

INTERNATIONAL INSTITUTE OF INFORMATION
TECHNOLOGY, HYDERABAD



MACHINE LEARNING, QUANTUM COMPUTATION

MACHINE LEARNING AND PHYSICS IN THE VEIN OF OUR PRIMARY WORK

Research Notes

Author: Animesh Sinha

June 21, 2020

Contents

1 Bayesian Machine Learning	4
1.1 What is Bayesian	4
1.2 Bayesian Linear Regression	4
1.3 Latent Variable Models	4
2 Reinforcement Learning	5
2.1 What is Reinforcement Learning	5
2.1.1 Constituents of a RL System	5
2.1.2 Terms relation to the state	5
2.2 Exploration and Exploitation	6
2.2.1 Lower Bound on Regret	7
2.2.2 Upper Confidence Bound	7
2.2.3 Model-Based Bayesian Bandits	8
2.2.4 Policy Based Methods	9
2.3 Dynamic Programming to solve MDPs	10
3 Basic Neural Network Architectures in Keras	11
3.1 Common Code for Neural Nets	11
3.1.1 Train models with Saved intermediates	11
3.2 AutoEncoders	12
3.2.1 Convolutional Autoencoders	12
3.3 Recurrent Units	13
3.3.1 Long Short Term Memory (LSTM)	13
3.3.2 Gated Recurrence Unit (GRU)	14
4 Particle Physics for Data Scientists	15
4.1 The Preliminaries	15
4.1.1 Problems with the standard model	15
4.1.2 The Particles we want to detect	15
4.1.3 The Experiments in LHC	16
4.1.4 Simulation Package	16
4.1.5 Feynman Diagrams	17
4.2 The Large Hadron Collider Setup	17
4.2.1 Tracking System	17

4.2.2	Ring Imaging Cherenkov Detector (RICH)	18
4.2.3	Calorimeter - Electromagnetic and Hadronic	19
4.2.4	Muon Chambers	19
4.2.5	How to make a classifier uniform	19
4.2.6	Adversarial Neural Networks	20
4.3	Discovering New Physics	20
4.3.1	NoEther's theorem	20
4.3.2		20
4.3.3	Doping	20
5	Group Theory and Graphs	21
5.1	Groups and Subgroups	21
5.1.1	What is a Group?	21
5.1.2	Cayley's Table	22
5.1.3	Subgroups and GCD	22
5.1.4	Cyclic Groups	22
5.1.5	Abelian Groups	23
5.2	Special Groups and Their Properties	24
5.2.1	Permutation Groups	24
5.2.2	Cosets and Lagrange's Theorem	25
5.2.3	Normal and Factor Subgroups	27
5.3	Isomorphism and Homomorphism	27
5.3.1	Isomorphisms	27
5.3.2	Homomorphism	28
5.4	Rings and Fields	28
5.4.1	Rings	28
5.4.2	The Integer Domain	29
5.4.3	Fields	29

Chapter 1

Bayesian Machine Learning

1.1 What is Bayesian

1.2 Bayesian Linear Regression

1.3 Latent Variable Models

Latent variable is one that we cannot observe.

Chapter 2

Reinforcement Learning

2.1 What is Reinforcement Learning

2.1.1 Constituents of a RL System

- Environment: Has an internal state \tilde{S}_t (true state), Takes an action A_t and returns Observation O_t .
- Reward Signal: Takes the Observation O_t to returns a reward R_t not mutable by the agent.
- Agent (Learning algorithm): Has a State S_t , Value Function, Policy, Model (optionally).

Return is the sum of rewards over all steps: $G_t = G_{t-1} + R_t$.

The goal is to maximize the following:

$$q(s, a) = E[G_t | S_t = s, A_t = a] = E[R_{t_1} + R_{t_2} + \dots | S_t = s, A_t = a] \quad (2.1)$$

2.1.2 Terms relation to the state

State of the Environment

History is defined to be the set of all actions and observations with which we can recreate the current state: $H_t = A_0 O_0 A_1 O_1 \dots A_n O_n$

Fully Observable State are defined to have: $S_t = O_t$, that is the current observation represents the entire state of the game.

For Partially observable environments, we can represent the state using a Recurrent Neural Network with Memory to store some features of the history.

The state may be Markovian only to the environment and not to the agent, if there are additional unknowns for the agent.

Definition 1.1: Markov Decision Process

A decision process is Markovian if: $p(r, s \mid S_t, A_t) = p(r, s \mid H_t, A_t)$, i.e. the future is independent of the part given the present ($H_t \rightarrow S_t \rightarrow H_{t+1}$).

One of the reasons why an environment is not Markov is because it's non-stationary (and we don't want to or can't model how the environment evolved)

State of the Agent**Theorem 1.1: Bellman Equations**

$$v_\pi = E[E_t + 1 + \gamma G_t + 1 \mid S_t = s, A_t \sim \pi(s)] \quad (2.2)$$

$$= E[E_t + 1 + \gamma v_\pi(S_t + 1) \mid S_t = s, A_t \sim \pi(s)] \quad (2.3)$$

The same holds true for the optimal value function (this equation is independent of the policy):

$$v_* = \max_a [E_t + 1 + \gamma v_* G_t + 1 \mid S_t = s, A_t = a] \quad (2.4)$$

Agents often approximate the value function.

Model of the Environment

Transition probabilities:

$$P(s, a, s') = P(S_t + 1 = s' \mid S_t = s, A_t = a) \quad (2.5)$$

Immediate reward:

$$R(s, a) = E[R_t + 1 \mid S_t = s, A_t = a] \quad (2.6)$$

2.2 Exploration and Exploitation

In this section we are taking the example of the Multi-Arm Bandit, the major assumption is that the length of the episode is 1, our actions do not affect our state.

The true value of the action is: $q(a) = E[R_t \mid A_t = a]$.

Note 2.1: Q-Value Update

We have an approximation on the value, we can update as follows:

$$Q_t(a) = \frac{\sum_{n=1}^t R_n I(A_n = a)}{\sum_{n=1}^t I(A_n = a)} \quad (2.7)$$

$$Q_t(A_t) = Q_{t-1}(A_t) + \alpha_f (R_f - Q_{t-1}(A_t)) \quad (2.8)$$

2.2.1 Lower Bound on Regret

Regret is the value by which we got the rewards wrong:

$$G = \sum_t (\nu_* - q(A_t)) = \sum_a N(a) (\nu_* - q(a)) \quad (2.9)$$

We want to minimize the regret over the training-cum-test period (we can forever keep learning).

If we use methods like ϵ -Greedy (explore some times and exploit other times), our regret growth is still linear in the length of time, which is only constant factor better from learning nothing.

We can prove the lower bound that Asymptotic total regret is at least Logarithmic in the number of steps. This is formally described by the gap $\Delta_a = (\nu_* - q(a))$.

$$\lim_{t \rightarrow \infty} L_t \geq \log(t) \sum_{a|\Delta_a > 0} \frac{\Delta_a}{KL(p(r|a) \| p(r|a_*))} \quad (2.10)$$

The regret is still unbounded, but it's lower-bounded by logarithm in t .

2.2.2 Upper Confidence Bound

Hoeffding's Inequality is a concentration bound result. Given independent and identically distributed random variables X_1, X_2, \dots, X_n , let $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ be the sample mean. Then:

$$p((\mathbb{E}[X] - \bar{X}) < u) = e^{-2nu^2} \quad (2.11)$$

If we assume that are rewards are in the range $[0, 1]$, then we can derive the relation for the upper confidence bound for probability p (the maximum range we can be sure the rewards are in with probability atleast p):

$$\begin{aligned} p(q(a) \geq Q_t(a) + U_t(a)) &\leq e^{-2N_t(a)U_t(a)^2} \\ \Rightarrow U_t(a) &= \sqrt{\frac{-\log(p)}{2N_t(a)}} \end{aligned}$$

The root over the number of observations is due to the assumptions that the variables are IID, using central limit theorem our confidences are gaussians, the standard deviation for it is proportional to root of n .

We want to set p increasing proportionally to time. So we can replace p with n .

Algorithm 2.1: UCB algorithm

$$a_t = \operatorname{argmax}_{a \in \mathcal{A}} Q_t(a) + c \sqrt{\frac{\log t}{N_t(a)}} \quad (2.12)$$

c is a hyperparameter which represent learning rate of some sort (usually 1 to 2).

Setting $c = \sqrt{2}$, the Logarithmic lower bound (upto constant) is achieved. $L_t \leq 8 \sum_{a|\Delta_a > 0} \frac{\log(t)}{\Delta_a} + O(\sum_a \Delta_a)$.

2.2.3 Model-Based Bayesian Bandits

Bayesian Bandit models are parametrized distributions over rewards $P(R_t|\theta, a)$, for each action a given some parametrization of our model θ (like mean and variance of Gaussian, etc, based on belief state of model).

$$p_t(\theta|a) \propto p(R_t|\theta, a)p_{t-1}(\theta|a) \quad (2.13)$$

Parameters Θ are being used to learn the true value of the reward, and we can inject our prior knowledge into $p_0(\theta|a)$.

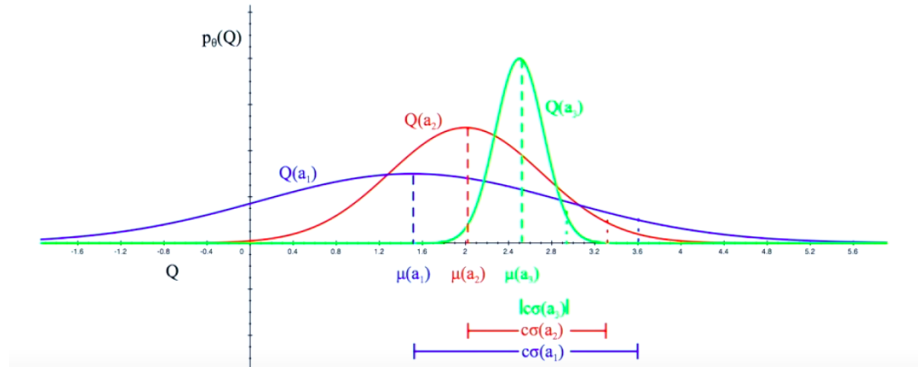


Figure 2.1: Example of the Bayesian Beliefs over Reward

If we have the case of Bernouli bandits, we can use the Beta distributions, which will be parametrized by the number of times we have gotten a zero and a one, i.e. $Q_t(a) = \beta(N_t(R = 0|a), N_t(R = 1|a))$.

Computing the posterior may be a lot harder for general examples. We have several ways to solve for the correct action given the posterior distributions.

Upper Confidence Bound - Bayesian

- We compute the Posterior over the action values ($p(R_t|\theta, a)$)
- We pick an upper confidence bound, eg. standard deviation ($U_t(a) = \sigma$)
- We act in a way that maximizes $Q_t(a) + U_t(a)$

Probability Matching

We will select each action with the probability of it being optimal.

$$\pi_t(a) = p(q(a) = \max_{a'} q(a') | H_{t-1}) \quad (2.14)$$

The probability may be difficult in general to get this probability analytically.

Thomson Sampling

We try to sample from the probability over the reward of each actions, and then select greedily.

$$\pi_t(a) = \mathbb{E}[\mathcal{I}(Q_t(a) = \max_{a'} Q_t(a'))] = p(q(a) = \max_{a'} q(a')) \quad (2.15)$$

This achieves the Logarithmic lower bound of regret.

Information State Space

If we can quantify the value of information gained in terms of future rewards that we gain given the information. This is a Markov Decision Process over the internal observation states. (We will need to know $p(\tilde{s}'|a, \tilde{s})$).

- Solvable using Model Free Reinforcement Learning (Q-Learning)
- Solvable using Model Based Bayesian Reinforcement Learning (Gittens Indices)

This can be unwieldy and not scalable.

2.2.4 Policy Based Methods

Below is the definition of the policy, soft-max over action values, where the values H are just preferences and do not bear a direct correspondance to the action/state values.

$$\pi(a) = \frac{e^{H_t(a)}}{\sum_b e^{H_t(b)}} \quad (2.16)$$

Now we want to do gradient ascent over the value of the rewards, updating our preferences.

$$\theta = \theta + \alpha \nabla_{\theta} \mathbb{E}[R_t | \theta] \quad (2.17)$$

Theorem 2.1: Log Likelihood trick

We are computing the gradient for the Bandits, this is also known as the Rein-

force trick (Williams, 1992).

$$\nabla_{\theta} \mathbb{E}[R_t | \theta] = \nabla_{\theta} \sum_a \pi_{\theta}(a) \mathbb{E}[R_t | A_t = a] \quad (2.18)$$

$$= \sum_a q(a) \nabla_{\theta} \pi_{\theta}(a) \quad (2.19)$$

$$= \sum_a q(a) \frac{\pi_{\theta}(a)}{\pi_{\theta}(a)} \nabla_{\theta} \pi_{\theta}(a) \quad (2.20)$$

$$= \sum_a \pi_{\theta}(a) q(a) \frac{\nabla_{\theta} \pi_{\theta}(a)}{\pi_{\theta}(a)} \quad (2.21)$$

$$= \mathbb{E}[R_t \frac{\nabla_{\theta} \pi_{\theta}(A_t)}{\pi_{\theta}(A_t)}] \quad (2.22)$$

$$= \mathbb{E}[R_t \nabla_{\theta} \log \pi_{\theta}(A_t)] \quad (2.23)$$

$$(2.24)$$

We have converted a gradient over expectation into an expectation over a gradient, so now we can sample from the gradient of the distribution to be able to perform stochastic gradient descent.

Stochastic Gradient Ascent is now $\theta = \theta + \alpha \mathbf{R}_t \nabla_{\theta} \log(\pi_{\theta}(\mathbf{A}_t))$

In the case of SoftMax over action preferences, we have:

$$H_{t+1}(a) = H_t(a) + \alpha R_t \frac{\partial \log \pi_t(A_t)}{\partial H_t(a)} \quad (2.25)$$

$$= H_t(a) + \alpha R_t (\mathcal{J}(a = A_t) - \pi_t(a)) \quad (2.26)$$

$$\Rightarrow H_{t+1}(a) = H_t(a) + \alpha R_t \times \begin{cases} (1 - \pi_t(a)) & \text{if } a = A_t \\ (-\pi_t(a)) & \text{if } a \neq A_t \end{cases} \quad (2.27)$$

So we are increasing the preferences of actions with higher rewards and pushing down the preferences of the actions even without taking them, therefore it's not the same as value.

We also note that the sum of probabilities over all actions is always 1, so we can add a baseline (constant independent of θ and a) to the reward without changing the preferences. But it does change the variance of the update. So **performance can improve by adding a baseline like negative of the mean till now.**

2.3 Dynamic Programming to solve MDPs

The discount factor affects the goal we are going after.

Chapter 3

Basic Neural Network Architectures in Keras

3.1 Common Code for Neural Nets

3.1.1 Train models with Saved intermediates

Following is the basic code to run when training the neural Network.

```
model.fit(x=X_train, y=X_train, epochs=25,
validation_data=[X_validation, Y_validation],
callbacks=[keras_utils.ModelSaveCallback(
    model_filename),
    keras_utils.TqdmProgressCallback()],
verbose=0,
initial_epoch=last_finished_epoch or 0)
```

And this is the code to load a presaved checkpoint of the model and start training from the last completed epoch. The `keras_utils` file is available in the snippets, and makes these callbacks available.

```
def load_checkpoint(last_epoch)
    model_filename = 'model.{0:03d}.hdf5'
    last_finished_epoch = None
    if last_epoch is not None:
        s = keras_utils.reset_tf_session()
        last_finished_epoch = 4
        model = keras.models.load_model(model_filename
            .format(last_finished_epoch))
```

Following this, we always save our weights in a file, and then load from it, as follows:

```
encoder.save_weights("encoder.h5")
decoder.save_weights("decoder.h5")
```

3.2 AutoEncoders

3.2.1 Convolutional Autoencoders

Here is the code for the setting up a Convolutional AutoEncoders. Follows a 4 layer Conv-Pool and then 1 Dense layer architecture to encode, then a dense layer followed by Transpose Convolutional layers to decode.

```
import tensorflow as tf
import numpy as np

def build_deep_autoencoder(img_shape, code_size):
    """
    Makes a Deep Autoencoder (Encoder-Decoder pair)
    :param img_shape: Shape of the image that is being
        encoded
    :param code_size: Number of the neurons in the
        bottleneck layer
    :returns: The full autoencoder model compiled with
        MSE loss.
    """
    encoder = tf.keras.models.Sequential([
        tf.keras.layers.InputLayer(img_shape),
        tf.keras.layers.Conv2D(filters=32, kernel_size
            =(3, 3), padding='same', activation='elu'),
        tf.keras.layers.MaxPooling2D(pool_size=(2, 2),
            padding='valid'),
        tf.keras.layers.Conv2D(filters=64, kernel_size
            =(3, 3), padding='same', activation='elu'),
        tf.keras.layers.MaxPooling2D(pool_size=(2, 2),
            padding='valid'),
        tf.keras.layers.Flatten(),
        tf.keras.layers.Dense(code_size),
    ], name='Encoder')
    decoder = tf.keras.models.Sequential([
        tf.keras.layers.InputLayer((code_size,)),
        tf.keras.layers.Dense((img_shape[0] // 4) * (
            img_shape[1] // 4) * 32),
        tf.keras.layers.Reshape((img_shape[0] // 4,
            img_shape[1] // 4, 32)),
        tf.keras.layers.Conv2DTranspose(filters=32,
            kernel_size=(3, 3), strides=2, activation='
            elu', padding='same'),
        tf.keras.layers.Conv2DTranspose(filters=
            img_shape[2], kernel_size=(3, 3), strides
            =2, activation='elu', padding='same'),
    ], name='Decoder')
    input_image = tf.keras.layers.Input(img_shape,
        name='Image_Input')
    code = encoder(input_image)
    output_image = decoder(code)
    autoencoder = tf.keras.models.Model(inputs=
        input_image, outputs=output_image)
```

```

    autoencoder.compile(optimizer="adamax", loss='mse')
    return autoencoder
model = build_deep_autoencoder(img_shape=(28, 28, 1),
                               code_size=100)
model.summary()
(x_train, y_train), (x_test, y_test) = tf.keras.
    datasets.fashion_mnist.load_data()
model.fit(x=x_train, y=x_train, validation_data=[
    x_test, x_test], epochs=5)

```

3.3 Recurrent Units

3.3.1 Long Short Term Memory (LSTM)

Schematic Sketch

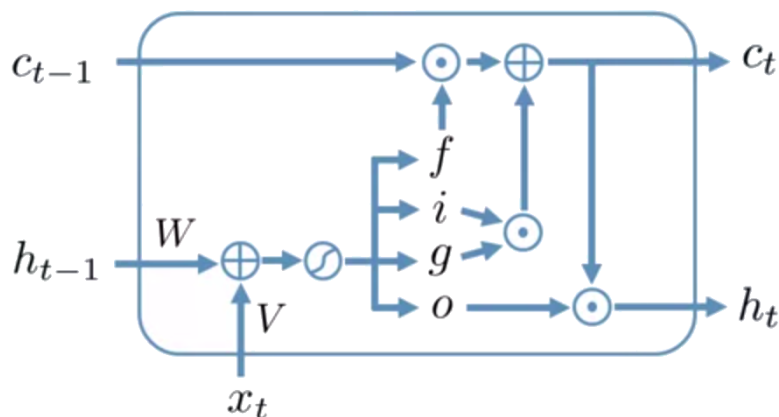


Figure 3.1: Schematic Diagram of LSTM

Update Equations

Following are the update equations for an LSTM. There are 3 gates, the Input Gate (which saves to cell memory), and Output Gate (which loads from cell memory) and a Forget Gate (which clears up cell memory). In addition a Gate Value is generated, which is what gets stored in the memory during input. All these get computed from the catenation of the previous hidden state and the current input.

$$g_t = \tilde{f}(W_g \cdot \text{Concat}[x_t, h_{t-1}] + b_g) \quad (3.1)$$

$$f_t = \sigma(W_f \cdot \text{Concat}[x_t, h_{t-1}] + b_f) \quad (3.2)$$

$$i_t = \sigma(W_i \cdot \text{Concat}[x_t, h_{t-1}] + b_i) \quad (3.3)$$

$$o_t = \sigma(W_o \cdot \text{Concat}[x_t, h_{t-1}] + b_o) \quad (3.4)$$

$$c_t = f_t \cdot c_{t-1} + i_t \cdot g_t \quad (3.5)$$

$$(3.6)$$

The hidden state is the output, which can be available as the entire sequence or only at the end. A dense layer or any other can be used on top of it to compute some other output y_t .

Management of Gradients

Since the memory cell connections are like a skip connection if there is no update/-forget operation, the gradients do not vanish as easily.

3.3.2 Gated Recurrence Unit (GRU)

Schematic Sketch

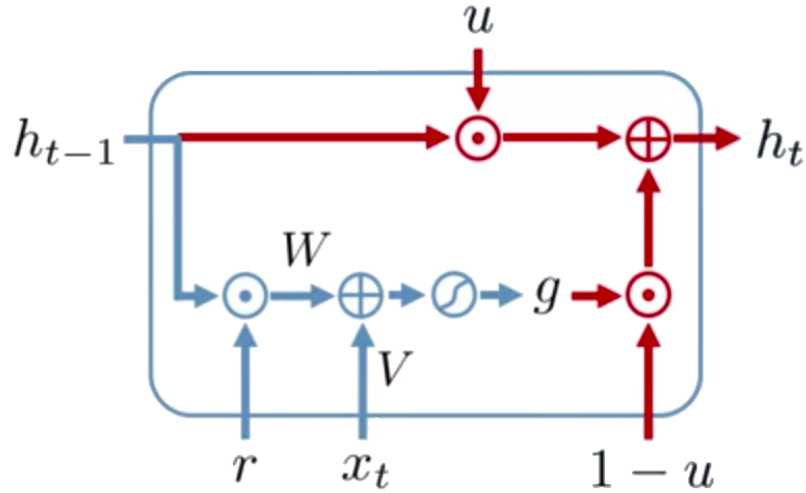


Figure 3.2: Schematic Diagram of GRU

Update Equations

Chapter 4

Particle Physics for Data Scientists

4.1 The Preliminaries

4.1.1 Problems with the standard model

What makes us unhappy?

- Matter and Antimatter inequivalence.
- 19 Arbitrary constants.
- Why is Gravity so weak.

4.1.2 The Particles we want to detect

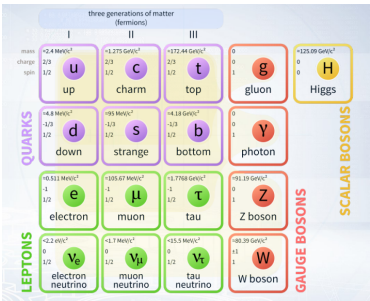


Figure 4.1: List of Particles of different types

The following are types of particles we want to detect.

- muon
- pion
- electron
- kon
- proton

4.1.3 The Experiments in LHC

There are 4 major detectors

- ALICE (A Large Ion Collider Experiment) is a heavy-ion detector on the Large Hadron Collider (LHC) ring. It is designed to study the physics of strongly interacting matter at extreme energy densities, where a phase of matter called quark-gluon plasma forms.
- ATLAS (A Toroidal LHC ApparatuS) is one of two general-purpose detectors at the Large Hadron Collider (LHC). It investigates a wide range of physics, from the search for the Higgs boson to extra dimensions and particles that could make up dark matter. Although it has the same scientific goals as the CMS experiment, it uses different technical solutions and a different magnet-system design. It has a cylindrical structure and measures particles in all directions.
- CMS (Compact Muon Solenoid) is a general-purpose detector at the Large Hadron Collider (LHC). It has a broad physics programme ranging from studying the Standard Model (including the Higgs boson) to searching for extra dimensions and particles that could make up dark matter. Although it has the same scientific goals as the ATLAS experiment, it uses different technical solutions and a different magnet-system design.
- LHCb (Large Hadron Collider Beauty) experiment specializes in investigating the slight differences between matter and antimatter by studying a type of particle called the "beauty quark", or "b quark". It is a single arm forward spectrometer.

We smash bunches of protons ('events'), record the pixels ('hits'), reconstruct trajectories ('jets', 'showers', 'tracks'), and we perform Statistical analysis on them.

A Trigger System is a system that uses criteria to rapidly decide which events in a particle detector to keep when only a small fraction of the total can be recorded.

4.1.4 Simulation Package

- <http://www.genie-mc.org/>
- <http://home.thep.lu.se/Pythia/>: Nutrino Simulations
- GEANT4: <http://geant4.web.cern.ch/>: Particles interacting with matter.
- FLUKA: <http://www.fluka.org/fluka.php>

4.1.5 Feynman Diagrams

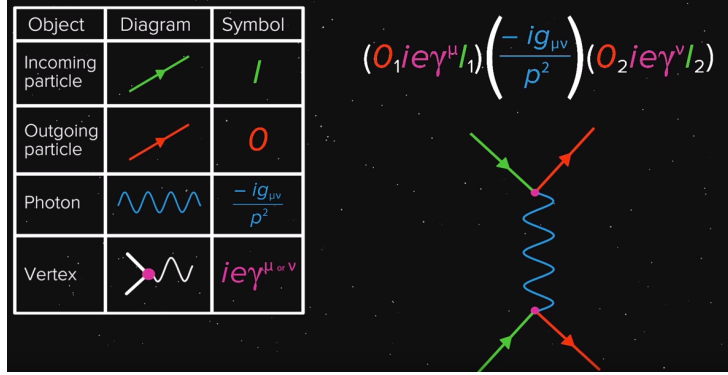
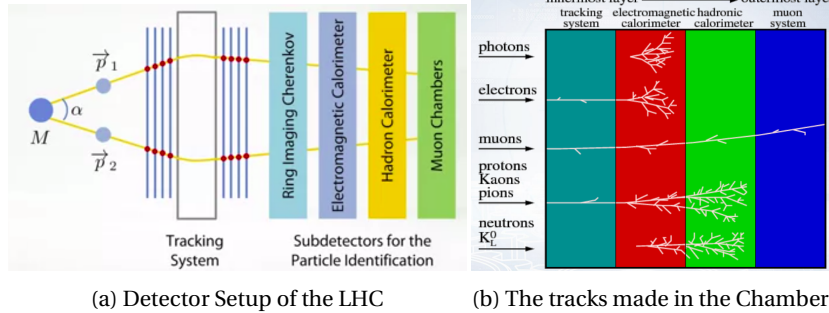


Figure 4.2: A Sample Feynman Diagram for scattering of Two electron by the transfer of one Photon

4.2 The Large Hadron Collider Setup



4.2.1 Tracking System

The first system of detectors, stands before particle collision area. Important for parameter estimation.

There are several layers of sensors that measure hits, and allow us to recognise particle trajectories. There is also a Magnetic Field that allows us to measure momentum (using radius of curvature in the field).

We have the following conservation equations when a particle D^0 with mass m breaks into a K^- with mass m_1 and a π^+ with mass m_2 , and they go away from each

other at angle α :

$$E_m = E_1 + E_2 \quad (4.1)$$

$$\hat{p}_m = \hat{p}_1 + \hat{p}_2 \quad (4.2)$$

$$E^2 = p^2 c^2 + m^2 c^4 \quad (4.3)$$

$$M^2 = m_1^2 + m_2^2 + \frac{2}{c^4} (E_1 E_2 - p_1 p_2 c^2 \cos \alpha) \quad (4.4)$$

Problem: Track Pattern Recognition

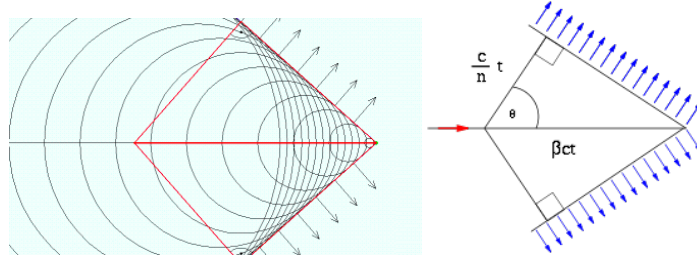
Recognizing hits that belong to the same track. Currently we have the following methods.

- Half Transform and Kalman Filtering. (Statistical, computationally cheaper.)
- Hopfield Neural Networks. (Denby Peterson and Cellular Automaton.)
- Convolutional Neural Networks (classify result as correct or wrong). Recurrent Neural Networks (predict the next hit location).

We can then combine these particle tracks into decays.

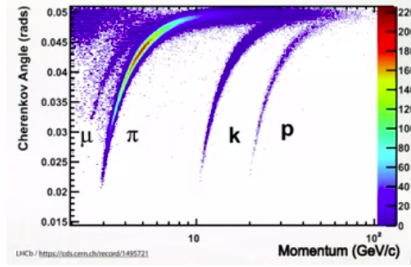
4.2.2 Ring Imaging Cherenkov Detector (RICH)

This is the first detector because it does not affect the flight of the particle.



(a) Cherenkov Simulation

(b) Cherenkov Angles



(c) Cherenkov Radiation Graphs

Figure 4.4: Registers in Processor Design

Here is the angle of the Chernokov radiation, derivable using simple geometric means and some special relativity in terms of the momentum.

$$p = \frac{mc\beta}{\sqrt{1 - v^2/c^2}} \quad (4.5)$$

$$\cos\theta = \frac{1}{n\beta} = \frac{\sqrt{p^2 + m^2c^2}}{np} \quad (4.6)$$

From the figure above, and the equation for it's analytical feel, that we can figure out the mass of the particle and thereby which particle it is.

4.2.3 Calorimeter - Electromagnetic and Hadronic

Electromagnetic stops all but than muons and quarks. Hadron calorimeter stops the quarks.

$$E_C = E_0 e^{\frac{x}{X_0}} \quad (4.7)$$

$$X_{max} = X_0 \ln\left(\frac{E_0}{E_c}\right) \quad (4.8)$$

$$N = E_0/E_c \quad (4.9)$$

After this we have crystals and scintillation counters, that measure the count and energy of particles, The electrons and photons, and since we measure energy, we already have momenta from the tracking system., this gives us the Mass, so we can identify them.

4.2.4 Muon Chambers

Muons do not interact well with matter. So we have multiple layers of metal to slow down the muons, they activate the chambers, we will extrapolate the data to get which particles in the tracking system were muons. We can also get the energy from how far they go in the muon chambers.

4.2.5 How to make a classifier uniform

The problem

$$1 - \text{false positive rate} = \text{background rejection} \quad (4.10)$$

$$\text{true positive rate} = \text{signal efficiency} \quad (4.11)$$

We want to have no dependence of the ROC AUC on all the momenta.

$$L_{ada} = \sum e^{-\gamma_i S_i} \quad (4.12)$$

$$L_{flat} = \sum_b w_b \int |F_b(s) - F(s)|^2 ds \quad (4.13)$$

$$L_{adaflat} = L_{ada} + \alpha L_{flat} \quad (4.14)$$

This is imposing an additional loss over the ADA boost function to make it less dependent on the momentum, alpha is a parameter we can tune to weight flatness over fit quality.

4.2.6 Adversarial Neural Networks

Minimizing dependencies on Mass, Momentum, etc. can also be done using adversarial neural networks. One network performs classification, and feeds its output to the other which tries to guess the momentum/mass from the output, if it can do this well, the model is bad, i.e. not flat. So, our new loss is $Loss_{classification} + f(Loss_{adversarial})$.

4.3 Discovering New Physics

4.3.1 NoEther's theorem

NoEther's theorem claims that the conservation laws we have are based on certain symmetries. Here is a list of some of them.

- Time Inversion: Energy
- Space Translation: Linear momentum
- Space Rotation: Angular momentum
- Charge Conjugation and Parity: Time Isotropy
- Charge, Lepton Number: Gauge Symetries

4.3.2

CvM test

$$CvM = \sum_{region} \int |F_{region}(s) - F_{global}(s)|^2 dF_{global}(S) \quad (4.15)$$

4.3.3 Doping

Chapter 5

Group Theory and Graphs

5.1 Groups and Subgroups

5.1.1 What is a Group?

A group is defined over a Set A and an arbitrary operation \times , denoted as: $\langle A | \times \rangle$.

- **Closure:** If a and b are elements in A , then $a \times b$ is also in A .
- **Identity:** There exists e such that $a \times e = a$.
- **Invertability:** There exists a^{-1} such that $a \times a^{-1} = e$.
- **Associativity:** $(a \times b) \times c = a \times (b \times c)$

A Subgroup is a subset of the original group that is itself a group.

One Step Subgroup Test states that if ab^{-1} is in the group H , then H is a subgroup of G .

Two Step Subgroup Test states that if a^{-1} is in H whenever a is in H and ab is in H for all a, b in H , then H is a subgroup of G .

Finite Subgroup Test If H is a non-empty finite subset of a group G , and H is closed under the operation G , then H is a subgroup of G .

Operations that Hold in groups are:

- Uniqueness of Identity (If $x \cdot a = x$ and $x \cdot b = x$, $(\forall x)$, then $a = b = e$)
- Uniqueness of Inverse (If $x \cdot a = e$ and $x \cdot b = e$, $(\exists x)$, then $a = b = x^{-1}$)
- Left and Right Cancellation (If $ab = ac$ then $b = c$. If $ba = ca$ then $b = c$.)

- Socks-Shoes Property $((ab)^{-1} = b^{-1}a^{-1})$

5.1.2 Cayley's Table

Cayley's Table is a 2-D matrix of all members of the group a and b on both axis and $a \cdot b$

5.1.3 Subgroups and GCD

Theorem 1.1: Equivalent Cyclic Subgroups

Let a be an element of order n in a group and let k be a positive integer. Then

$$\langle a^k \rangle = \langle a^{gcd(n,k)} \rangle \text{ and } |a^k| = n / gcd(n, k).$$

Proof 1.1: Equivalent Cyclic Subgroups

$(a^{gcd(n,k)})^\alpha = a^k$, Since $gcd(n,k)$ divides k . Thereby $\langle a^k \rangle \subseteq \langle a^{gcd(n,k)} \rangle$. Also, $gcd(n, k) = \alpha n + \beta k$, so $a^{gcd(n,k)} = a^{\alpha n + \beta k} = a^{\alpha n} a^{\beta k} = e \cdot a^{\beta k}$, therefore we can state that, $\langle a^{gcd(n,k)} \rangle \subseteq \langle a^k \rangle$. So we proved that $\langle a^k \rangle = \langle a^{gcd(n,k)} \rangle$.

Next, using the proof in the first part, since the groups are equal their orders are the same, so $|a^k| = |a^{gcd(n,k)}| = \frac{n}{gcd(n,k)}$, since the gcd divides n , it is the least solution to $(a^{gcd(n,k)})^x = a^n$

This has the following crucial corollaries.

- In a finite cyclic group, the order of an element divides the order of the group.
- Let $|a| = n$. Then $\langle a^i \rangle = \langle a^j \rangle$ if and only if $gcd(n, i) = gcd(n, j)$, and $|a^i| = |a^j|$ if and only if $gcd(n, i) = gcd(n, j)$.
- Let $|a| = n$. Then $\langle a \rangle = \langle a^j \rangle$ if and only if $gcd(n, j) = 1$, and $|a| = |\langle a^j \rangle|$ if and only if $gcd(n, j) = 1$.
- An integer k in Z_n is a generator of Z_n if and only if $gcd(n, k) = 1$.

These facts help us count the number of subgroups in a given set.

5.1.4 Cyclic Groups

Theorem 1.2: Fundamental Theorem of Cyclic Groups

Every subgroup of a cyclic group is cyclic. Moreover, if $|\langle a \rangle| = n$, then the order of any subgroup of $\langle a \rangle$ is a divisor of n ; and, for each positive divisor k of n , the group $\langle a \rangle$ has exactly one subgroup of order k , namely, $\langle a^{n/k} \rangle$.

Theorem 1.3: Number of Elements of Order D

Let G_n be a finite Cyclic Group of order N , if d is a positive Divisor of N , then the number of elements of order D in G of order d is $\phi(d)$.

For any group (i.e. Including non-cyclic), the number of elements of order d is a multiple of $\phi(d)$.

Proof 1.2: Number of Elements of Order D

If d is a divisor of n , then there is only one subgroup of order D of the group G_n . The elements in that group are those which are of the form $\langle a^x \rangle$, such that D and x are coprime. Therefore there can only be exactly $\phi(d)$ elements that are of order d .

If there are no elements of order D in the group, $\phi(d) \nmid 0$. Now let $\langle a_1 \rangle$ be a subgroup of order d , it has $\phi(d)$ elements in the subgroup of order d . So on for each $\langle a_i \rangle$ for all i , therefore a multiple of $\phi(d)$.

Theorem 1.4: Fundamental Theorem of Abelian Groups

Every finite Abelian group is a direct product of cyclic groups of prime order power. Moreover, the number of terms in the product and the orders for the cyclic groups are uniquely determined by the group. This is valid **upto isomorphism**.

Proof 1.3: Fundamental Theorem of Abelian Groups

Lemma 1 Let G be a finite Abelian Group of order $p^n m$, where p is a prime that does not divide m . Then $G = H \times K$, where $H = \{x \in G \mid x^{p^n} = e\}$ and $K = \{x \in G \mid x^m = e\}$. Moreover $|H| = p^n$

Lemma 2 Let G be an Abelian group of prime-power order and let a be an element of maximum order in G . Then G can be written in the form $\langle a \rangle = K$

Lemma 3 A finite Abelian group of prime-power order is an internal direct product of cyclic groups.

Lemma 4 G be a finite Abelian group of prime-power order. If $G = H_1 \times H_2 \times \dots \times H_m$ and $G = K_1 \times K_2 \times \dots \times K_n$, where the H 's and K 's are nontrivial cyclic subgroups with $|H_1| \geq |H_2| \geq \dots \geq |H_m|$ and $|K_1| \geq |K_2| \geq \dots \geq |K_n|$, then $m = n$ and $|H_i| = |K_i|$ ($\forall i$).

Example 1.1: Using Cardinality and Generators

Prove that every abelian group of order that is a product of primes is also cyclic. Given any abelian Group of order $p \cdot q$. If there is no prime such that

5.2 Special Groups and Their Properties

5.2.1 Permutation Groups

A Dihedral Group (D_n) is a group of all permutations of a n -sided Regular Polygon. The number of elements in this group is 2^n .

Definition 2.1: Permutation Groups

Permutation Groups are a group of permutations, where a permutation is a bijective function from a group to itself. Eg.

$$ExamplePermutation = \begin{bmatrix} A_0 & A_1 & A_2 & A_3 & A_4 & \dots \\ \alpha(A_0) & \alpha(A_1) & \alpha(A_2) & \alpha(A_3) & \alpha(A_4) & \dots \end{bmatrix} \quad (5.1)$$

Disjoint Cycle Notation Any permutation can be written as a product of disjoint cycles. Each cyclic subgroup is expressed as a separate disjoint cycle.

$$(1, 3)(2, 7)(4, 5, 6)(8) = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ 3 & 7 & 1 & 5 & 6 & 4 & 2 & 8 \end{bmatrix} \quad (5.2)$$

These Permutation Cycles can be multiplied together (that is operated by the Group Operator, Function Composition). eg.

$$\begin{aligned} P_1 \times P_2 &= (1,3)(2,7)(4,5,6)(8) \times (1,2,3,7)(6,4,8)(5) \\ &= (1,3)(2,7)(4,5,6)(8)(1,2,3,7)(6,4,8)(5) \\ &= (1,7,3,2)(4,8)(5,6) \end{aligned}$$

We go from Right to Left when multiplying, as function composition applies from Right to left. We see that $1 \rightarrow 2 \rightarrow 7, 7 \rightarrow 1 \rightarrow 3, 3 \rightarrow 7 \rightarrow 2, 2 \rightarrow 3 \rightarrow 1$ so we get the cycle $(1,3,7,2)$. We continue in this fashion to multiply. By Convention, we can choose not to write down single element cycles.

A Symetric Group (S_n) is a group of all the permutations over the operator Function Composition. There are $n!$ elements in S_n

Every Cycle can be written as a product of Disjoint Cycles Each disjoint cycle is a permutation, and it can always be written as a composition of two or more permutations. Also, any **disjoint cycles can always commute**.

Order of a permutation is defined as the Least Common Multiple of the lengths of the disjoint cycles.

Every Cycle can be Written as a Product of 2 cycles Any permutation is a composition of flips. Also by direct computation, we can prove this. Here is an example: $(01234)(56)(789) = (43)(42)(41)(40)(56)(97)(98)$

The representation of any cycle as a composition of 2 cycles can be classified as even or odd, i.e. any cycle C will require either even number of 2-cycles in all possible breakdowns or odd 2-cycles in all possible breakdowns, but it can never be that some 2-cycle decompositions are odd order and some are even order.

5.2.2 Cosets and Lagrange's Theorem

A Coset of a subgroup H in group G is the set $\{ah \mid H \in a\}$ - Left Coset, or $\{ha \mid H \in a\}$ - Right Coset.

Properties of Cosets of subgroup H in group G are:

1. $a \in aH$, since $e \in H$, so $ae \in aH$.
2. $aH = H$ iff $a \in H$, If $ah = H$ then $a = ae \in aH = H$, Also if $a \in H$ then $aH \subset H$ due to closure, $H \subset aH$ as $h = eh = aa^{-1}h = a(a^{-1}h) \in aH$, as a^{-1} is in H (invertability of a in H and then closure).
3. $(ab)H = a(bH)$ and $H(ab) = (Ha)b$, as Associativity holds.

4. **$aH = bH$ iff $a \in bH$** , as if $aH = bH$ then $a = ae \in aH = bH$. The other way, if $a \in bH$ then $aH = (bh)H = b(hH) = bH$ so they are equal sets.
5. **$aH = bH$ or $aH \cap bH = \Phi$** , Using Property 4, if there exists c in the intersection, then $c \in aH$ and $c \in bH$, so $aH = cH = bH$ thereby $aH = bH$.
6. **$aH = bH$ iff $a^{-1}b \in H$** , If there is $c \in aH, bH$ so we can say that $a^{-1}c \in H$ and $b^{-1}c \in H \iff c^{-1}b \in H$ due to invertability. Multiplying both elements in the we get $a^{-1}c \cdot c^{-1}b = a^{-1}b \in H$, again by Closure.
7. **$|aH| = |bH|$** , We can show a one-one map $aH \rightarrow bH$ using the cancellation property, i.e. $ah = bh (\forall h \in H)$.
8. **$aH = Ha$ iff $H = aHa^{-1}$** , right-multiply both sides by a^{-1} , so $aHa^{-1} = Haa^{-1} = H$.
9. **aH is a subgroup of G , iff $a \in H$** , aH must have identity e to be a subgroup, since $aH \cap eH \neq \phi \implies aH = eH = H$. By Property 2, $a \in H$. Conversely if $a \in H$ then $aH = H$ is a subgroup.

Theorem 2.1: Lagrange's Theorem

If G is a finite group and H is a subgroup of G , then $|H|$ divides $|G|$. Moreover, the number of distinct (left/right) cosets of H in G are $|G|/|H|$.

Proof 2.1: Lagrange's Theorem

Using the Properties proven above, $|aH| = |bH| = |eH| = |H|$ and $aH = bH$ or $aH \cap bH = \phi$, therefore $|G|$ is divisible by $|H|$.

The fruitfulness of cosets, when applied to Permutation Groups is displayed here:

Stabilizer of a Point: Let G be a group of permutations of Set S . For each i in S , let $stab_G(i) = \{\phi \in G \mid \phi(i) = i\}$. We call $stab_G(i)$ the stabilizer of i in G .

Orbit of a Point: Let G be a group of permutations of a set S . For each s in S , let $orb_G(s) = \{\phi(s) \mid \phi \in G\}$. The set $orb_G(s)$ is a subset of S called the orbit of s under G .

Eg.: $G = \{(1), (132)(465)(78), (132)(465), (123)(456), (123)(456)(78), (78)\}$

Orbits of G

$orb_G(1) = \{1, 3, 2\}$
 $orb_G(2) = \{2, 1, 3\}$
 $orb_G(4) = \{4, 6, 5\}$
 $orb_G(7) = \{7, 8\}$

Stabilizers of G

$stab_G(1) = \{(1), (78)\}$
 $stab_G(2) = \{(1), (78)\}$
 $stab_G(4) = \{(1), (78)\}$
 $stab_G(7) = \{(1), (132)(465), (123)(456)\}$

Orbit-Stabilizer Theorem states that $|orb_G(i)| \cdot |stab_G(i)| = |G| (\forall i)$.

5.2.3 Normal and Factor Subgroups

Normal Subgroups are subgroups H of group G , if $aH = Ha$ ($\forall a \in G$), written as $H \triangleleft G$. i.e., $ah = h'a$, the commutations are fudged a bit, when commuting we are allowed to use a different elements from H .

Normal Subgroup Test: A subgroup H of G is normal in G if and only if $xHx^{-1} \subseteq H$ ($\forall x \in G$).

Factor Groups (or Quotient Groups) are the groups, such that H is a normal Subgroup of G , $G/H = \{aH \mid a \in G\}$ \parallel $(aH)(bH) = abH$. This is because the normal subgroups are such that their left or right cosets are themselves subgroups. The operation is the composition of the operation that generated the cosets.

Internal Direct Product: We say that G is the internal direct product of H and K and write $G = H \times K$ if H and K are normal subgroups of G , $G = HK$, $H \cap K = \{e\}$

5.3 Isomorphism and Homomorphism

5.3.1 Isomorphisms

Definition 3.1: Isomorphism

A Isomorphism is a bijective mapping from one group onto another wherein if $a \times b = c$ then $\Phi(a) \times \Phi(b) = \Phi(c)$. That is the mapping preserve the result of every operation, and every inverse.

Automorphism is an Isomorphism that exists from a group onto itself.

Example 3.1: Disproving Isomorphisms

The Set of Real numbers can never be isomorphic to proper subset of itself under the operation addition.

The homomorphisms of real numbers to its proper subsets under addition are highly constrained: $\phi(a) + \phi(a) = \phi(2a) \implies \phi(m) = qm$. Therefore there can only be a multiplicative map that holds the homomorphism.

Any subgroup of rational numbers can be generated by a subset of the series if prime inverses: $\{\frac{1}{2}, \frac{1}{3}, \frac{1}{5}, \frac{1}{7}, \frac{1}{11}, \frac{1}{13}, \dots\}$, the set of all prime denominators. Since both of our subgroups are generated by a subset of this set, there must exist a map, just a multiplicative factor, ϕ , from $A \rightarrow B$. This never exists, since there is no such factor that maps the set of prime numbers to another set of prime

numbers (or the reciprocals thereof).

Example 3.2: Proving Isomorphisms

From $\langle \mathbb{R} | + \rangle$ to $\langle \mathbb{R} | \times \rangle$, show an isomorphism exists.

We have $\phi(x) = 2^x$, the mapping is bijective, and $\phi(x) \times \phi(y) = 2^x \times 2^y = 2^{x+y} = \phi(x+y)$

Cayley's Theorem states that every group is isomorphic to a group of permutations.

5.3.2 Homomorphism

Definition 3.2: Group Homomorphism

A Homomorphism ϕ is a mapping from a group G to \bar{G} is a mapping that preserves the group operation, i.e. $\phi(ab) = \phi(a)\phi(b)$ ($\forall a, b \in G$).
(The mapping may not be bijective, as opposed to Isomorphism)

Definition 3.3: Kernel of a Homomorphism

The kernel of a homomorphism ϕ from a group G to a group with identity e is the set $\{x \in G \mid \phi(x) = e\}$. The kernel of ϕ is denoted by $\text{Ker}\phi$.
(All elements in G that map to the identity, for isomorphism it's just the identity element - trivial subgroup)

5.4 Rings and Fields

5.4.1 Rings

A Ring is a set with 2 binary operations on it such that the following properties hold true. (Z is a arbitrary set, and $+$ and \times are arbitrary operations in $\langle Z | +, \times \rangle$, also $+$ is the first and \times is second operation on the ring).

- Associativity over $+$: $(a + b) + c = a + (b + c)$
- Commutativity over $+$: $a + b = b + a$
- Identity over $+$: $a + 0 = a$ must exist for some 0
- Invertable over $+$: $a + (-a) = 0$ must exist for some $-a$
- Associativity over \times : $a \times (b \times c) = (a \times b) \times c$.

- Distributivity of \times over $+$: $a \times (b + c) = (b + c) \times a = (a \times b) + (a \times c)$.

In summary, A Ring is a Abelian Group over $+$, and Associative and Distributive over \times .

A subring S of a ring R is the same operations defined over a subset of elements such that they in themselves form a ring.

Subring Test A non-empty subset of a ring is a subring if it is closed under subtraction ($a - b$) and multiplication ($a \times b$).

5.4.2 The Integer Domain

Rings were invented to abstract the algebraic properties of Integers. However they lose essential features in this abstraction, those of **Existence of Unity, Commutativity, and Cancellation**

Zero Divisors are elements of a Commutative Ring such that there is a non-zero element $b \in R$ with $ab = 0$.

Integral Domain is a commutative ring with unity and no zero-divisors. (It may also be defined as a Commutative Ring with Cancellation, Equivalent)

5.4.3 Fields

A Field is a commutative ring with unity in which every nonzero element is a unit (i.e. Invertible).

Every finite integral domain is a field. (Finite Order: So every element is invertible)

Characteristic of a Ring is the least positive integer n such that $nx = 0$ for all x in the Ring. If no such integer exists then it is 0.

If R is a ring with unity (1), then the characteristic of n is the order of 1, unless its ∞ in which case characteristic is 0.

Characteristic of an integral domain can only be 0 or prime. Its because if the field is finite, then $0 = n \cdot 1 = (p \cdot 1)(q \cdot 1)$. So, either $p \cdot 1 = 0$, or $q \cdot 1 = 0$. Since p or q are smaller than n , then p or q is the characteristic, not n .