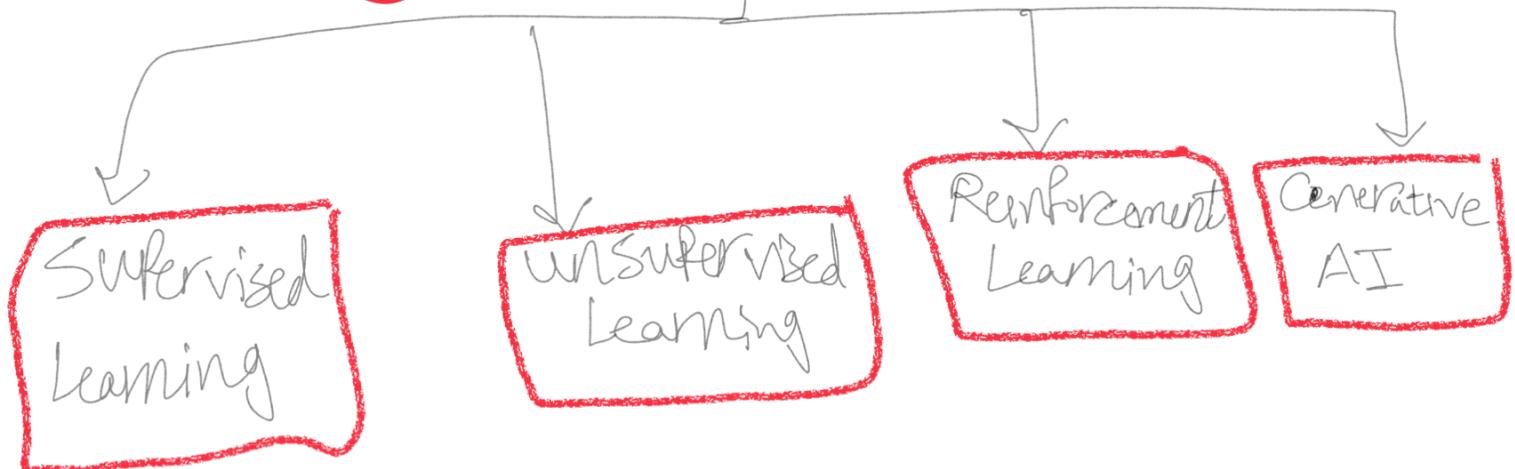


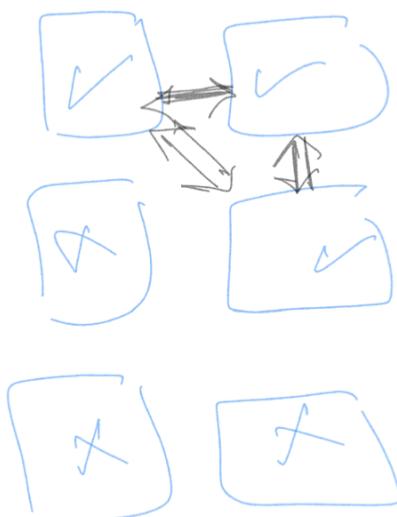
Machine Learning

Course Link: <https://developers.google.com/machine-learning/crash-course/>

Types of Machine Learning

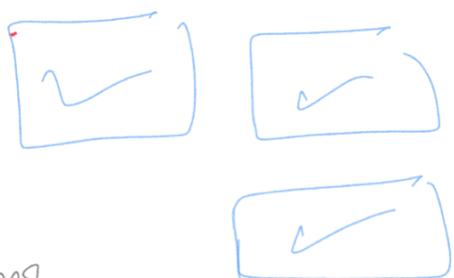


Supervised Learning:



Connecting

"Finding connections
in the data that
produce correct
answers"



Predictions

Data

"has correct
Answers"

Use cases of Supervised Learning

Regression

Classification

Data

number

* Output is number

Binary

Multi class

Data

Data

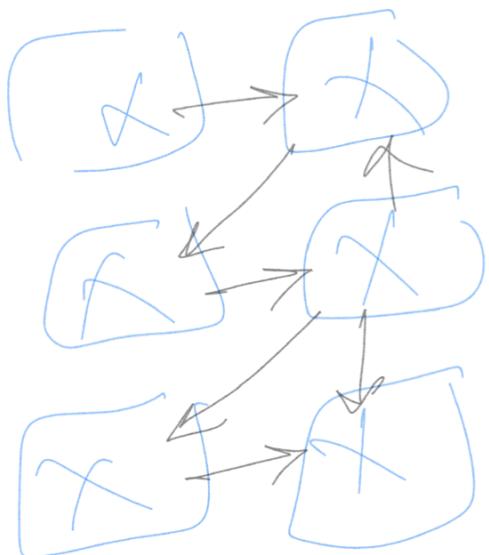
Class
1

Class
2

* Only 2 classes
in output

* more than
2 classes
in output

Unsupervised Learning:



clustering
 "Finding groupings
 in the data"



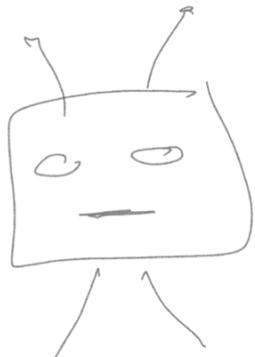
Predictions

Random Data

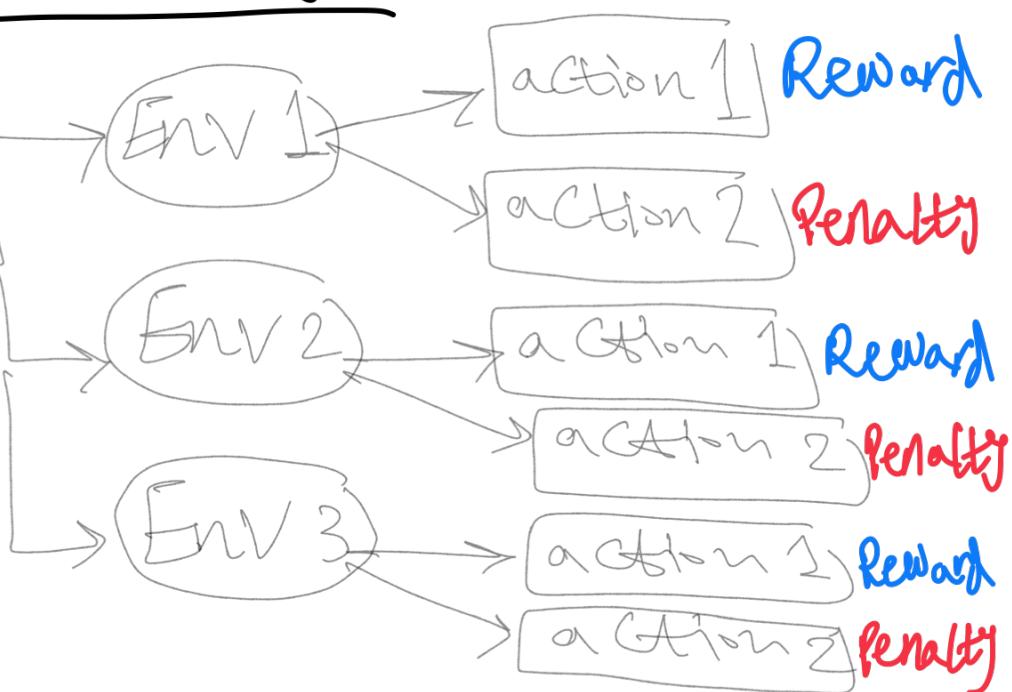
"has no correct
 answers"

"Defined by
 ML model"

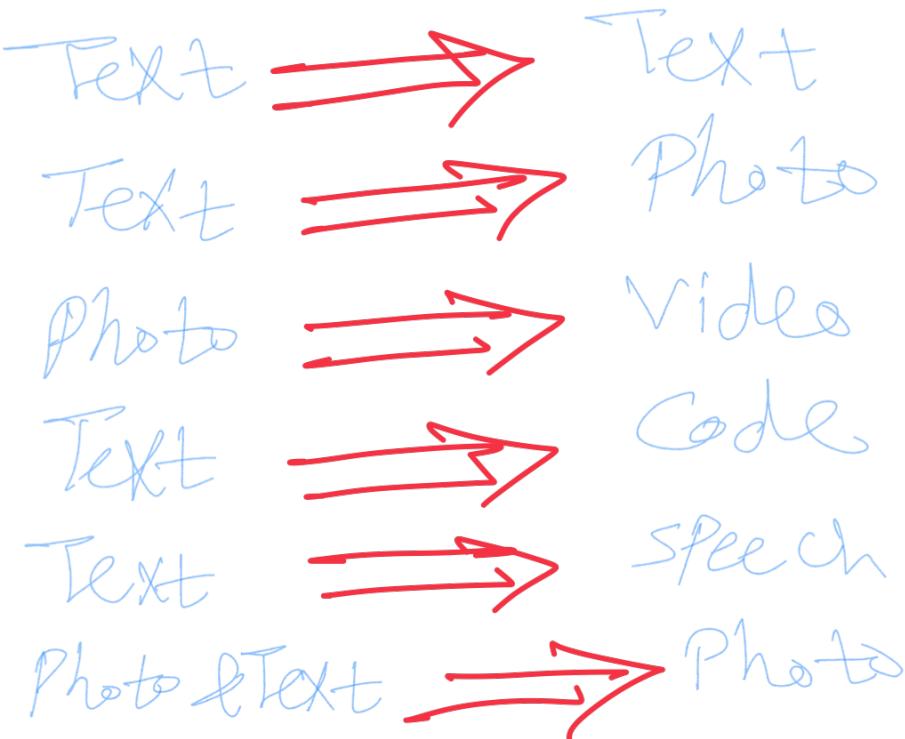
Reinforcement Learning:



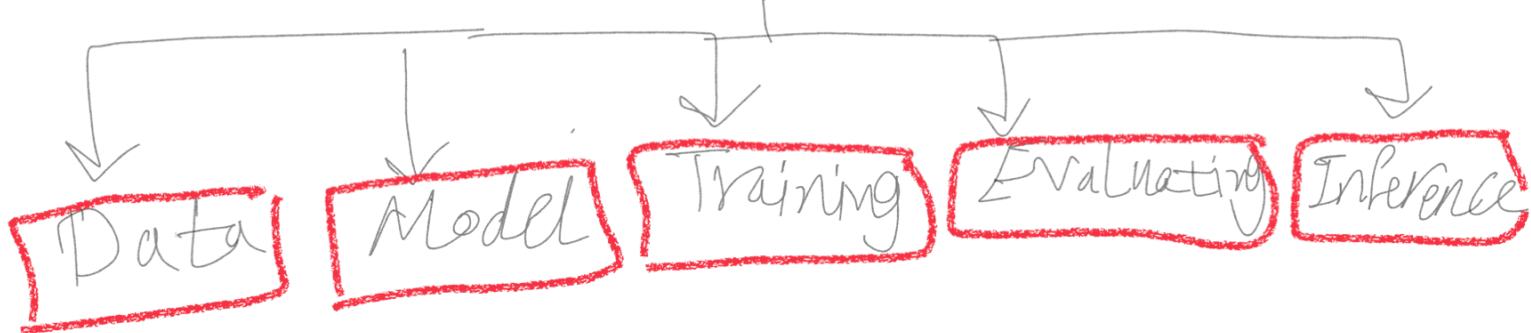
Model



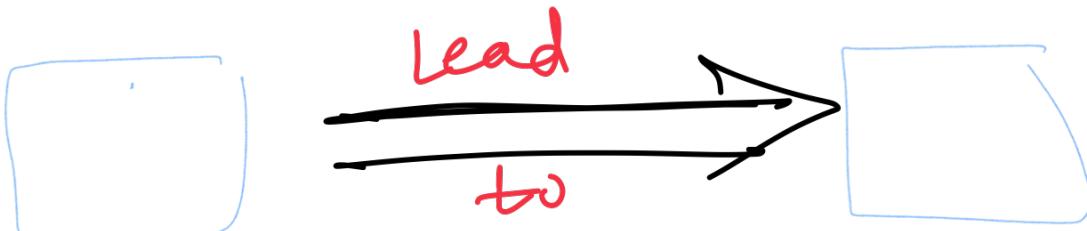
Generative AI :



Supervised Learning Concepts



Data :



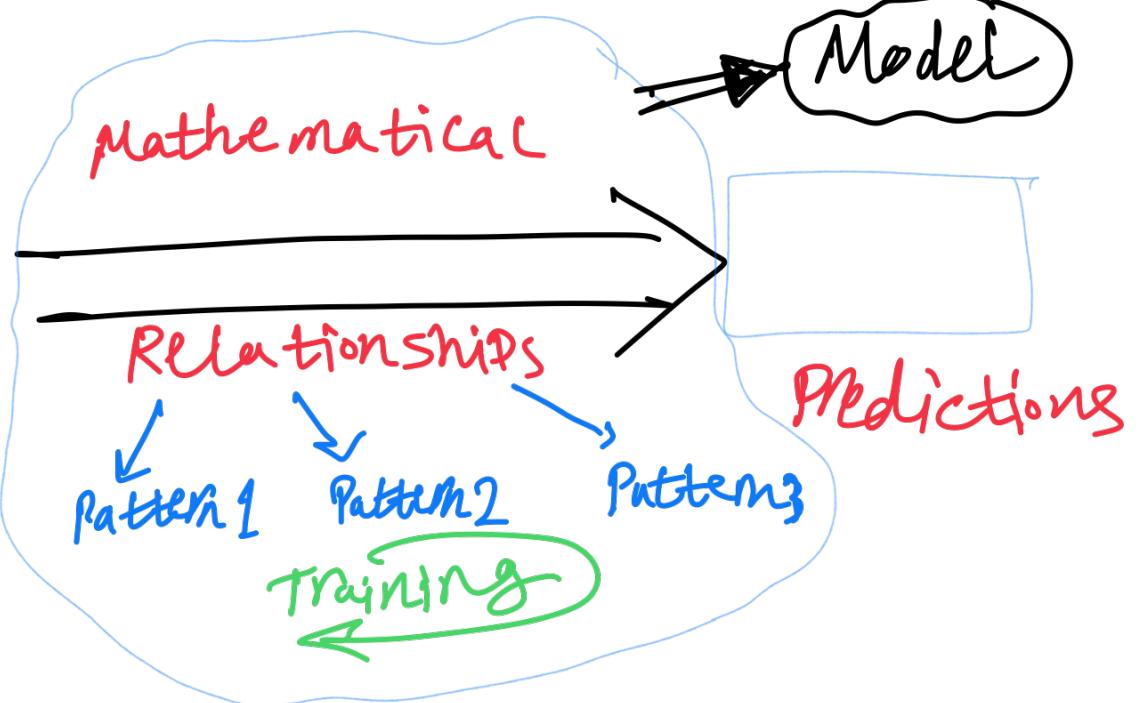
Large size Data
& High Diversity

Precision
Predictions

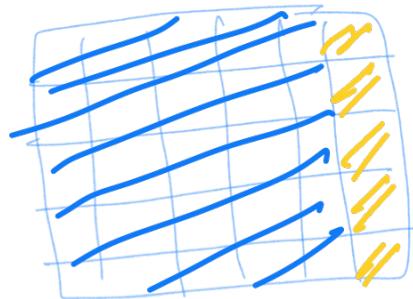
Model :



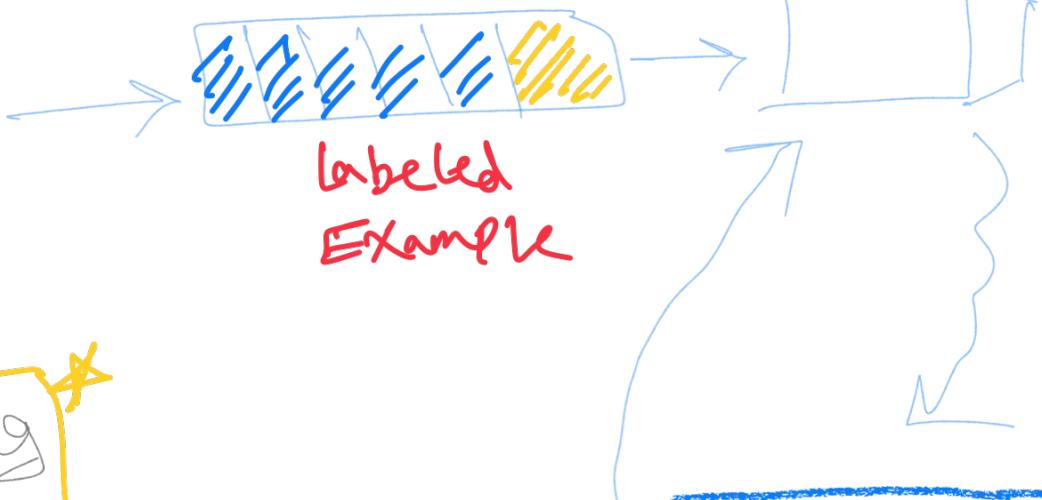
Data



Training :



Dataset

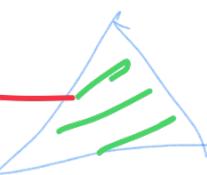


During Training
We can make
subtle adjustments
to the configs, &
features the Model
uses to get better
predictions

Prediction

update
Model

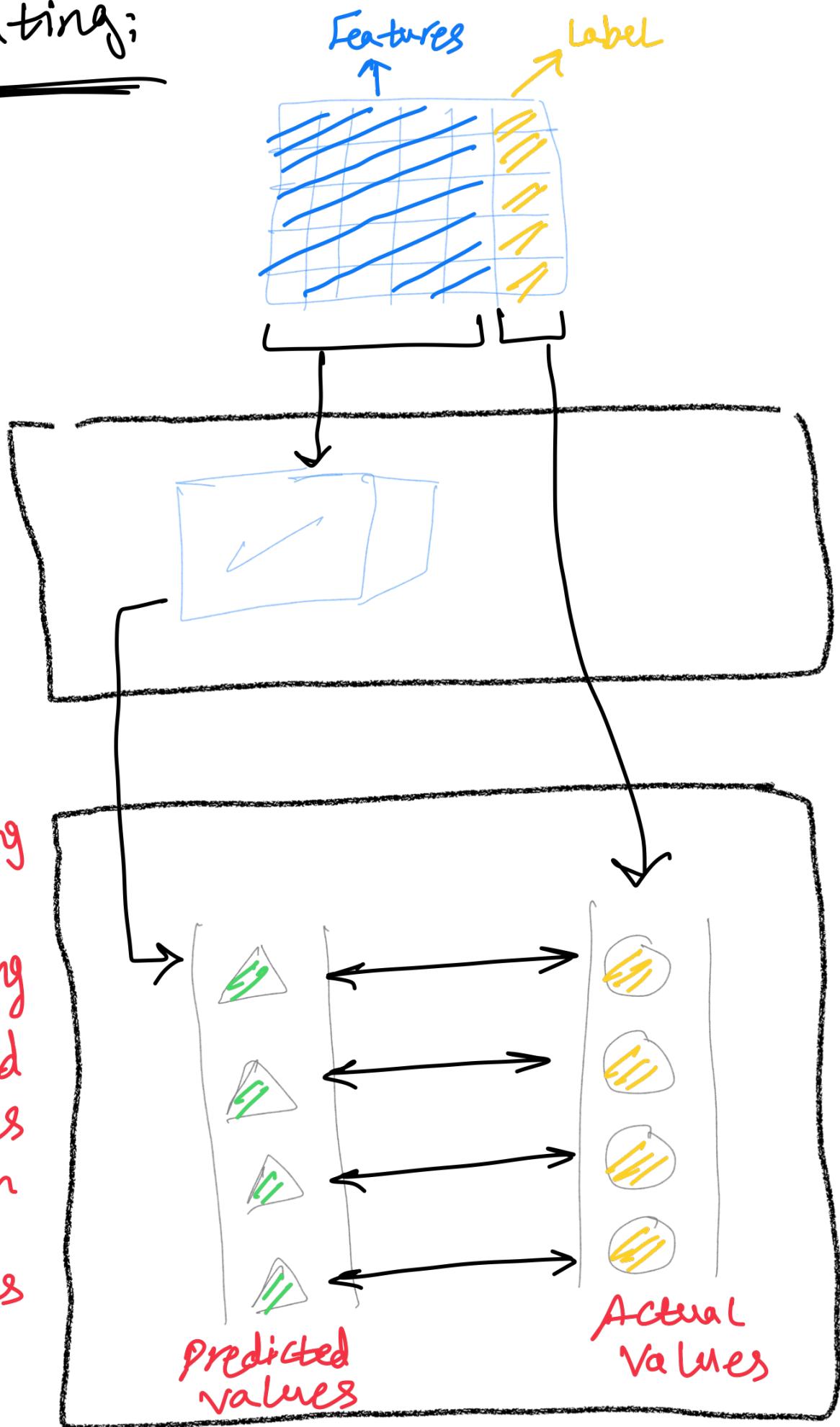
Actual



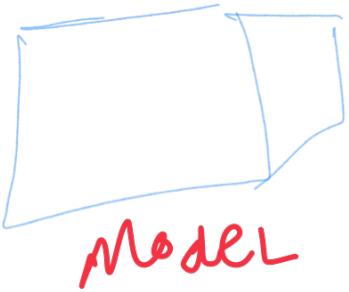
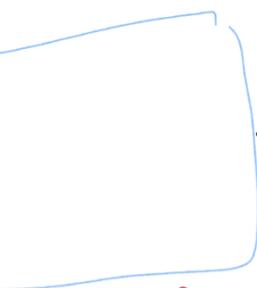
VS.



Evaluating:

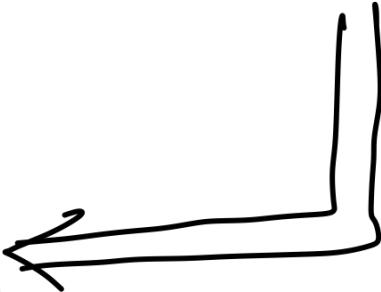


Inference:



unlabeled
Example

Satisfying
Precision
Predictions

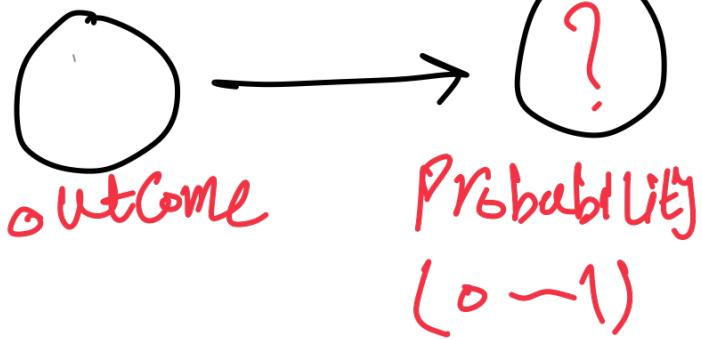
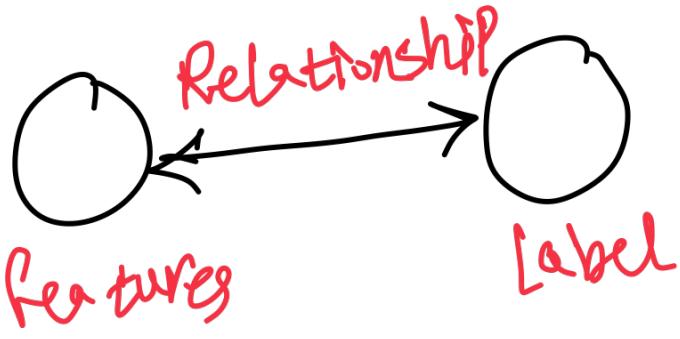


Predictions
OR Inferences

Regression Types

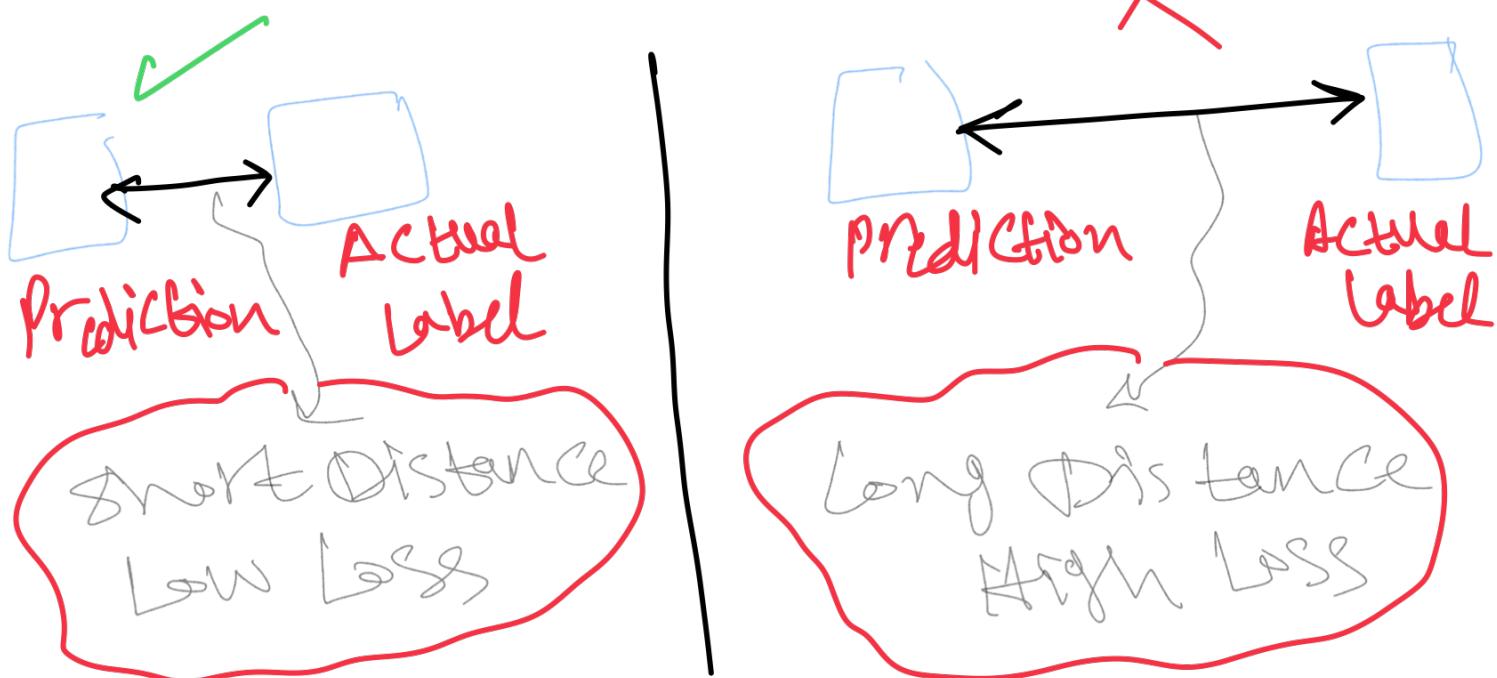
Linear

Logistic



Linear Regression:

(Loss): Numeric metric that describes how wrong a model's predictions are.



Types of Loss

Mean Absolute Error **MAE**

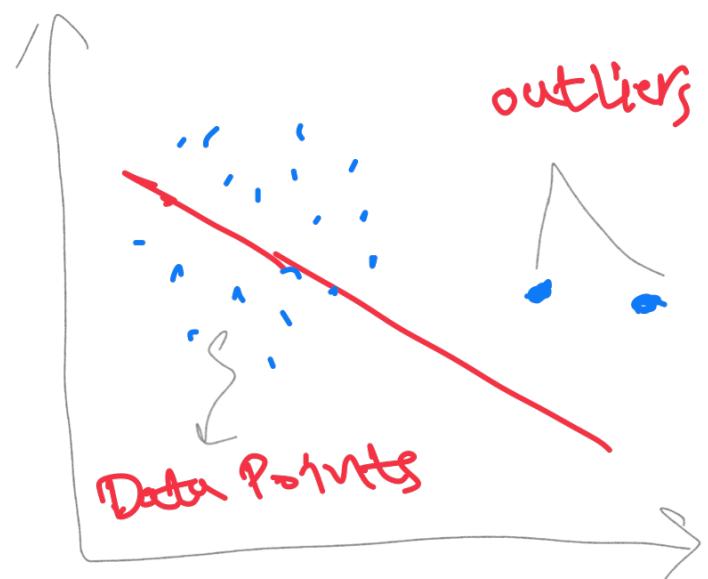
Mean Squared Error **MSE**

MAE: The Model is further away from the outliers but closer to most of other data points.

MSE: The Model is closer to the outliers but further away from most of other data points.



Model trained
using **MAE**



Model trained
using **MSE**

Equation:

$$y' = b + w_1 x_1$$

Prediction

Bias

Weight Feature
value

Calculated from
Training

Gradient Descent:

Weight

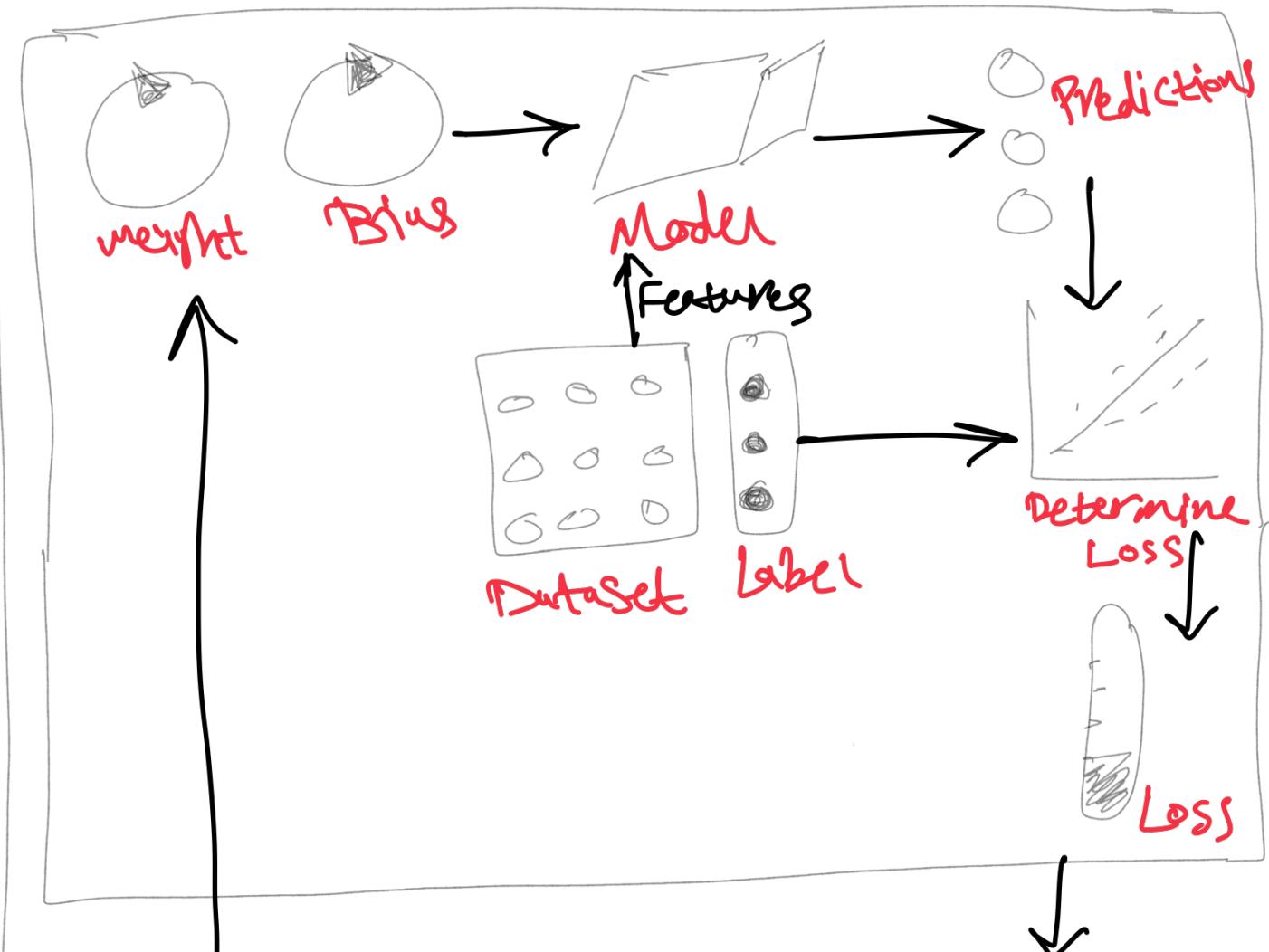
Technique

Bias

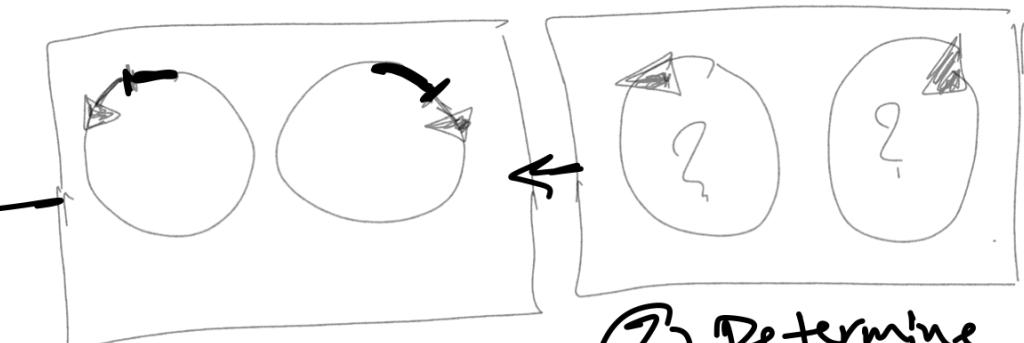
Produce

Lowest
Loss

① Calculate Loss

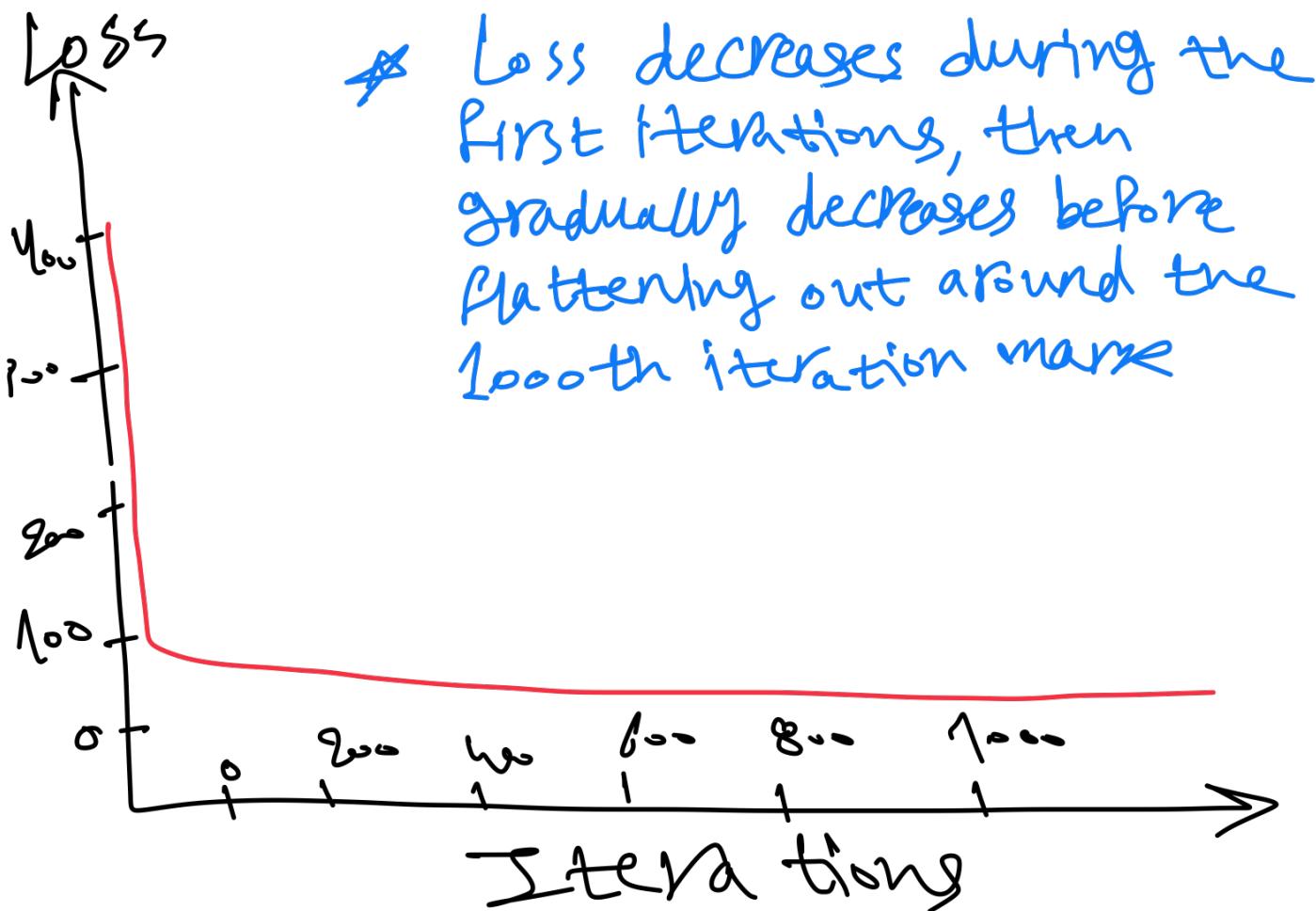


④ Repeat process till loss can't be reduced



③ Move a small amount in the direction that reduces loss

② Determine direction to move weight & bias

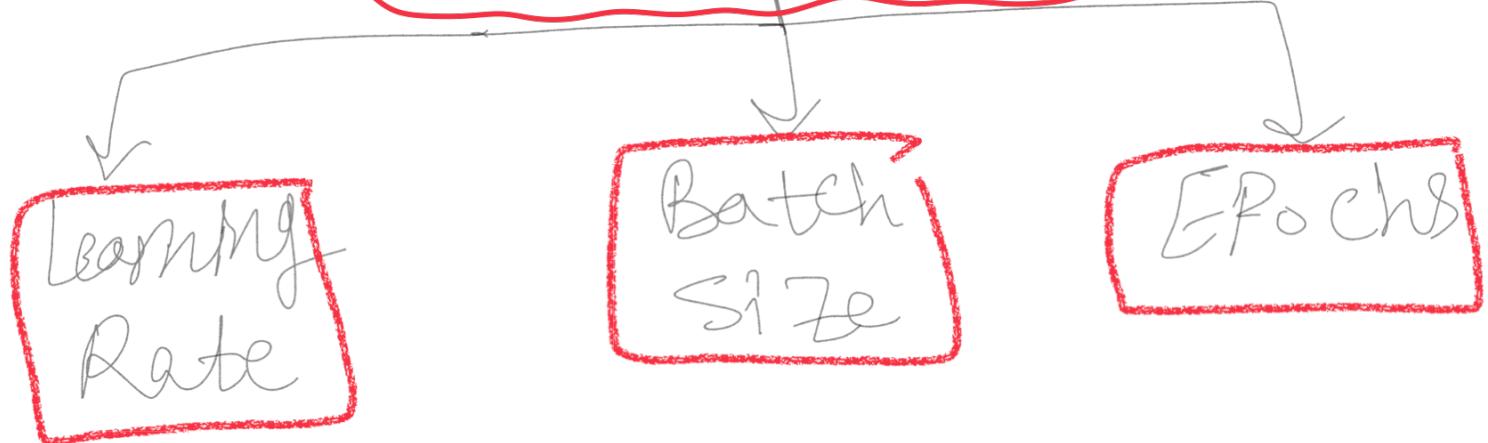


* After 1000th iterations, we can be mostly certain that the model has converged.

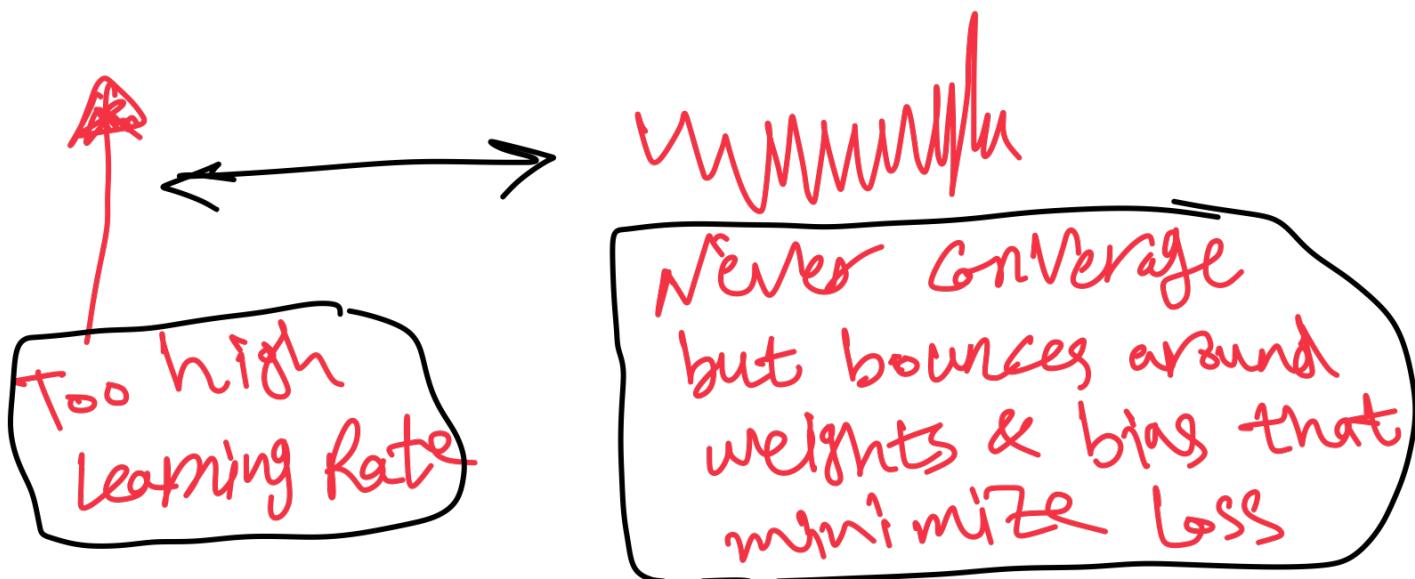
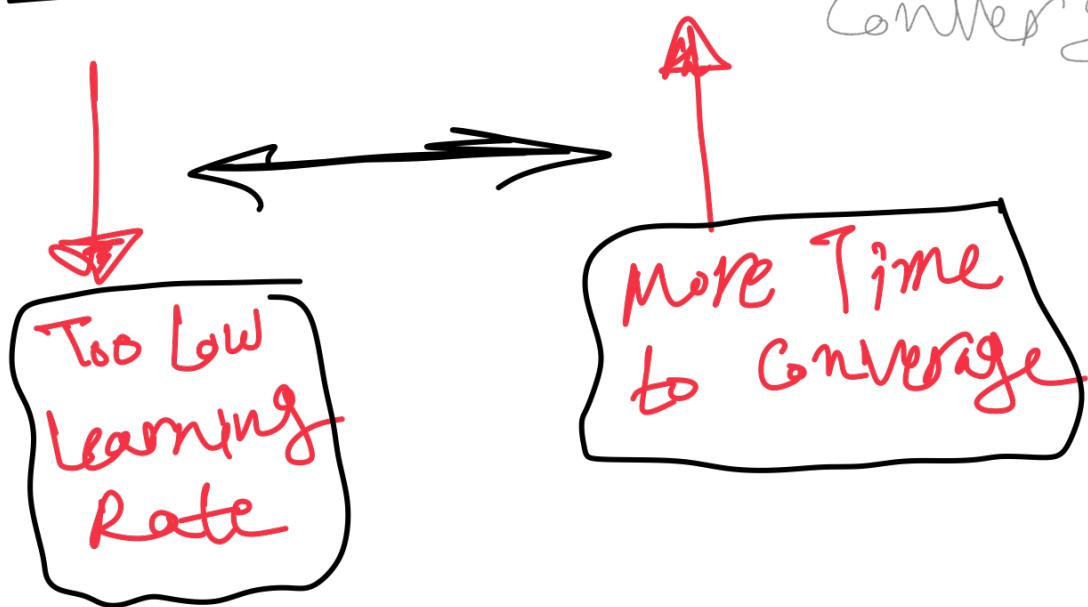
HyperParameters: are variables

that control different aspects of training.

HyperParameters



Learning Rate: Number refers to how quickly the model converges.



↑ ↓
Not too High
Not too Low

Converge quickly

Batch Size: Number of Examples the model processes before updating weights & bias.

Common Techniques

Stochastic gradient descent
SGD

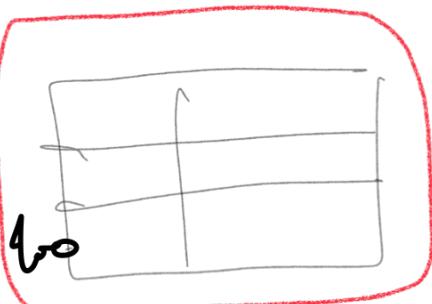
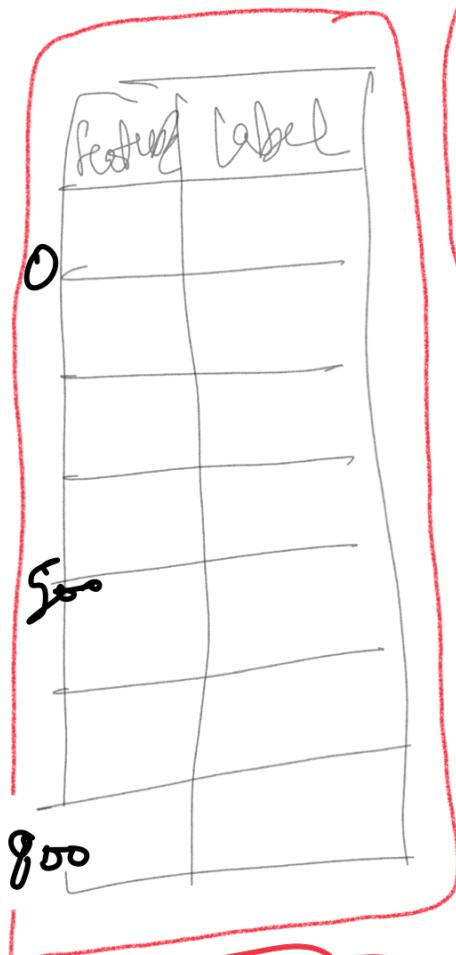
only one example per iteration

very noisy

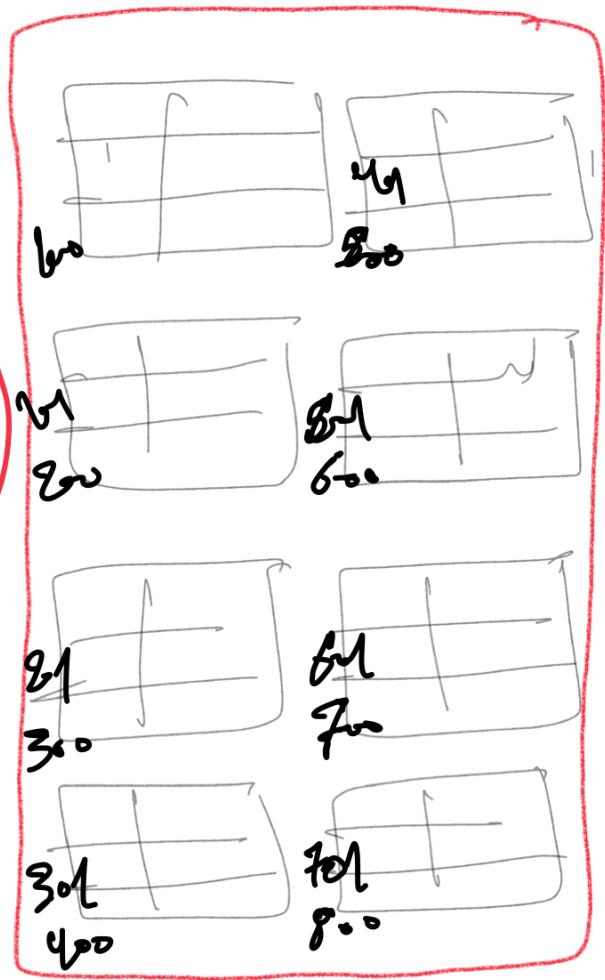
MINI-batch stochastic gradient descent

update weights & bias once per iteration

Epochs : The Model Processed every example in the training set once.



Single batch
Mini batch



Full batch
An entire
Data Set

one epoch ←
Full batch
composed of 8
batches

Batch type	when weight & bias updates occur
full batch	After the model looks at all the examples in dataset
SGD	After the model looks at single example from dataset
mini-SGD	After the model looks at the examples in each batch

Classification

Confusion Matrix: There are four possible outcomes for each output from a **binary classifier.**

→ class 1 (+)
→ class 2 (-)

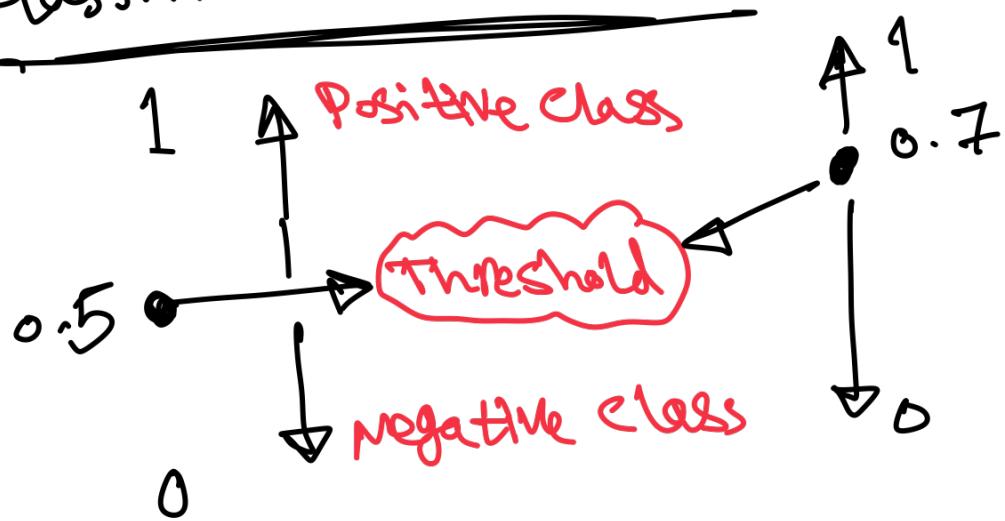
* In the Spam Email Spam/NotSpam

- Spam Email (Positive Class)
- NotSpam Email (Negative Class)

	Actual Positive	Actual Negative
Predicted Positive	True Positive (TP) Spam email correctly classified as Spam	False Positive (FP) A not Spam Email misclassified as Spam
Predicted Negative	False Negative (FN) Spam Email misclassified as not Spam	True Negative (TN) A not Spam Email correctly classified as not Spam

- * $TP + FP \Rightarrow$ All Predicted Positives
- * $FN + TN \Rightarrow$ All Predicted Negatives
- * $TP + FN \Rightarrow$ All Real Positives
- * $FP + TN \Rightarrow$ All Real Negatives

Classification Threshold:



Accuracy: The proportion of all classifications that were correct, whether positive or negative.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

Recall, or True Positive Rate: The proportion of all actual positives that were classified correctly as positives.

$$\text{Recall (TPR)} = \frac{TP}{TP + FN}$$

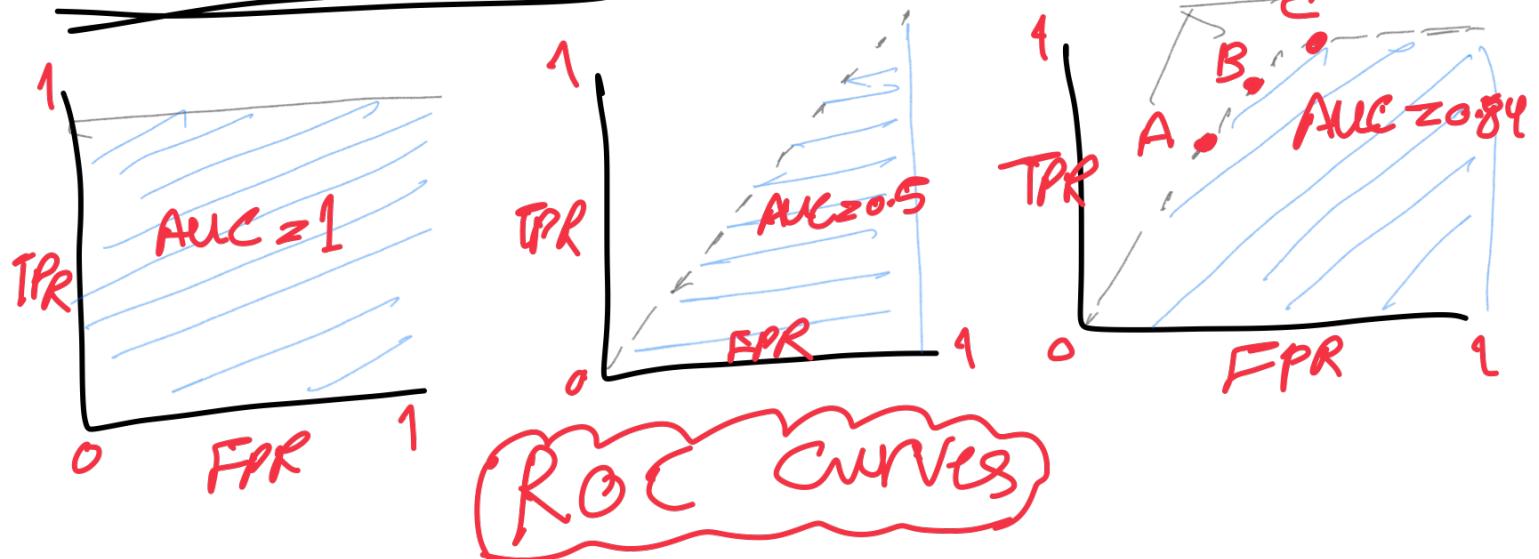
False Positive Rate: The proportion of all actual negatives that were classified incorrectly as positives.

$$\text{FPR} = \frac{FP}{FP + TN}$$

Precision: The proportion of all positive classifications that were actually positive.

$$\text{Precision} = \frac{TP}{TP + FP}$$

ROC & AUC:



@ A threshold: $\downarrow FP \& \downarrow TP$

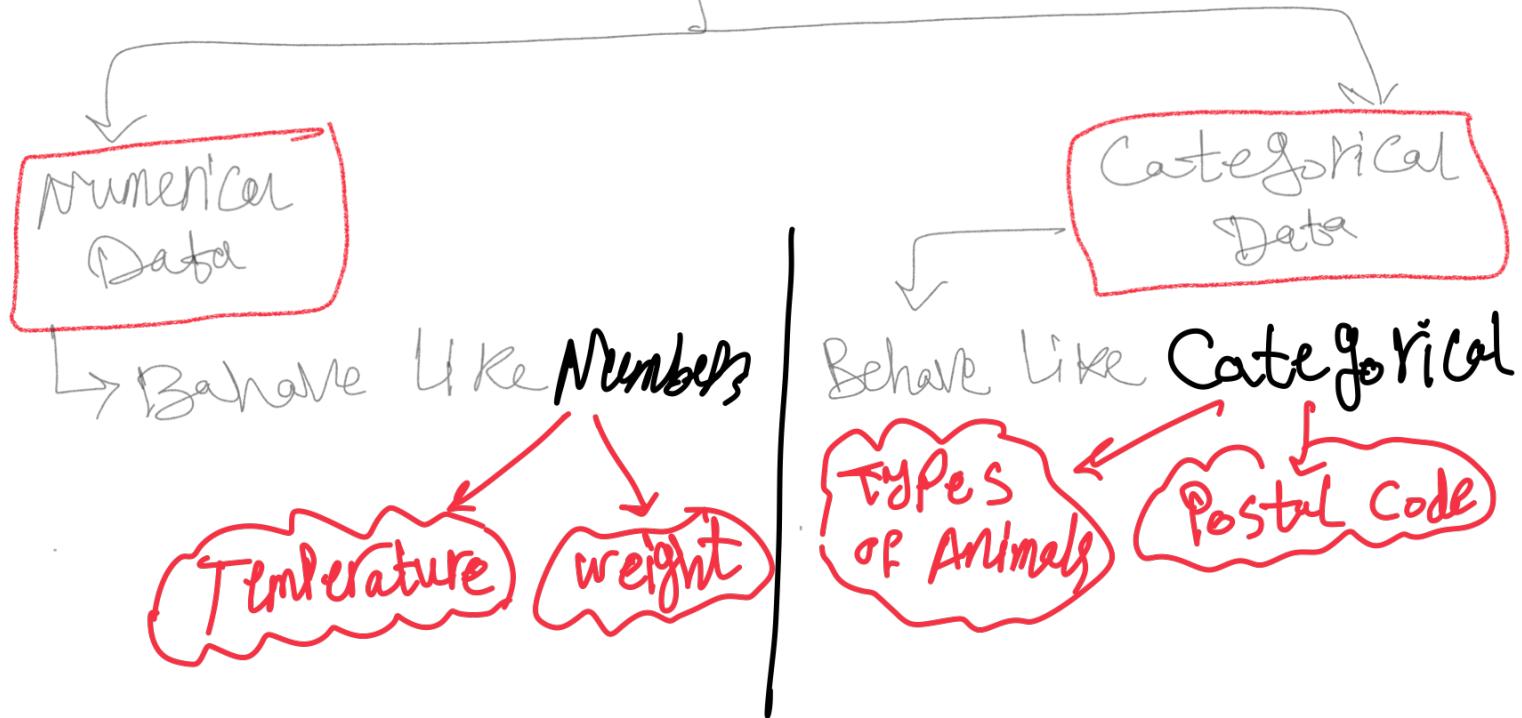
@ C threshold: $\uparrow FP \& \uparrow TP$

Prediction Bias: The difference between the mean of model's predictions and the mean of ground-truth labels in the data.

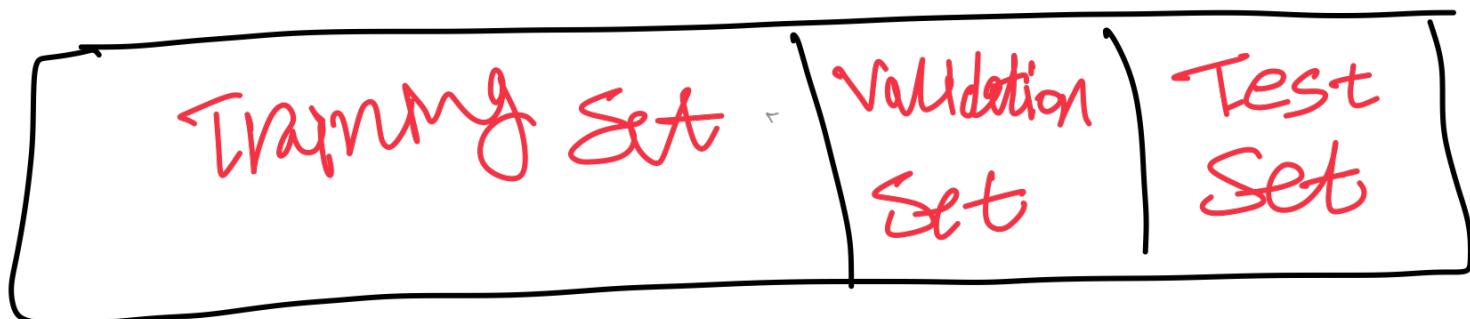
* Caused by:

- Biases or noise in the data.
- Too strong regularization.
- Bugs in the model training pipeline.
- The features provided to the model being insufficient for the task.

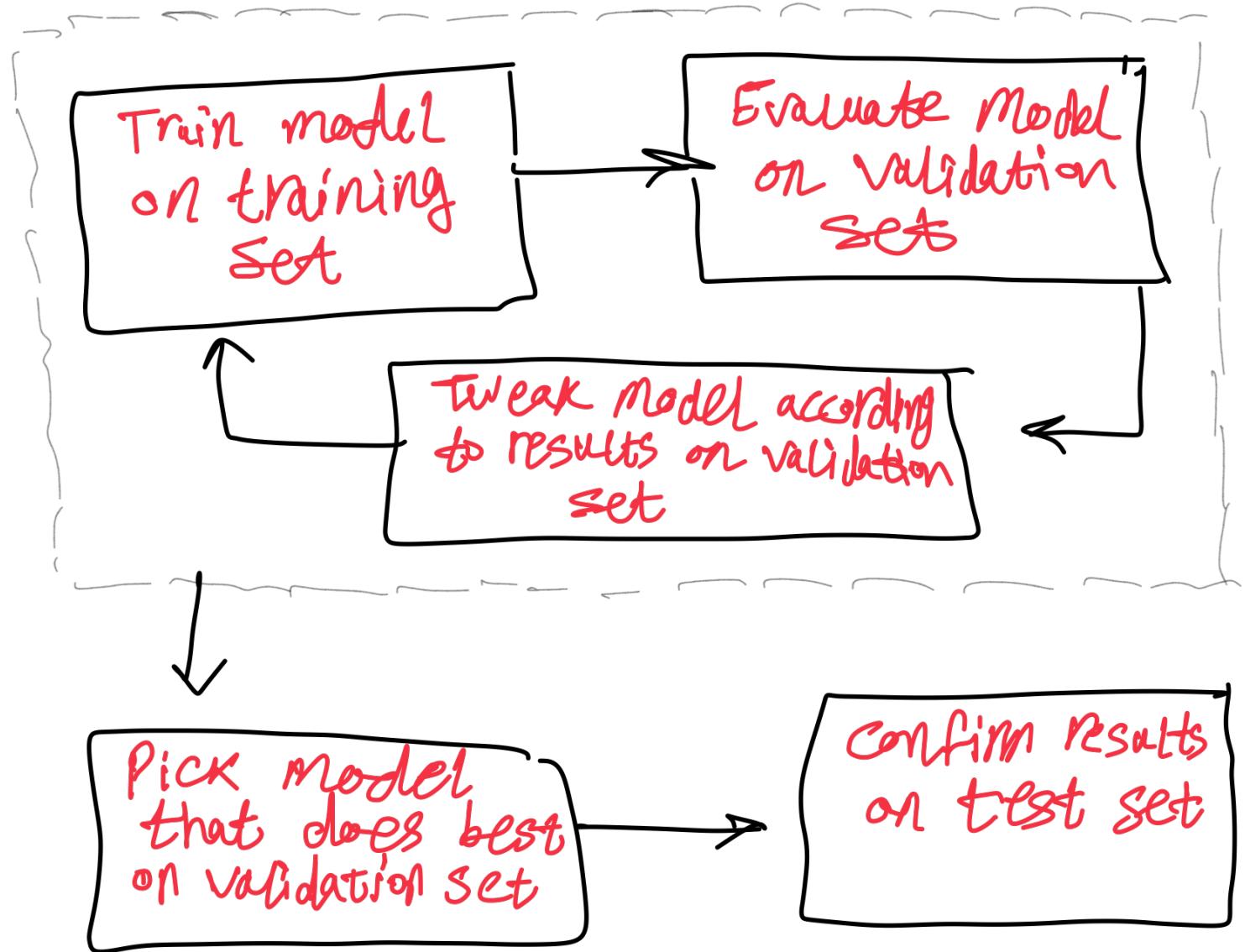
Data



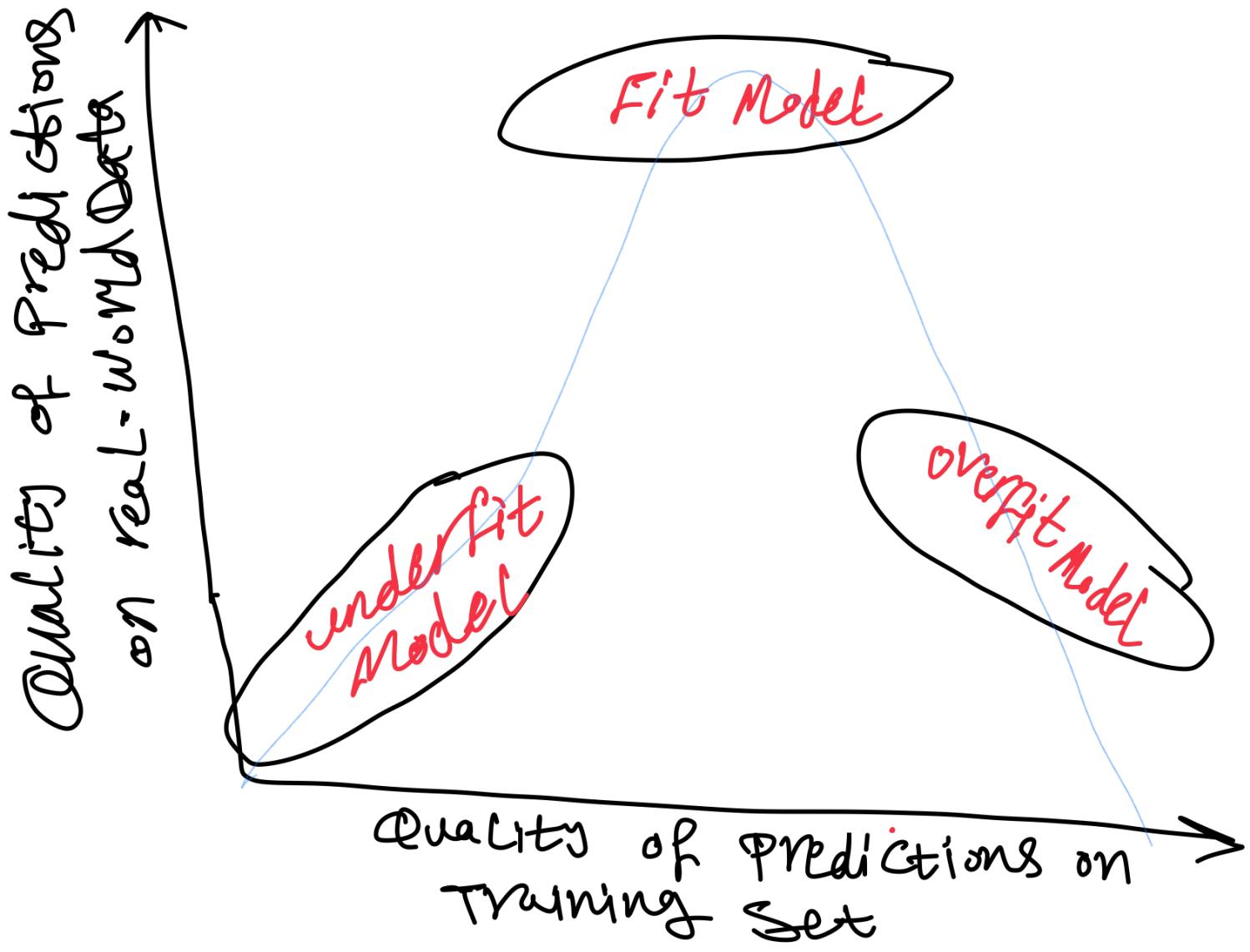
Dividing Dataset:



Workflow of Development & Testing:



overfitting: means creating model that memorizes the training set so closely but the model fails to make correct predictions on new data.



- * **overfit Model**: Best Prediction in training set, but bad predictions in new data
- * **underfit Model**: Bad predictions in both training set & new data

Generalization Curve that implies overfitting:



* After a number of iterations, loss declines or holds steady for training set, but increases for the validation set, this suggests **overfitting**

Overfitting Reasons:

1. Training set doesn't represent real life data.
2. The model is too complex.

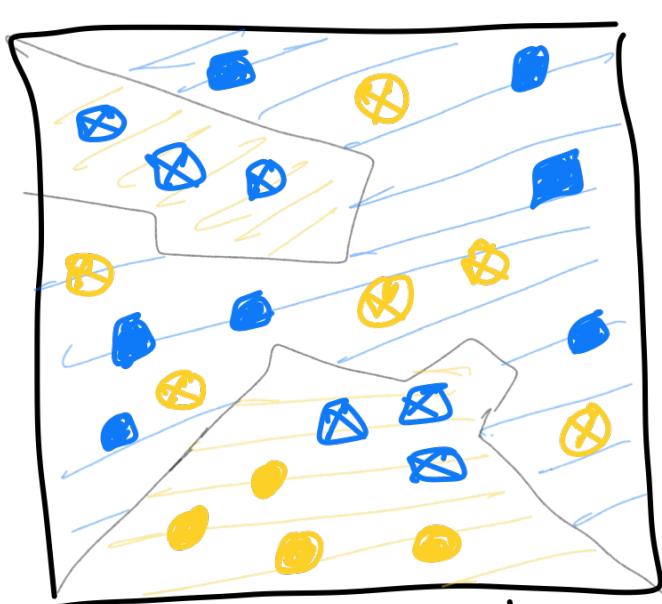
Generalization Conditions:

Model trains on training set, but the real test of a model's worth is how well it makes predictions on new examples on real-world data. Training a model that generalizes well implies the following dataset conditions:-

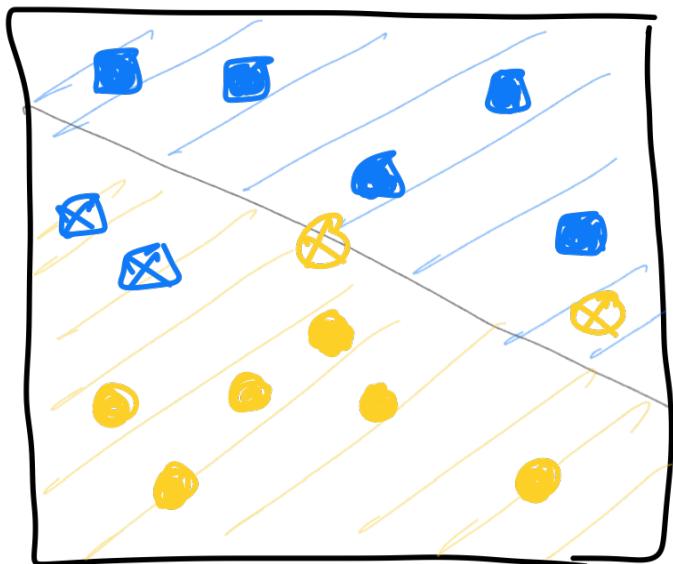
- Examples must be independently and identically distributed
- Dataset is stationary, meaning that dataset doesn't change significantly over time.
- Dataset Partitions have the same distribution, [Examples in training set are similar to examples in validation set, test set and real-world set]

* To have a similar distribution to the examples in all dataset types, you should shuffle the examples in the dataset before partitioning them.

Model Complexity:



Complex Model



Simple Model

Model predicts healthy trees in this region

Model Predicts sick trees in this region

- Model correctly predicted healthy tree.
- ✗ Healthy tree, but mistakenly model predicted it as sick tree.

- Model correctly predicted sick tree.
- ✗ Sick tree, but mistakenly model predicted it as healthy tree.

* Simple model generalizes better than complex model on new data, simple model made better predictions on the test set than complex model.

Regularization: Machine Learning

models must simultaneously meet two conflicting goals:

- Fit data well.
- Fit data as simply as possible.

* one approach to keep model simple is to penalize complex models, forcing & penalizing complex models during training is one form of regularization to become simpler.

Loss & Complexity: Loss & Complexity

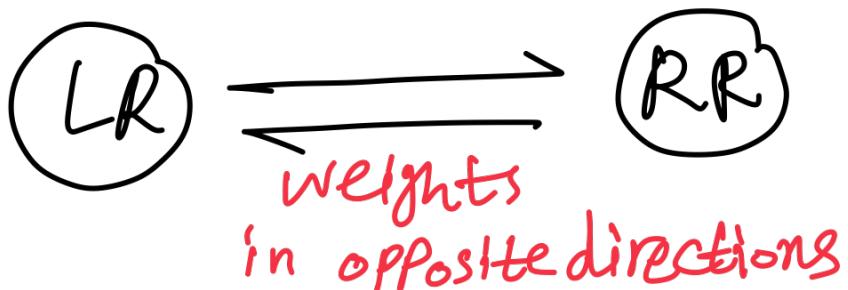
are inversely related, as complexity increases, loss decreases, as complexity decreases, loss increases, you should find a reasonable middle ground where model makes good predictions on both training data and real-world data.

Regularization Rate:

$$\text{minimize} (\text{loss} + \lambda \text{ complexity})$$

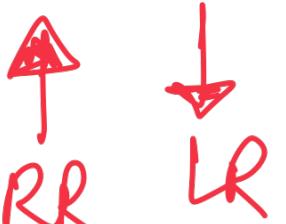
Reguralization Rate

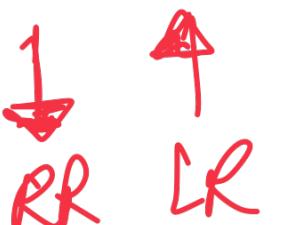
Equilibrium Between Learning Rate & Regularization Rate:



* High \underline{LR} \rightarrow pulls weights away from zero.

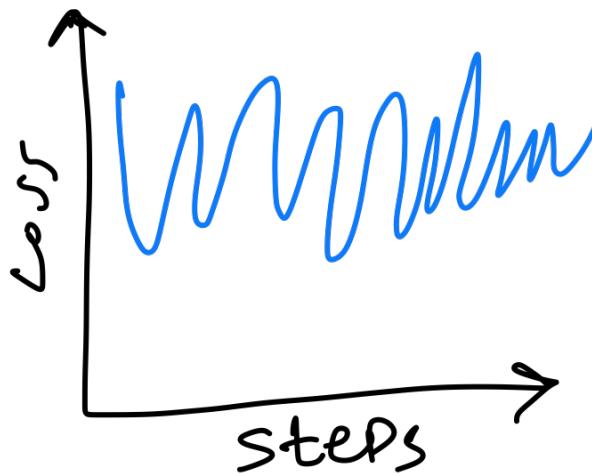
* High \underline{RR} \rightarrow pulls weights towards zero.

 Weak weights tend to produce a model that makes poor predictions.

 Strong weights tend to produce an overfit model

* Optimal case is to find balance between RR & LR

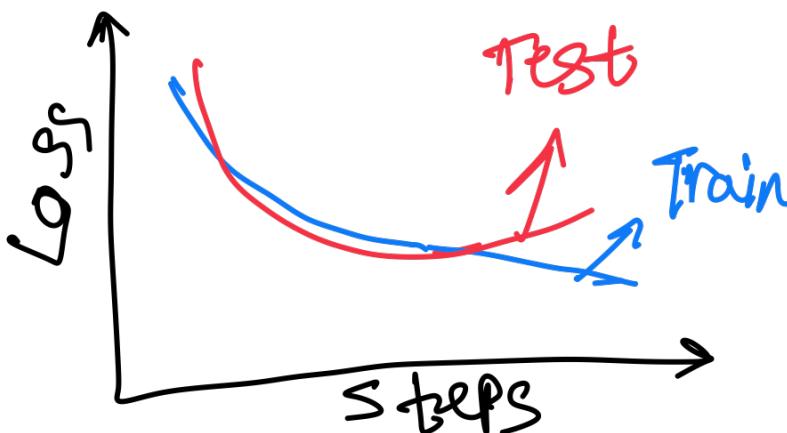
Interpreting Loss Curves:



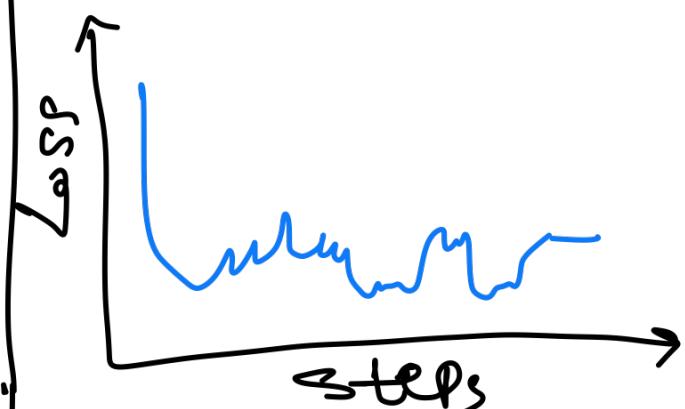
- * Training set has large number of untrusted examples.
- * High LR
- * There is bad examples in the data



- * Input data contains a burst of outliers
- * There is one or more NaNs in the input data (a value caused by a division by zero)



- * The model is overfitting the training set.



- * The training set contains repetitive sequences of examples

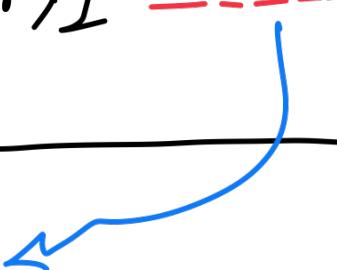
Neural Networks

Neural Networks are a family of model architectures designed to find nonlinear patterns in data.

Large Language Model

Language Model estimates the probability of a token or sequence of tokens occurring within a longer sequence of tokens. A token could be a word, subword or even a single character.

when I hear rain on my roof, I ---
in my kitchen



Probability	Token
9.4%	Cook soup
5.2%	warm up kettle
3.6%	Cover
2.5%	nap

- * A language model determines probabilities of different tokens or sequence of tokens to complete that blank in the previous sentence.
- * In some situations, the sequence of tokens could be an entire sentence, paragraph or even an entire essay.
- * An application can use the probability table to make predictions. Prediction might be the highest probability, for example (Cook soup) or random selection from tokens having probability greater than a certain threshold.
- * Estimating probability of what fills in the blank in a text sequence can be extended to more complex tasks, including:-
 - Generating text from one language to another
 - Translating text from one language to another
 - Summarizing documents.

N-gram Language Model: N-gram are ordered sequences of words used to build LM, where N is the number of words in the sequence.

You are very nice

2-grams based Model

- you are
- are very
- very nice

You are very nice

3-grams based Model

- you are very
- are very nice

* Given two words as input, LM based on 3-grams can predict the likelihood of the third word.

orange is

→ LM examines all the different 3-grams derived from its training corpus that start with orange to determine the most likely third word.

→ orange is ripe → is about orange fruit.

→ orange is cheerful → is about color orange.

Context:

Context is helpful information before or after the target token. It can help LM determine whether "orange" refers to fruit or color.

* Context Can help LM make better predictions, the only context the 3-gram provides is the first two words.

For example the two words **orange is** doesn't provide enough context for the LM to predict the third word. Due to lack of context, LM based on 3-grams make a lot of mistakes.

* Longer N-grams would certainly provide more context than shorter N-grams.

Recurrent Neural Networks:

Recurrent Neural Networks provide more context than N-grams. It is a type of neural network that trains on a sequence of tokens, for example, it can gradually learn selected context from each word in a sentence, kind of like you would.

Listening to someone speak: Recurrent neural network can gain context from a passage of several sentences.

- * Although recurrent neural networks learn more context than N -grams, the amount of useful context that can intuit is still relatively limited.
- * Recurrent Neural Networks evaluate information "token by token". In contrast the LM can evaluate "The whole context at once".

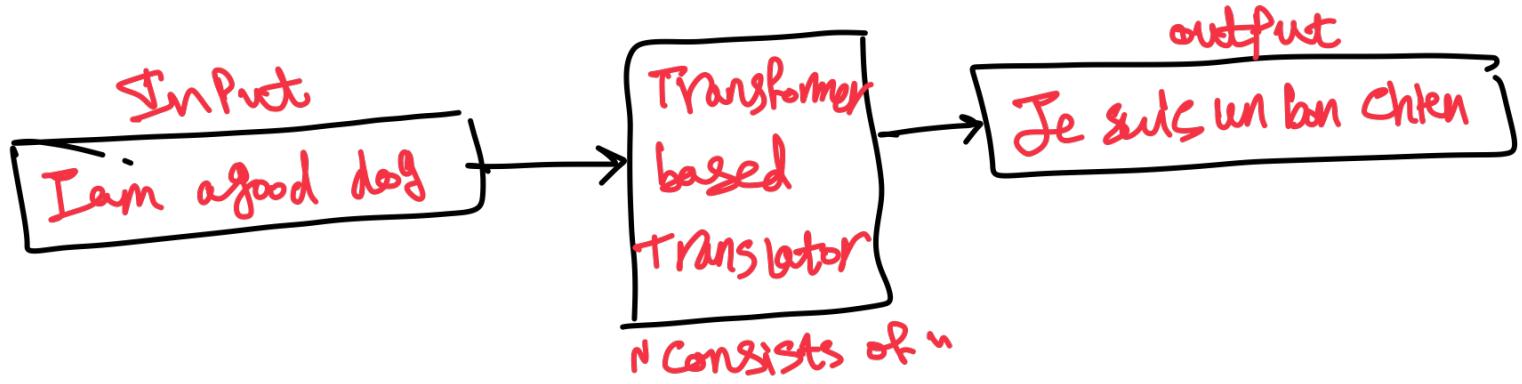
Large Language Model: LLMs predict token or sequence of tokens, sometimes many paragraphs worth of predicted tokens. Token can be a word, subword or even a single character.

- * LLMs make much better predictions than N -gram language models or recurrent networks.

- * LLMs contain more parameters than recurrent models.
- * LLMs gather far more context.

Transformer:

is the state-of-the-art architecture for a wide variety of LM applications.

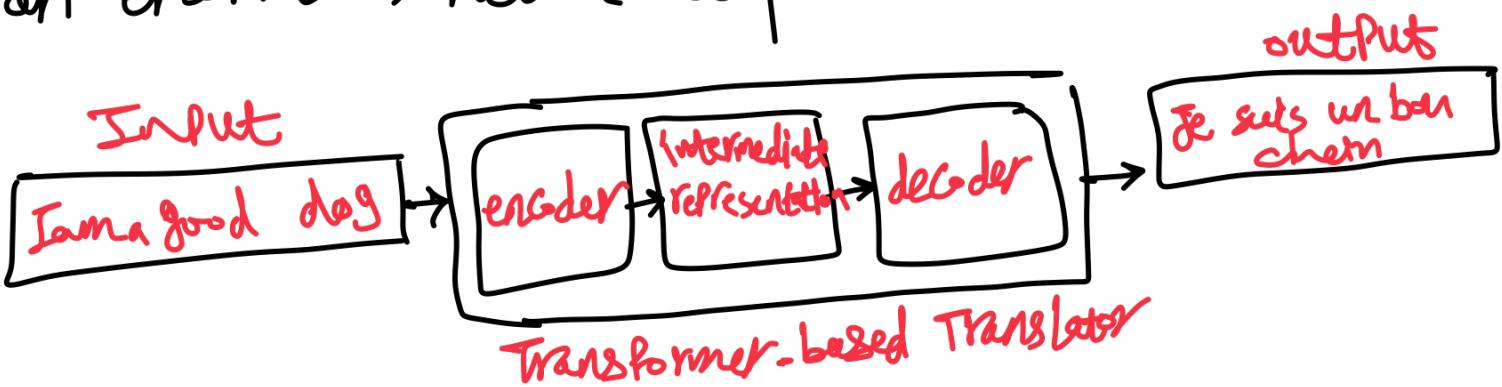


encoder

decoder

* Encoder converts input text into intermediate representation. Encoder is an enormous neural net.

* Decoder converts that intermediate representation into useful text. Decoder is an enormous neural net.



Self Attention:

Transformers rely heavily on concept called self-attention.

* Effectively on behalf of each token of input, Self Attention asks this question:-

How much does each other token of input affect the interpretation of this token.

The Animal didn't cross the street because it was too tired

→ Self Attention determines the relevance of each nearby word to the pronoun it.

[Animal is more important than street to the pronoun it]

The Animal didn't cross the street because it was too wide

→ [Street is more relevant than animal to the pronoun it]

How does LLM generate text?

LLMs are autocomplete mechanisms that can automatically predict thousands of tokens.

My dog, Max, knows how to perform many traditional dog tricks.

3.1%
2.9%

→ he can sit, stay and roll over.

→ he knows how to sit, stay and roll over.

Problems with LLMs:

- 1- Gathering an enormous training set.
- 2- Consuming multiple months and enormous computational resources and electricity.
- 3- Solving parallelism challenges.

Prompt Engineering: prompt Engineering enables

LLM's end users to customize model's output. That is, end users clarify how LLM should respond to their prompt.

* LLMs learn well from examples, showing one example to LLM is called **one-shot prompting**.

* Suppose you want model to use the following format to output a fruit's family

Peach: drupe
apple: ----

→ apple: Pome

One-shot Prompt shows LLM single example of a format and then asks LLM to complete a new based on that example

* In other situations, single example is insufficient, the user must show LLM multiple examples.

Plum: drupe
Pear: Pome
Lemon: ----

Providing multi examples is called **few-shot prompting**, you can think of the first two lines of the preceding prompt as training examples.

* Can LLM provide useful predictions with no examples **zero-shot prompting**? sometimes, but LLMs like context, without context, the zero shot prompt might return information about the technology company rather than fruit.

apple: ---- → Fruit ??
→ Company ??

{Production ML Systems}

Static Training (offline Training):

Train model only once. Then serve that same trained model for awhile.

Dynamic Training (online Training):

train model frequently. Then serve recently trained model.

	Static Training	Dynamic Training
Adv's	Simpler: You only need to develop and test model once.	More Adaptable: Your model keeps up with any changes to the relationship between features and labels.
Dis - Adv's	Sometimes Staler: If relationship between features and labels changes over time, your model's predictions will degrade.	more work: You must build, test and release a new product all the time.

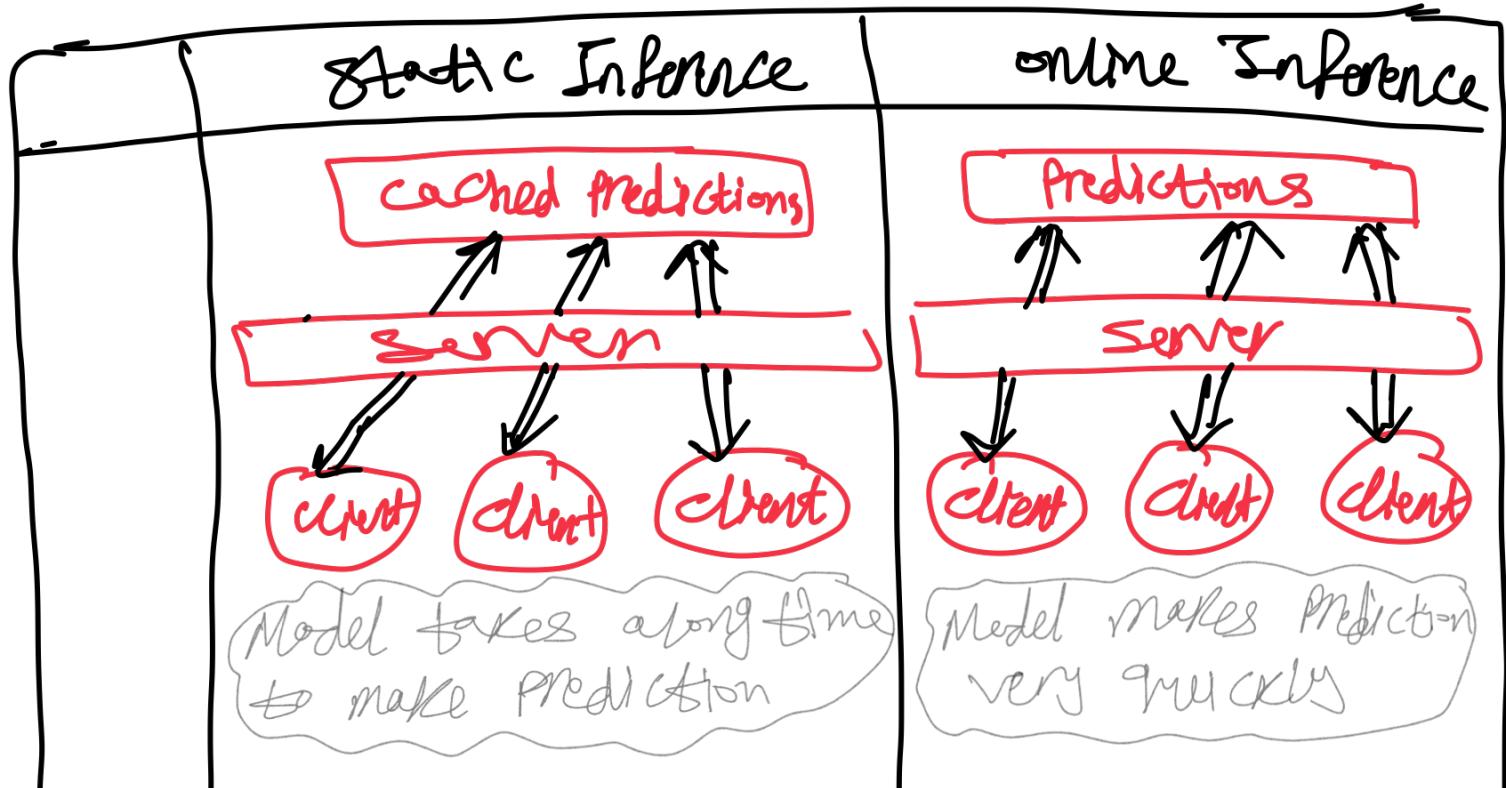
Interface: Interface is the process of making predictions by applying trained model to unlabeled examples.

Static Inference (offline Inference):

Model makes predictions on a bunch of common unlabeled examples and then caches those predictions somewhere.

Dynamic Inference (online Inference):

Model makes predictions on demand, for example when a client requests a prediction.



Advantages:

- Don't need to worry about cost.
- Can do post-verification of Predictions before pushing.

- Can infer prediction on any new item as it comes in, which is great for long tail Predictions.

	Static Inference	Dynamic Inference
DIS - ADM	<ul style="list-style-type: none"> - Can only serve cached predictions, so the system might not be able to serve predictions for uncommon input examples. - Update latency is likely measured in hours or days. 	<ul style="list-style-type: none"> - Compute intensive and latency sensitive - Monitoring needs are more intensive.

(Fairness)

Types of Bias:

ML models are not inherently objective. ML practitioners train models by feeding them a dataset of examples, and human involvement in the provision and curation of this data can make model's predictions susceptible to bias.

* When building models, it is important to be aware of common human biases that can manifest in your data, so you can take steps to mitigate their effects.

Reporting bias

DEFINITION

Reporting bias occurs when the frequency of events, properties, and/or outcomes captured in a dataset does not accurately reflect their real-world frequency. This bias can arise because people tend to focus on documenting circumstances that are unusual or especially memorable, assuming that the ordinary does not need to be recorded.

Click ➤ for an example.



Historical bias ↴

DEFINITION

Historical bias occurs when historical data reflects inequities that existed in the world at that time.

Click ➤ for an example



Automation bias

DEFINITION

Automation bias is a tendency to favor results generated by automated systems over those generated by non-automated systems, irrespective of the error rates of each.

Click > for an example



Selection bias

Selection bias occurs if a dataset's examples are chosen in a way that is not reflective of their real-world distribution. Selection bias can take many different forms, including coverage bias, non-response bias, and sampling bias.

Coverage bias

DEFINITION

Coverage bias occurs if data is not selected in a representative fashion.

Click > for an example



Non-Response bias

DEFINITION

Non-response bias (also known as **participation bias**) occurs if data ends up being unrepresentative due to participation gaps in the data-collection process.

Click > for an example



Sampling bias

DEFINITION

Sampling bias occurs if proper randomization is not used during data collection.

Click ➤ for an example



Group attribution bias

Group attribution bias is a tendency to generalize what is true of individuals to the entire group to which they belong. Group attribution bias often manifests in the two following forms.

In-group bias

DEFINITION

In-group bias is a preference for members of your own group *you also belong*, or for characteristics that you also share.

Click ➤ for an example



Out-group homogeneity bias

DEFINITION

Out-group homogeneity bias is a tendency to stereotype individual members of a group to which *you do not belong*, or to see their characteristics as more uniform.

Click ➤ for an example



DEFINITION

Implicit bias occurs when assumptions are made based on one's own model of thinking and personal experiences that don't necessarily apply more generally.

Click ➤ for an example



Confirmation bias

DEFINITION

Confirmation bias occurs when model builders unconsciously process data in ways that affirm pre-existing beliefs and hypotheses.

Click ➤ for an example



Experimenter's bias

DEFINITION

Experimenter's bias occurs when a model builder keeps training a model until it produces a result that aligns with their original hypothesis.

Click ➤ for an example

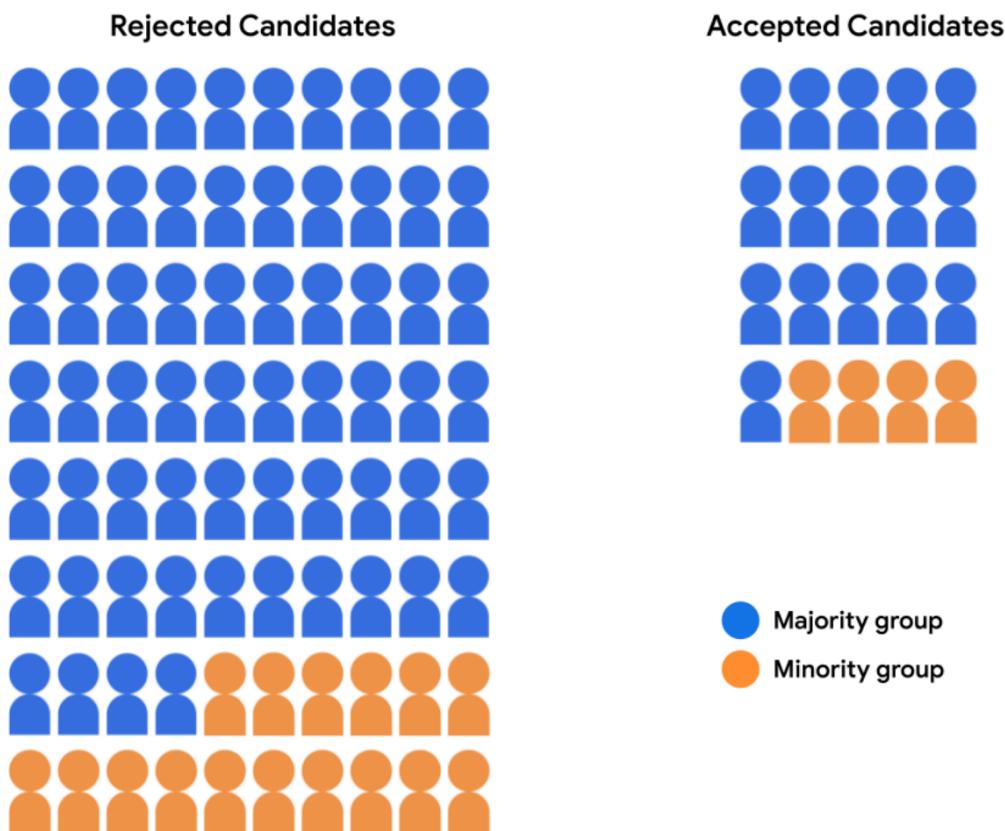


Fairness Metrics!



Demographic Parity: one method

we can use to evaluate model's predictions for fairness is to compare admissions rate for the majority group and the minority group. If the 2 admissions rates are equal, then the model's predictions exhibit demographic parity.



	Majority group	Minority group
Accepted	16	4
Rejected	64	16
Acceptance Rate	20%	20%

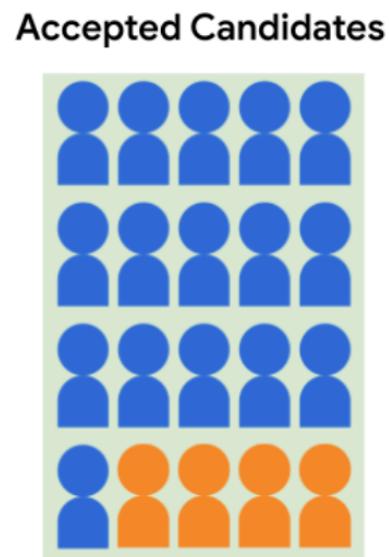
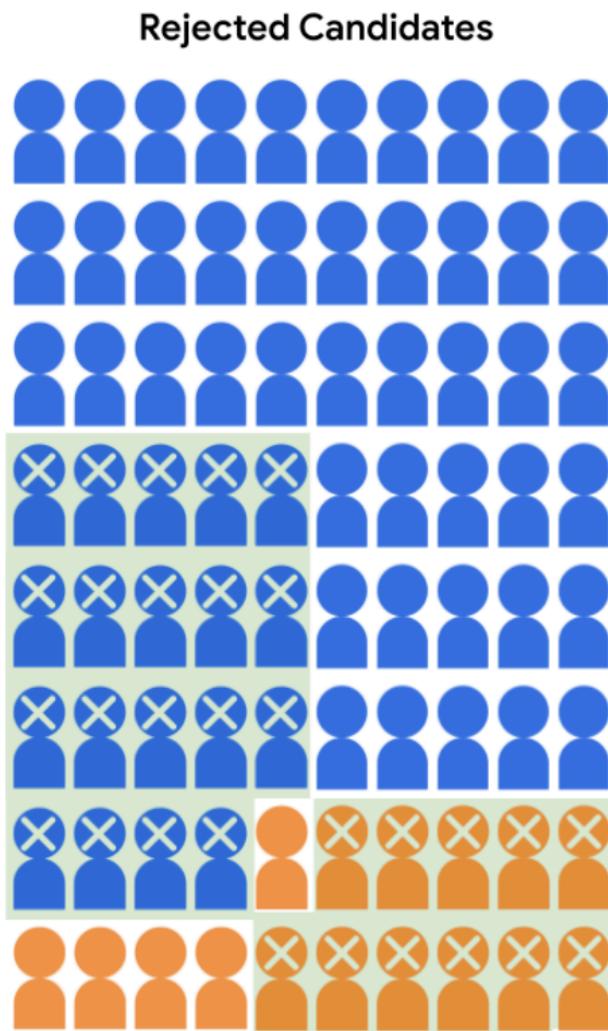
Benefits & Drawbacks :

Benefits

- Ensures that both Majority & Minority are represented in the admitted class of students in the same proportions as they are in the candidate pool.

Drawbacks

- It doesn't take the distribution of predictions of each demographic group (Qualified v.s Unqualified) into account when evaluate how 20 admission slots should be allocated.



- Majority group
- Minority group
- Actually qualified
- ✗ Qualified, but rejected

	Majority group		Minority group	
	Accepted	Rejected	Accepted	Rejected
Qualified	16	19	4	11
Unqualified	0	45	0	5

Majority acceptance Rate = $\frac{16}{35} = 46\%$

Minority acceptance Rate = $\frac{4}{15} = 27\%$

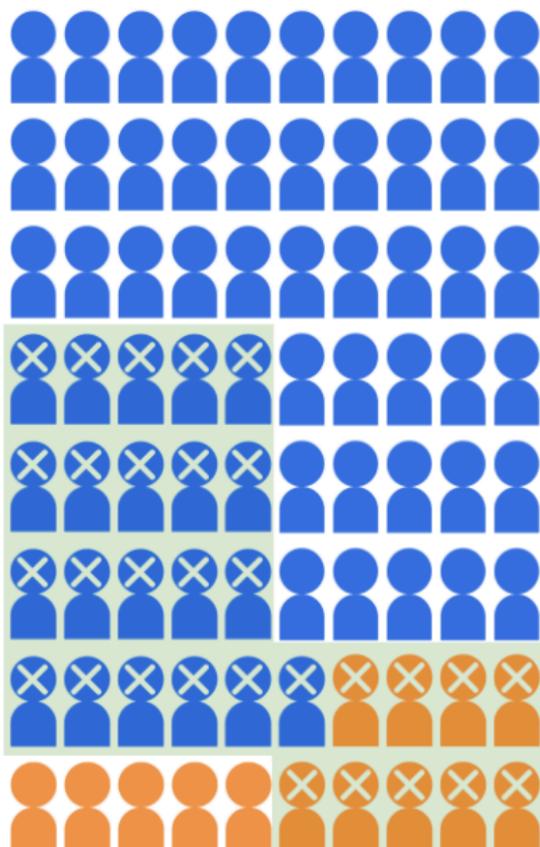
Equality of opportunity: We could

compare the acceptance rates for just the qualified candidates in the majority and the minority group. If the acceptance rates for qualified students in both groups are equal, the model exhibits

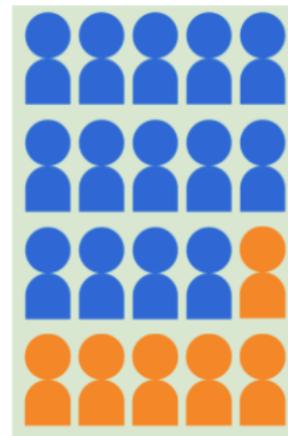
equality of opportunity

	Majority group	Minority group
Qualified	35	15
Unqualified	45	5

Rejected Candidates



Accepted Candidates



- Majority group
- Minority group
- Actually qualified
- ✗ Qualified, but rejected

	Majority group		Minority group	
	Accepted	Rejected	Accepted	Rejected
Qualified	14	21	6	9
Unqualified	0	45	0	5

Benefits & Drawbacks:

Benefits

- Allows model's ratio of positive to negative predictions to vary across demographic groups, provided that model is equally successful at predicting the preferred label for both groups.

Drawbacks

- Designed for use cases where there is clear-cut preferred label.
- It assesses fairness by comparing error rates in aggregate for demographic groups, which may not always be feasible.

* $\frac{14}{80} = 17.5\%$ "being Accepted in Majority"

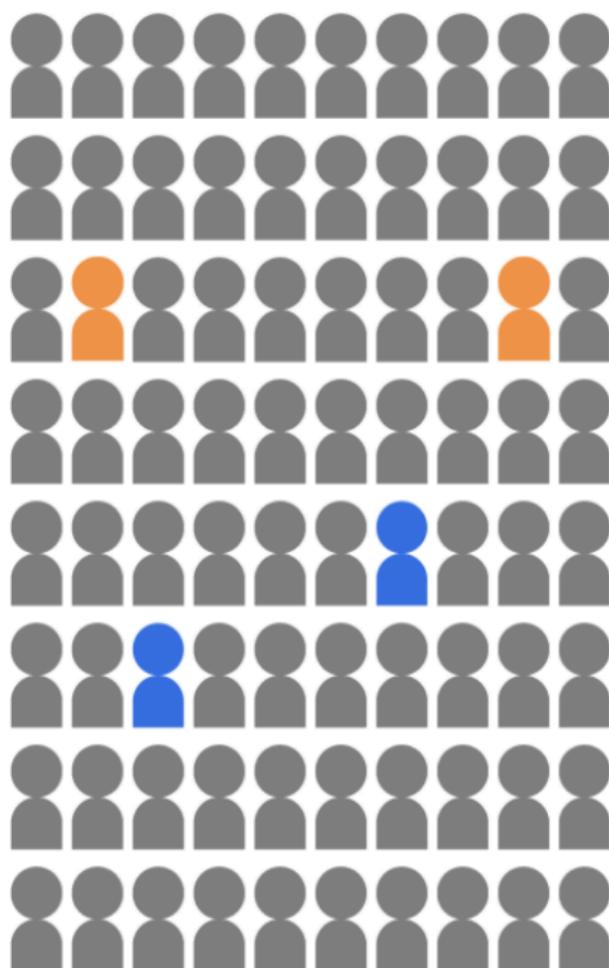
* $\frac{6}{20} = 30\%$ "being Accepted in Minority"

* 40% "being Accepted in Qualified in both groups"

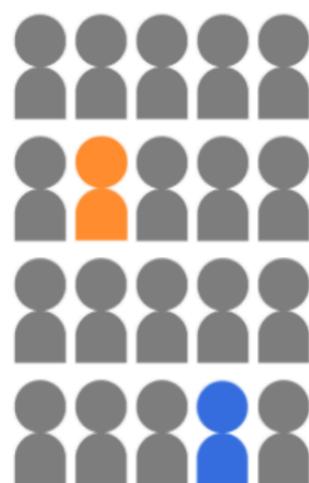
Counterfactual Fairness? Suppose that

our admissions dataset doesn't contain complete demographic data; instead, demographic group membership is recorded for just a small percentage of examples, such as students who opted to self-identify which group they belonged to.

Rejected Candidates



Accepted Candidates

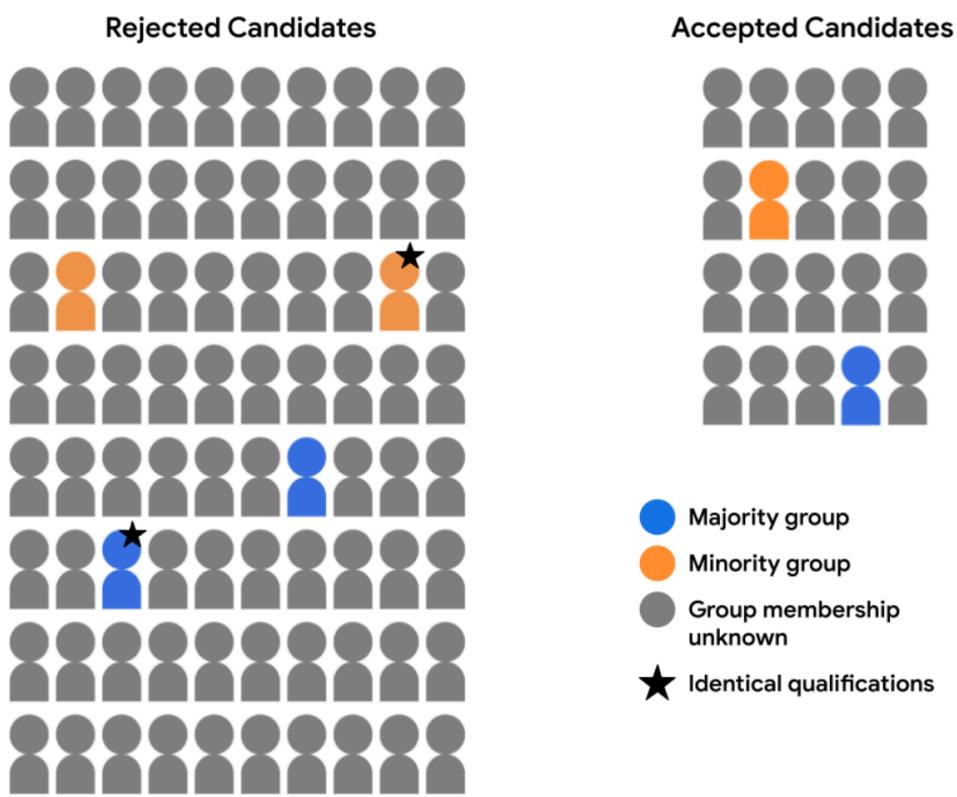


- Majority group
- Minority group
- Group membership unknown

* If we have thoroughly reviewed the feature data available for 2 candidates (annotated with a star in the image below), and have determined that they are identically qualified for admission in all respects. If the model makes same prediction for both of these candidates, (either rejects both or accepts both), it is said to satisfy Counterfactual Fairness,

For these examples.

* Counterfactual Fairness stipulates that 2 examples that are identical in all respects, except a sensitive attribute (Demographic group membership), should result in the same model prediction.



Benefits & Drawbacks:

Benefits

- Can be used to evaluate Predictions for Fairness in many cases where using other metrics wouldn't be feasible.

Drawbacks

- It doesn't provide as holistic a view of bias in model Predictions.