World ML

Model is small part of whole large ecosystem.



ML Model is heart of whole Machine learning production system.

So overall Model constitutes 5% of whole ecosystem. Major of efforts goes in collecting the data, validating it, and extracting features from it.

Static vs Dynamic Training

① Static Training → which is also called "offline training"
means you train model only once, then serve the model for a while. [pros: simple, disadvantage: static, it depends and are not up-to-date]

② Dynamic Training → "Online training" you train the model continuously or atleast frequently, serve the recently trained model. [pros: mostly up-to-date, cons: more work and resources]

→ Data don't change then choose static training
cheaper to create & maintain
data which is tend to constant can be changed

hence always better practice to monitor input values & features.

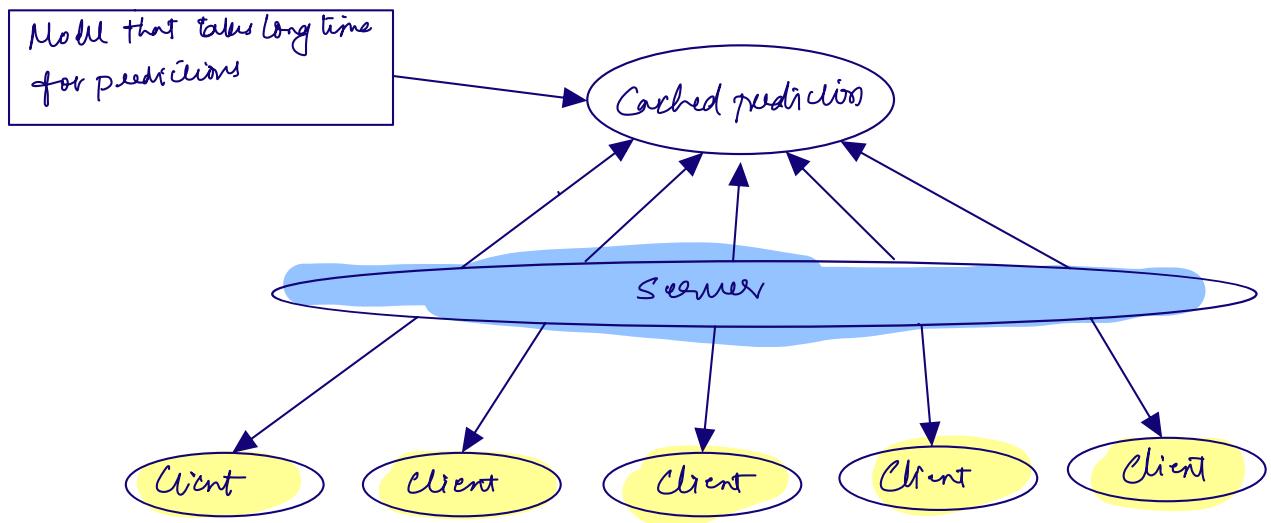
Static vs Dynamic Inference

Inference is the process of making predictions by applying trained model on unlabeled examples.

Model can infer predictions in two ways.

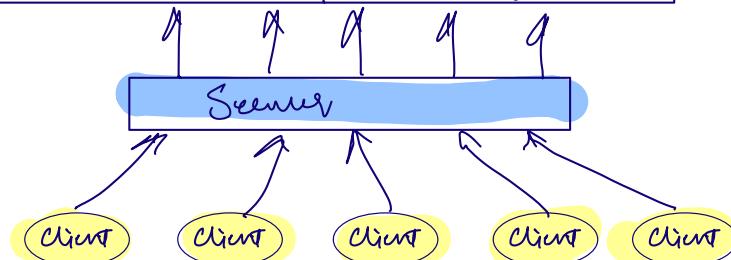
Static Inference is also called as offline inference or batch inference
↳ Means model does the inference on bunch of common unlabeled examples & then caches the predictions somewhere

Dynamic Inference also called as online inference or real-time inference → means Model makes predictions on demand
eg: When client demands the predictions



In static inference, a model generates predictions, which are then cached in a server.

model which does predictions very quickly



Transforming Data

Raw data must be feature engineered (transformed)

Transforming data can be done at any of the below phases →

When data is feature engineered then this process is called transformed

- ① Before training the model
- ② While training the model

Before Training the model

⇒ Then you write the code or use specialized tools to transform the raw data

⇒ Store the transformed data on the disk or somewhere from which model can pick/ingest

⇒ Systems transform the data only once.

⇒ System can analyze the entire data set to determine the next transformation strategy

Disadvantage → Training-Serving Skew.

It is more dangerous when system performs dynamic (online) inference

System use dynamic inference,

⇒ Software transforms the raw data set

usually differs from software that serves predictions
It causes training - serving skew

In context system must be inference can
sometimes run the same software.

Transforming data while training

Transformation is part of model code. The model ingests
raw data and transforms it.

→ you can use the same raw data file if you
change the transformations

* you are ensured the same transformation at training
and at predictions time

→ Complicated transforms can increase model latency
* transformations occur for each and every batch

Transforming the data per batch can be tricky, for example
Suppose you want to use Z-score normalization do transformations
on numerical data

Z-score normalization requires the mean and standard
deviations of the feature.

transformations per batch mean you only have
access to one batch only. not full dataset.

So if batches are highly variant, a Z-score of say -2.5 in one batch won't have the same meaning as -2.5 in another batch

Workaround your systems can precompute the mean & standard deviation across the entire dataset & then use them as constants in the model.

Deployment Testing

Model which is ready to be deployed

- ① Validating input data
- ② Validating feature data
- ③ Validating the quality of new model versions
- ④ Validating & scaling infrastructure
- ⑤ Testing integration between pipeline components.

TDD → Test Driven Deployment



It is different in Machine learning

Eg: Before training the model you can write test to validate the loss, instead you must first define the achievable loss during the model development. Then test new model versions against the achievable loss.

Test model updates, with reproducible training

When wanted to improve your model, for example

do some feature engineering and train the model again hoping to get better results (or same results). It is sometimes difficult to reproduce model training.

To improve reproducibility below steps to be followed:

- ⇒ Deterministically seed the random number generator
- ⇒ Initialize model components in a fixed order to ensure the components get the same random from the random number generator on every run. ML libraries typically handle this requirement automatically.
- ⇒ Take the average of several runs on the model
- ⇒ Use version control, even for preliminary iterations so that you can pinpoint code and parameters when investigating your model.

Test calls to machine learning API

Write unit test to generate random I/P data & run single step gradient descent.

Write integration tests for pipeline component.

Write integration test, where components work together.

Basically when you change in one component cause error in another component.

To run integration test faster, train on a subset of the data or with simpler model.

Validate model quality before Scoring

Test below two test degradation

① Sudden degradation → if any in the new revisions
could cause significantly

② Slow degradation

lower quality

validate new revisions

↓
your test for sudden degradation might
not detect slow degradation in model quality

by checking this quality
against the previous run

Same above Model predictions on a validation

dataset meet a fixed threshold. If your validation dataset deviates
from quality threshold , update validation dataset.

Validate model - Infrastructure compatibility before Scoring

if your model is updated faster than your server
the check software dependencies , mainly causing incompatibilities

Monitoring Pipelines

① Write "data schema" to validate raw data.

↳ to monitor your data, you should continuously
check it against expected statistical values by writing rules that
the data must satisfy. This collection of rules is called a
data schema. Below steps-

① Understand the range and distributions of your feature.

② Encode your understanding into the data schema

→ Ensure that user-submitted ratings are
always in range 1 to 5

→ Check that word which occurs more frequently.

= check each categorical feature is set to a value from fixed set of possible values.

③ Test your data against the data schema. Your schema should catch data errors such as

→ Anomalies

→ Unexpected values of categorical variables

→ Unexpected data distributions

⇒ Write unit tests to validate feature engineering

Data schema is nothing but certain rules to pass the data - in this scenario it might pass the data schema model cannot used to train the model.

Data which is feature engineered is the one model is trained on.

Eg: Model will be trained on normalized feature values, feature engineered data can be very different from raw input data

With unit tests based on understand. Eg the conditions such as

① Numeric features are scaled, for example $0 \rightarrow 1$

② One-hot-encoded vectors only, contains single 1 and N-1 zeros.

③ data distributions after transformation contains a single 1 & N-1 zeros

④ Outliers are handled, such as by scaling or clipping

Check Metrics for important data slices

A successful cohort sometimes observes an unsuccessful subset. In other words, a model with great overall metrics might still make terrible predictions for certain situations. For example

AUC → Area Under the Curve

If you're satisfied with great AUC then you might not notice the model issues in the sahara desert. If making good predictions for every region is important then track performance for every regions.

Subset of data like sahara desert is

called data slices

Identify data slices of interest. Then compare model metrics for these data slices against metrics of your entire dataset.

Use real World Metrics

Model Metrics don't necessarily measure the real-world impact of your model. For example changing a hyperparameter might increase a model AUC, but how did change affect user experience. To measure real world impact, you need to define separate metrics. For example you could survey

Training Skew

input data during training differs from your input data in seeing.

① Schema Skew

Definition: Training & seeing input data do not conform to the same schema.

Ex: The format or distribution of the seeing data changes while your model continues to train on old data.

Solution: Use same schema to validate training and seeing data

→ Separability check the statistic rules not checked by schema such as the failure of mining val

② Feature Skew: Engineered data differs from training & seeing

Ex: Feature engineering code differs between training and seeing, producing different engineered data.

Solution: Similar to schema skew, apply the same statistical rules across training and seeing engineered data.

⇒ Track the number of detected skewed features and ratio of skewed examples per feature.

Causes: Training & Seeing Skew can be subtle. Always consider what data is available to your model at prediction time.

Underfitted → when you have model to predict daily revenue based on no of customers and also

→ Problem → you don't have customers at prediction time before the daily sales are complete → same This feature is not useful

Check for Label Leakers

Label Leakers meant ground truth labels that you are trying to predict have inadvertently entered your training features.

Some times very difficult to detect.

Monitor Model Age

If the seeing data evolves with time and your model is retrained sequentially. Then you will see decline in quality of model.

Rebuild the linear since the model was retrained on new data and set threshold age for delete. You should monitor the age throughout the pipe line.

⇒ Test the model weights & outputs are numerically stable.
write test to check NaN, inf values.

Monitor Model Performance

⇒ Track Model performance by revisions of code

involve 0 data.

⇒ Catch Memory leaks.

⇒ Monitor API response time

⇒ Number of queries answered per second.

Test the quality of live model on Seed data

- ① generate values using human rates.
- ② Investigate models that show significant statistical bias in predictions.
- ③ Track real world metrics for model.
if Confusing Spam , compare with much expected spam.
- ④ Mitigate potential divergence between training & testing data setting a new model.

Randomization

data generation pipeline reproducible

Eg: Suppose you want to add feature to check model quality. your dataset should be identical except for this new feature. In that spirit make sure any randomization in data generation can be made deterministic.

- ① Seed your Random Number Generators (RNG's)

Seeding console can be set of values in same order whenever you run it.

② Use invariant hash keys

Hashing is common way to split data

Eg:

hash each example & use the resulting integer to decide in which split to place the example
Input to your hash function shouldn't change every time

Considerations for Hashing

imagine you are using search queries & using hashing to include or exclude queries.

⇒ your training set will see a less diverse set of queries :-

⇒ Evaluation time will be artificially hard, because they won't overlap much with your training data, so your evaluation should reflect that

Considerations before your data & model
in production Systems

① Is each feature helpful

⇒ Continue you should monitor your features if they are contributing little or nothing to the model's prediction

ability

⇒ If input data for that feature abruptly changes
model's behaviour will also change undetectably.

② Does the usefulness of the feature justify the cost of including it?

⇒ It is always tempting to add more features to the model,

Eg: You find a feature which result you better predictions, which is neither than slightly worst prediction but that adds maintenance burden.

③ Is your data source reliable?

⇒ Is the signal always going to be available or is it coming from an unreliable source?

⇒ Is the signal coming from a server that works under heavy load.

⇒ If signal coming from human who go to vacation every August.

⇒ Does the system that computes gone models input data ever change? Then how often

⇒ How to know when system changes!

→ Never consider creating your own copy of the data you receive from the upstream process. Then only advance to the next version of the upstream data when you are certain that it is safe.

* Is your Model part of a feedback loop?

Model can affect its own training data.

For example

Result of some model can be become directly or indirectly feature to that same model

→ Model can affect another model for example, consider two models for predicting stock prices.

Eg:

→ Model A, (bad predictive model)

→ Model B

If Model A is buggy, it mistakenly decides to buy stock in stock X. Then purchases drive up the price of stock X. Model B uses the price of stock X as an input feature, thus model B can come to false conclusion.