

[Open in app](#)

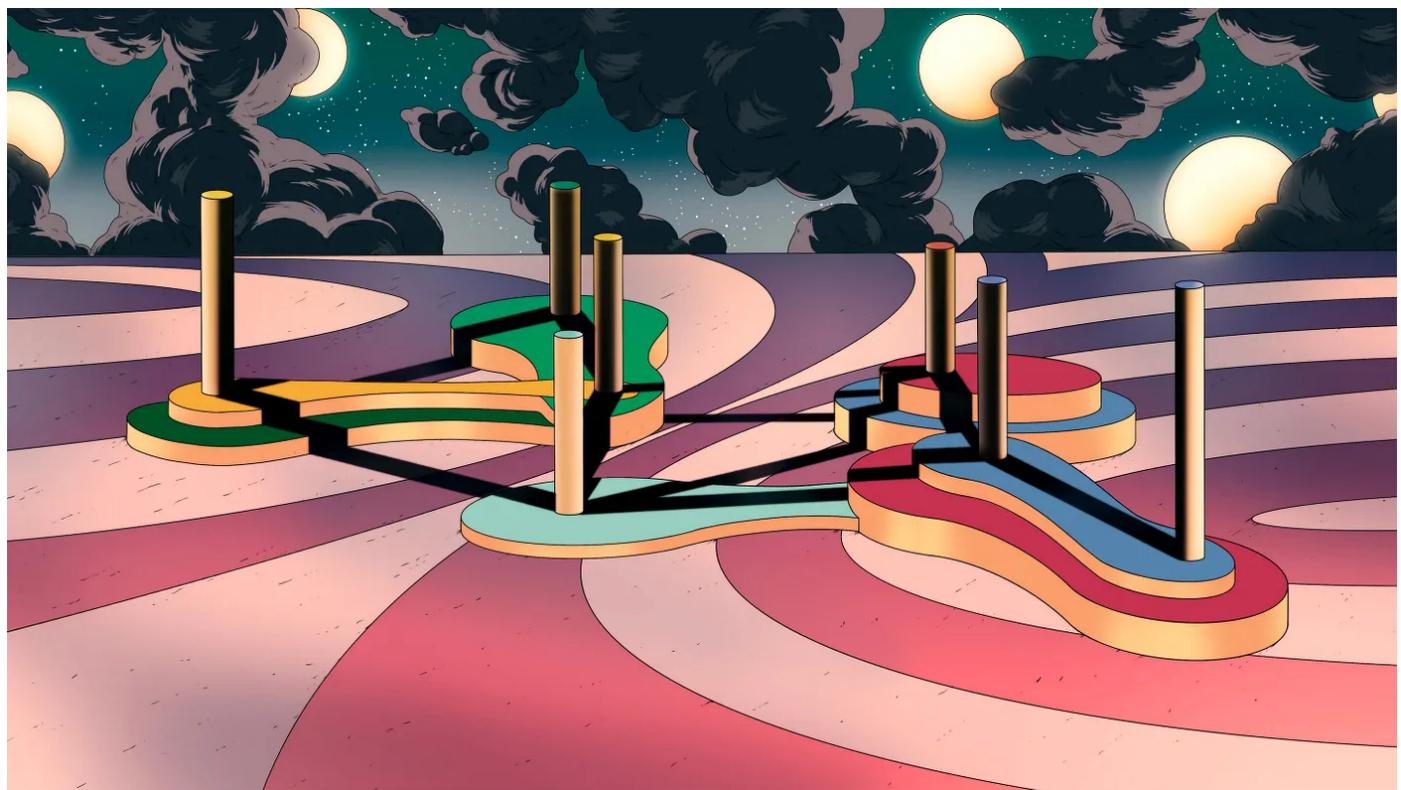
Search Medium



Machine Learning Coding Problems

Abhijit Mondal · [Follow](#)

10 min read · Dec 25, 2021

[Listen](#)[Share](#)[More](#)

source: quanta magazine

In many ML interview rounds, candidates are asked to demonstrate their coding skills w.r.t. some of the common and basic ML problems. These kind of rounds helps to identify both coding skills as well as ML skills required to be a top notch ML engineer in some of the top companies.

Let's look at some of the common ML coding problems asked in such interviews:

Problem 1:

Given the API `rand7()` that generates a uniform random integer in the range `[1, 7]`, write a function `rand10()` that generates a uniform random integer in the range `[1, 10]`. You can only call the API `rand7()`, and you shouldn't call any other API. Please do not use a language's built-in random API.

Solution:

`rand7()-1` generates random integer between 0 and 6.

To generate random integers between [1, 10], we need to generate 10 numbers which are equally probable or a range which is equally divisible into 10 parts.

Let $a = \text{rand7}() - 1$ and $b = \text{rand7}()$, then $7a + b$ will generate random integers between `[1, 49]` each with equal probability of $1/49$.

If the value of $7a + b$ is in the range `[1, 40]`, then we can use these range to equally divide into 10 parts each of size 4 as below:

`[1,4] [5,8] [9,12] [13,16] [17,20] [21,24] [25,28] [29,32] [33,36] [37,40]`

Thus if the value of $7a + b$ lies between `[1,4]`, we return 1, if the value lies between `[5,8]` return 2 and so on.

Thus return $\text{int}((7a+b-1)/4)+1$ if $1 \leq 7a+b \leq 40$.

```

1  class Solution(object):
2      def rand10(self):
3          while True:
4              a = (rand7()-1)
5              b = rand7()
6
7              c = 7*a+b
8
9              if c <= 40:
10                 return int((c-1)/4)+1

```

rand10.py hosted with ❤ by GitHub

[view raw](#)

Probability that the while loop returns on each call is $40/49$ because out of 49 possibilities for $7a+b$, only 40 are valid. Thus the expected number of times the while loop executes for each call is $49/40$.

For an arbitrary range $[a, b]$ using `rand7()`, first normalize the range to $[1, b-a+1]$, then find the smallest P such that $b-a+1 \leq 7^P$. Then find the largest K such that $K^*(b-a+1) \leq 7^P$.

*Let $Q = 7^{(P-1)} * (\text{rand7}() - 1) + 7^{(P-2)} * (\text{rand7}() - 1) + \dots + 7 * (\text{rand7}() - 1) + \text{rand7}()$*

Q generates integers in the range $[1, 7^P]$.

Consider only the range $[1, K^*(b-a+1)]$ from the values of Q . Among this range, if we divide it into $b-a+1$ equal parts each of size K , then for any integer in the 1st part i.e. $[1, K]$ we return 1, for 2nd part $[K+1, 2K]$ return 2 and so on.

Thus $\text{int}((Q-1)/K)+1$ will return random integer in the range $[1, b-a+1]$. To return in the range $[a, b]$, add $a-1$ to the output i.e. $a+\text{int}((Q-1)/K)$.

```

1 def rand7():
2     return random.randint(1, 7)
3
4 def generate_random(a, b):
5     while True:
6         p = int(math.log(b-a+1)/math.log(7))+1
7         q = 0
8
9         for r in range(p-1, 0, -1):
10            x = rand7()-1
11            q += x*(7**r)
12
13         q += rand7()
14         k = int((7**p)/(b-a+1))
15
16         if q <= k*(b-a+1):
17             return a+int((q-1)/k)

```

generate_random.py hosted with ❤ by GitHub

[view raw](#)

Expected number of times while loop executes is $7^P/K^*(b-a+1)$.

Problem 2:

You are given a **0-indexed** array of positive integers `w` where `w[i]` describes the **weight** of the `i`th index.

You need to implement the function `pickIndex()`, which **randomly** picks an index in the range `[0, w.length - 1]` (**inclusive**) and returns it. The **probability** of picking an index `i` is `w[i] / sum(w)`.

Solution:

Store the cumulative sums of weights till index i for each index i in the array. Generate a random integer M between [1, sum(w)].

Find the smallest index j, such that cum_sum[j] ≥ M. Return j.

Since the cumulative sums array is a sorted array, we can do **binary search** to find the index j s.t. $M \leq \text{cum_sum}[j]$.

```

1  class Solution(object):
2
3      def __init__(self, w):
4          self.cum_sum = [0]*len(w)
5          s = 0
6          for i in range(len(w)):
7              s += w[i]
8              self.cum_sum[i] = s
9
10
11     def pickIndex(self):
12         u = random.randint(1, self.cum_sum[-1])
13
14         left, right = 0, len(self.cum_sum)-1
15         p = -1
16
17         while left <= right:
18             m = int((left+right)/2)
19
20             if self.cum_sum[m] >= u:
21                 p = m
22                 right = m-1
23             else:
24                 left = m+1
25
26         return p

```

random_pick.py hosted with ❤ by GitHub

[view raw](#)

Time complexity is $O(\log N)$ for each call to `pickIndex()`.

This problem have very wide applications such as:

- Randomly selecting a server by load balancer based on the negative of current load as weights.
- Randomly sampling data from a discrete probability distribution.
- Generating artificial data for simulating traffic.

Problem 3:

You are given an integer `n` and an array of **unique** integers `blacklist`. Design an algorithm to pick a random integer in the range `[0, n - 1]` that is **not** in `blacklist`. Any integer that is in the mentioned range and not in `blacklist` should be **equally likely** to be returned.

Optimize your algorithm such that it minimizes the number of calls to the **built-in** random function of your language.

Solution:

Sort the integers in the `blacklist`.

Then generate intervals of the form `[blacklist[i-1]+1, blacklist[i]-1]`, i.e. **divide the entire range of integers such that it excludes the blacklisted integers**

For example, if the range of integers is `[1, 100]` and sorted blacklisted integers are `[10, 25, 40, 90]`, then the generated ranges are `[1,9] [11,24] [26,39], [41,89] [91,100]`.

Then based on the size of each range, pick a range with probability = size of range/sum of sizes of all ranges.

This can be done using the algorithm provided in problem 2.

After picking a range, randomly pick an integer within that range.

Probability of picking a number = Probability of picking a range * Probability of picking a number in that range.

If size of range is `S`, then probability = `S/sum(size of ranges)*1/S = 1/sum(size of ranges)`

Since `sum(size of ranges) = N-len(blacklisted)`, thus each number not in `blacklist` is picked uniformly.

```
1  class Solution(object):
2
3      def __init__(self, n, blacklist):
4          self.arr = []
5          self.cum_sum = []
6
7          if len(blacklist) == 0:
8              self.arr = [[0, n-1]]
9              self.cum_sum = [n]
10
11         else:
12             blacklist = sorted(blacklist)
13
14             for i in range(len(blacklist)):
15                 start = blacklist[i-1]+1 if i > 0 else 0
16                 end = blacklist[i]-1
17
18                 if end >= start:
19                     self.arr.append([start, end])
20
21                 if len(self.cum_sum) > 0:
22                     self.cum_sum.append(self.cum_sum[-1]+end-start+1)
23                 else:
24                     self.cum_sum.append(end-start+1)
25
26             start = blacklist[-1]+1
27             end = n-1
28
29             if end >= start:
30                 self.arr.append([start, end])
31
32             if len(self.cum_sum) > 0:
33                 self.cum_sum.append(self.cum_sum[-1]+end-start+1)
34             else:
35                 self.cum_sum.append(end-start+1)
36
37     def pick(self):
38         k = random.randint(1, self.cum_sum[-1])
39
40         left, right = 0, len(self.cum_sum)-1
41         p = -1
42         while left <= right:
43             m = int((left+right)/2)
44
45             if self.cum_sum[m] >= k:
```

```

46         p = m
47         right = m-1
48     else:
49         left = m+1
50
51     return random.randint(self.arr[p][0], self.arr[p][1])

```

random_pick_blacklisted.py hosted with ❤ by GitHub

[view raw](#)

Problem 4:

Given a method `biased_toss()` that generates 0 or 1. It generates 1 with some probability p (not necessarily 0.5). Write a method that generates 0 or 1 with probability 0.5.

Solution:

If we call `biased_toss()` method twice, then the possibilities are 00, 01, 10, 11. The probabilities of 00 is $(1-p)^2$, 01 is $(1-p)*p$, 10 is $p*(1-p)$ and 11 is p^2 .

Since the probabilities of 01 and 10 are equal i.e. $p(1-p)$, we can return 0 if we obtain 01 and 1 if we obtain 10.*

```

1 def get_unbiased_toss():
2     while True:
3         a = biased_toss()
4         b = biased_toss()
5
6         if a^b == 1:
7             return a
8

```

biased_toss.py hosted with ❤ by GitHub

[view raw](#)

Expected number of times while loop executes is $1/p*(1-p)$.

Problem 5:

Given an infinite stream of numbers. Write a method that would return K random samples from the stream without replacement.

Solution:

This can be solved using reservoir sampling.

For each number in the stream, we track its index.

If the size of the sample array is less than K, then append the number at the end of the array.

If the size of the sample array is more than equal to K, then generate a random integer between 0 and current index of the number in the stream. If the random integer is between 0 and K-1, then replace the index in the array corresponding to the random integer with the current number.

```

1  class ReservoirSampling:
2      def __init__(self, k):
3          self.out = []
4          self.max_length = k
5          self.index = 0
6
7      def add(self, num):
8          if self.index < self.max_length:
9              self.out.append(num)
10         else:
11             r = random.randint(0, self.index)
12             if r < self.max_length:
13                 self.out[r] = num
14
15         self.index += 1
16
17     def get_samples(self):
18         return self.out

```

reservoir_sampling.py hosted with ❤ by GitHub

[view raw](#)

We have to prove that the probability of each number in the stream, to be in the sample is equal to K/N where K is the length of samples array and N is the length of stream.

Assuming the length of stream is N .

For numbers with $\text{index} < K$, it is part of the sample if it is not replaced by any number with $\text{index} \geq K$.

Probability that number with index $J < K$ is part of sample = Probability that random integer generated at index K is not equal to J and Probability that random integer generated at index $K+1$ is not equal to J and ... Probability that random integer generated at index $N-1$ is not equal to J .

Probability that random integer generated at index K is not equal to $J = K/K+1$

Probability that random integer generated at index $K+1$ is not equal to $J = K+1/K+2$

Probability that random integer generated at index $N-1$ is not equal to $J = N-1/N$

Thus

Probability that number with index $J < K$ is part of sample = $K/(K+1)(K+1)/(K+2)*...*(N-1)/N=K/N.$*

For numbers with index $J \geq K$, it is part of the sample if the random integer X generated between 0 and J is $< K$ and the random integers generated at index $> J$ is not equal to X.

Probability that number with index $J \geq K$ is part of sample = Probability that random integer X generated between 0 and J is $< K$ and Probability that random integer generated at index $J+1$ is not equal to X and Probability that random integer generated at index $J+2$ is not equal to X and ... Probability that random integer generated at index $N-1$ is not equal to X.

Probability that random integer X generated between 0 and J is $< K = K/(J+1)$

Probability that random integer generated at index $J+1$ is not equal to X = $J+1/(J+2)$

Probability that random integer generated at index $J+2$ is not equal to X = $J+2/(J+3)$

Probability that random integer generated at index $N-1$ is not equal to X = $N-1/N$

Thus

Probability that number with index $J \geq K$ is part of sample = $K/(J+1)(J+1)/(J+2)*(J+2)/(J+3)*...*(N-1)/N=K/N.$*

Problem 6:

Randomly shuffle an array.

Solution:

Randomly sample an index j from i to Len(arr)-1. Then swap arr[j] with arr[i] and increment i by 1.

```

1 def shuffle(arr):
2     for i in range(len(arr)):
3         j = random.randint(i, len(arr)-1)
4         temp = arr[j]
5         arr[j] = arr[i]
6         arr[i] = temp
7

```

random_shuffle.py hosted with ❤ by GitHub

[view raw](#)

Problem 7:

Flip a fair coin repeatedly until you get two heads in a row (HH). What is the probability of getting HH in at-most N ($N > 1$) tosses ? Can you write a recursive function and implementation to calculate the probability ?

Solution:

Let $F(n, H)$ be the number of sequences of length n ending with HH and beginning with H and $F(n, T)$ be the number of sequences of length n ending with HH and beginning with T.

Thus $F(n+1, T) = F(n, H) + F(n, T)$ because if the sequence begins with T then the remaining sequence can begin with either H or T.

$F(n+1, H) = F(n, T)$ because if the sequence begins with H then the remaining sequence can begin only with T (HH is forbidden).

Let $G(n+1) = F(n+1, H) + F(n+1, T) = F(n, T) + F(n, H) + F(n, T) = F(n-1, H) + F(n-1, T) + F(n, H) + F(n, T)$

$$G(n+1) = G(n) + G(n-1)$$

$$G(1) = 0, G(2) = 1$$

$$\text{Probability } P(n) = (1/2^2)*G(2) + (1/2^3)*G(3) + (1/2^4)*G(4) + \dots + (1/2^n)*G(n)$$

```
1 def probability(n):
2     dp = [0]*n
3     dp[1], dp[2] = 0, 1
4
5     for i in range(3, n+1):
6         dp[i] = dp[i-1] + dp[i-2]
7
8     p = 0
9     power = 1.0/2.0**2
10
11    for i in range(2, n+1):
12        p += power*dp[i]
13        power /= 2.0
14
15    return p
```

hh.py hosted with ❤ by GitHub

[view raw](#)

Problem 8:

Compute nCr where $n \geq r$ and $1 \leq n \leq 1000$, $1 \leq r \leq n$.

Solution:

$$nCr = n-1Cr + n-1Cr-1$$

Thus we will use dynamic programming to compute the values:

```

1  def ncr(n, r):
2      dp = [[0]*(r+1) for i in range(n+1)]
3
4      for i in range(1, n+1):
5          for j in range(1, min(r+1, i+1)):
6              if i == j:
7                  dp[i][j] = 1
8              elif j == 1:
9                  dp[i][j] = i
10             else:
11                 dp[i][j] = dp[i-1][j] + dp[i-1][j-1]
12
13     return dp[n][r]

```

ncr.py hosted with ❤ by GitHub

[view raw](#)

Time and space complexity is $O(N^2)$.

Problem 9:

Implement Stratified sampling.

Given training data, class labels and the fraction of validation data per class ‘test_size’. Write a method to sample train data, test data, train labels and test labels such that each class has ‘test_size’ fraction of data/labels for validation and 1-test_size fraction of data/labels for training.

Solution:

One of the solution is the first group by based on class labels, then for each group call the sampling method to get the validation samples and then use the remaining as training samples.

Since the size of the data is fixed here, instead of using reservoir sampling as above (for streaming data), we will use another more cleaner method:

Generate a uniform random number between 0 and 1. If it is $\leq test_size$, then add the current data and label in validation set else add it in the training set.

```

1  def stratified_sampling(data, labels, test_size=0.2):
2      indices_per_label = {}
3
4      for i in range(len(labels)):
5          if labels[i] not in indices_per_label:
6              indices_per_label[labels[i]] = []
7              indices_per_label[labels[i]].append(i)
8
9      test_data, train_data, test_labels, train_labels = [], [], [], []
10
11     for label, indices in indices_per_label.items():
12         for i in indices:
13             k = random.uniform(0, 1)
14
15             if k <= test_size:
16                 test_data += [data[i]]
17                 test_labels += [labels[i]]
18             else:
19                 train_data += [data[i]]
20                 train_labels += [labels[i]]
21
22     return train_data, test_data, train_labels, test_labels

```

stratified1.py hosted with ❤ by GitHub

[view raw](#)

But do we really need to first group by the class labels ?

If we iterate over all data, then what is the probability of me selecting a point for adding in the validation set s.t. the probability per class is $test_size$.

Let $P(y=1)$ be the probability that a data /label is selected and $P(C_i)$ be the probability of class C_i or fraction of class C_i in the data.

Then $P(y=1) = P(y=1|C_1)*P(C_1) + P(y=1|C_2)*P(C_2) + \dots + P(y=1|C_m)*P(C_m)$ where there are ' m ' classes.

We know that $P(y=1|C1) = P(y=1|C2) = \dots = test_size$ i.e. the probability of selecting a data/label for validation given the class label is `test_size`.

Thus $P(y=1) = test_size * (P(C1) + P(C2) + \dots + P(Cm)) = test_size$, since $P(C1) + P(C2) + \dots + P(Cm) = 1$

Thus instead of grouping by class label, we can iterate over all data and select a data/label if $\text{uniform}(0,1) \leq test_size$.

```

1  def stratified_sampling(data, labels, test_size=0.2):
2      test_data, train_data, test_labels, train_labels = [], [], [], []
3
4      for i in range(len(data)):
5          k = random.uniform(0, 1)
6
7          if k <= test_size:
8              test_data += [data[i]]
9              test_labels += [labels[i]]
10         else:
11             train_data += [data[i]]
12             train_labels += [labels[i]]
13
14     return train_data, test_data, train_labels, test_labels

```

[startified2.py](#) hosted with ❤ by GitHub

[view raw](#)

It may not be intuitive at first since it might appear that we are sampling `test_size` fraction over all data instead of per class, but observe that the distribution of the classes is already handled when I am looping over all data. Thus a class with more samples will have more chance.

Problem 10:

You are given the true labels and the predicted probabilities from logistic regression model for N test examples. Approximately compute the AUC scores for ROC and PR curves.

Solution:

For ROC we need the TPR and FPR values.

$$\text{TPR} = \text{TP}/(\text{Number of all data with true label}=1)$$

$$\text{FPR} = \text{FP}/(\text{Number of all data with true label}=0)$$

For each probability threshold, find the TPR and FPR values.

Given a threshold T , $\text{TP} = \text{number of data with probability} \geq T \text{ having true label} = 1$ and $\text{FP} = \text{number of data with probability} \geq T \text{ having true label} = 0$. Also when T decreases both TPR and FPR increases and vice versa.

We can solve this efficiently by sorting the data on probability.

After sorting on probability, starting from the end, if at index i , true label = 1, then we increment TP by 1 else we increment FP by 1, because we are assuming the threshold to be probability[i].

To compute the AUC, we make use of the trapezoid formula: $0.5 * \text{width} * (\text{height}_1 + \text{height}_2)$ where width is the difference is FPR between 2 consecutive thresholds and height1 is the TPR at threshold $T(i)$ and height2 is TPR at threshold $T(i+1)$.

```

1  def auc_roc(true_labels, predictions):
2      data = zip(predictions, true_labels)
3      data = sorted(data, key=lambda k:k[0])
4
5      positives = sum([1 for x, y in data if y == 1])
6      negatives = sum([1 for x, y in data if y == 0])
7
8      tp, fp = 0, 0
9
10     prev_tpr, prev_fpr = 0, 0
11     area = 0
12
13     for i in range(len(data)-1, -1, -1):
14         if data[i][1] == 1:
15             tp += 1
16         if data[i][1] == 0:
17             fp += 1
18
19         tpr = tp/positives if positives > 0 else 0
20         fpr = fp/negatives if negatives > 0 else 0
21
22         if prev_fpr is not None and prev_tpr is not None:
23             area += 0.5*(fpr-prev_fpr)*(prev_tpr+tpr)
24
25         prev_fpr = fpr
26         prev_tpr = tpr
27
28     return area

```

[roc_auc.py](#) hosted with ❤ by GitHub

[view raw](#)

Similarly we can compute the AUC for PR curve.

Precision = TP/(Number of data predicted to be positive i.e. count of points with probability \geq threshold).

Recall = TP/(Number of all data with true label=1).

```

1  def auc_pr(true_labels, predictions):
2      data = zip(predictions, true_labels)
3      data = sorted(data, key=lambda k:k[0])
4
5      positives = sum([1 for x, y in data if y == 1])
6
7      tp = 0
8
9      prev_p, prev_r = 0, 0
10     area = 0
11
12     for i in range(len(data)-1, -1, -1):
13         if data[i][1] == 1:
14             tp += 1
15
16         p = tp/(len(data)-i) if len(data)-i > 0 else 0
17         r = tp/positives if positives > 0 else 0
18
19         if prev_p is not None and prev_r is not None:
20             area += 0.5*(r-prev_r)*(prev_p+p)
21
22         prev_r = r
23         prev_p = p
24
25     return area

```

auc_pr.py hosted with ❤ by GitHub

[view raw](#)

Problem 11:

Implement a basic binary logistic regression classifier. It should perform the following operations:

- Train model
- Predict on testing data

Assume the training and testing data is in the form of 2D matrix and the class labels is a 1-D array.

Solution:

Important methods to consider:

- Computing sigmoid scores
- Computing the logistic loss
- Gradient descent (batch or stochastic)

```
1 def sigmoid_probs(data, weights, bias):
2     probabilities = [0]*len(data)
3
4     for i in range(len(data)):
5         h = bias
6         for j in range(len(data[i])):
7             h += weights[i]*data[i][j]
8
9     probabilities[i] = 1.0/(1.0+math.exp(-h))
10
11    return probabilities
12
13 def loss(data, labels, weights, bias, reg_lambda):
14     probs = sigmoid_probs(data, weights, bias)
15     l = 0
16     for i in range(len(data)):
17         l += -labels[i]*math.log(probs[i])-(1-labels[i])*math.log(1-probs[i])
18
19     l += reg_lambda*(sum([w*w for w in weights]) + bias*bias)
20
21    return l
22
23 def gradient_descent(data, labels, weights, bias, num_epochs, learning_rate, reg_lambda):
24     for epoch in epochs:
25         probs = sigmoid_probs(data, weights, bias)
26
27         for j in range(len(data[0])):
28             s = 0
29             for i in range(len(data)):
30                 s += data[i][j]*(probs[i]-labels[i])
31
32             s += 2*reg_lambda*weights[j]
33             weights[j] -= learning_rate*s
34
35             s = 0
36             for i in range(len(data)):
37                 s += probs[i]-labels[i]
38
39             s += 2*reg_lambda*bias
40
41             bias -= learning_rate*s
42
43             curr_loss = loss(data, labels, weights, bias, reg_lambda)
44             print(epoch, loss)
45
```

```
46     return [weights, bias]
47
48 def init_weights(n):
49     return [random.random()]*n
50
51 class LogisticRegression:
52     def __init__(self, reg_lambda, epochs, batch_size, learning_rate):
53         self.weights = []
54         self.bias = 0
55         self.reg_lambda = reg_lambda
56         self.epochs = epochs
57         self.batch_size = batch_size
58         self.learning_rate = learning_rate
59
60     def train(self, data, labels):
61         n = len(labels)
62         #feature transformation and selection
63
64         self.weights = init_weights(len(data))
65         self.weights, self.bias = gradient_descent(data, labels,
66                                         self.weights, self.bias,
67                                         self.epochs,
68                                         self.learning_rate, self.reg_lambda)
69
70         #hyperparameter optimization with CV
71
72     def predict_proba(self, data):
73         output = [0]*len(data)
74
75         for i in range(len(data)):
76             h = self.bias
77             for j in range(len(data[i])):
78                 h += self.weights[j]*data[i][j]
79
80             output[i] = 1.0/(1.0+math.exp(-h))
81
82         return output
```

logistic_regression.py hosted with ❤ by GitHub

[view raw](#)

Problem 12:

Implement a simple decision tree classifier. It should perform the following operations:

- Train model
- Predict on testing data

Assume the training and testing data is in the form of 2D matrix and the class labels is a 1-D array.

Solution:

Important methods to consider:

- Computing information gain
- Splitting a node on a feature and feature value.

```
1  class Node:
2      def __init__(self):
3          self.data = []
4          self.labels = []
5          self.feature = None
6          self.feature_value = None
7
8      def get_entropy(indices, labels):
9          class_map = {}
10
11         for index in indices:
12             if labels[index] not in class_map:
13                 class_map[labels[index]] = 0
14                 class_map[labels[index]] += 1
15
16         ent = 0
17
18         for label, cnt in class_map.items():
19             p = float(cnt)/len(indices)
20             ent += -p*math.log(p)
21
22         return ent
23
24     def get_information_gain(left_indices, right_indices, labels):
25         left_ent = get_entropy(left_indices, labels)
26         right_ent = get_entropy(right_indices, labels)
27
28         total_entropy = get_entropy(left_indices+right_indices, labels)
29
30         return total_entropy-(left_ent+right_ent)
31
32     def construct_tree(data, labels, max_samples_leaf):
33         if len(data) == 0:
34             return None
35
36         elif len(data) <= max_samples_leaf or len(set(labels)) == 1:
37             node = Node()
38             node.data = data
39             node.labels = labels
40             return node
41
42         max_gain = -float()
43         best_feature = None
44         best_feature_value = None
45
```

```
46     for j in range(len(data[0])):
47         temp = []
48         for i in range(len(data)):
49             temp.append((data[i][j], i))
50
51         temp = sorted(temp, key=lambda k:k[0])
52
53         for k in range(len(temp)):
54             x, index = temp[k]
55
56             # unique values of features
57             if k == 0 or temp[k][0] != temp[k-1][0]:
58                 left_indices = [temp[h][1] for h in range(k+1)]
59                 right_indices = [temp[h][1] for h in range(k+1, len(temp), 1)]
60
61                 ig = get_information_gain(left_indices, right_indices, labels)
62
63                 if ig > max_gain:
64                     max_gain = ig
65                     best_feature = index
66                     best_feature_value = x
67
68             if best_feature is not None:
69                 node = Node()
70                 node.feature = best_feature
71                 node.feature_value = best_feature_value
72
73                 left_data, right_data = [], []
74                 left_labels, right_labels = [], []
75
76                 for i in range(len(data)):
77                     if data[i][best_feature] <= best_feature_value:
78                         left_data.append(data[i])
79                         left_labels.append(labels[i])
80                     else:
81                         right_data.append(data[i])
82                         right_labels.append(labels[i])
83
84                 node.left = construct_tree(left_data, left_labels, max_samples_leaf)
85                 node.right = construct_tree(right_data, right_labels, max_samples_leaf)
86
87             return node
88
89         return None
```

```
90
91 def predictions(root, data):
92     if root.data and len(root.data) > 0:
93         return max_class(root.labels)
94
95     f = root.feature
96     v = root.feature_value
97
98     if data[f] <= v:
99         return predictions(root.left, data)
100
101    return predictions(root.right, data)
102
103 class DecisionTree:
104     def __init__(self, max_samples_leaf=10):
105         self.max_samples_leaf = max_samples_leaf
106         self.root = None
107
108     def train(self, data, labels):
109         self.root = construct_tree(data, labels, self.max_samples_leaf)
110
111     def predict(self, data):
112         out = []
113         for d in data:
114             out.append(predictions(self.root, d))
115
116         return out
```

decision_tree.py hosted with ❤ by GitHub

[view raw](#)

Date Coding

Leetcode

Probability

LinkedIn

Machine Learning

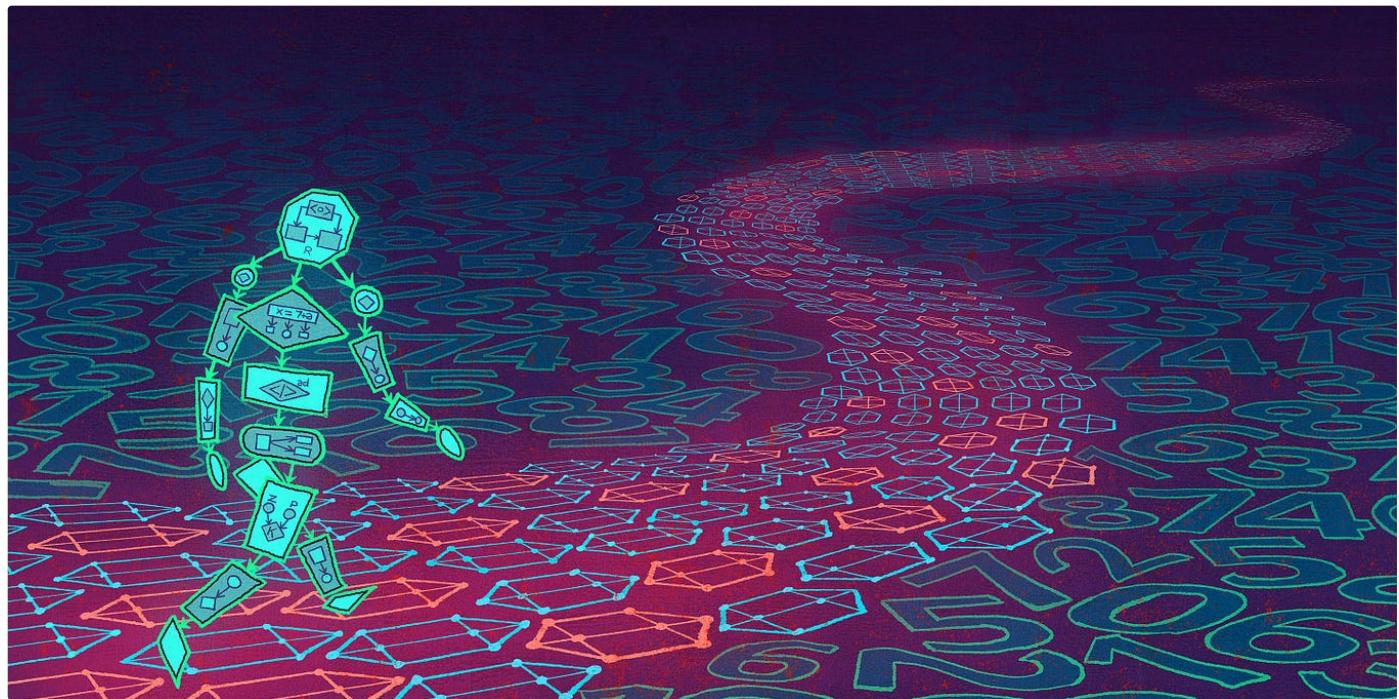
[Follow](#)

Written by Abhijit Mondal

1.6K Followers

Engineer

[More from Abhijit Mondal](#)



 Abhijit Mondal

Think twice before using asyncio in Python

I have been using multithreading and multiprocessing in Python for quite some time now. Very recently I came across this library called...

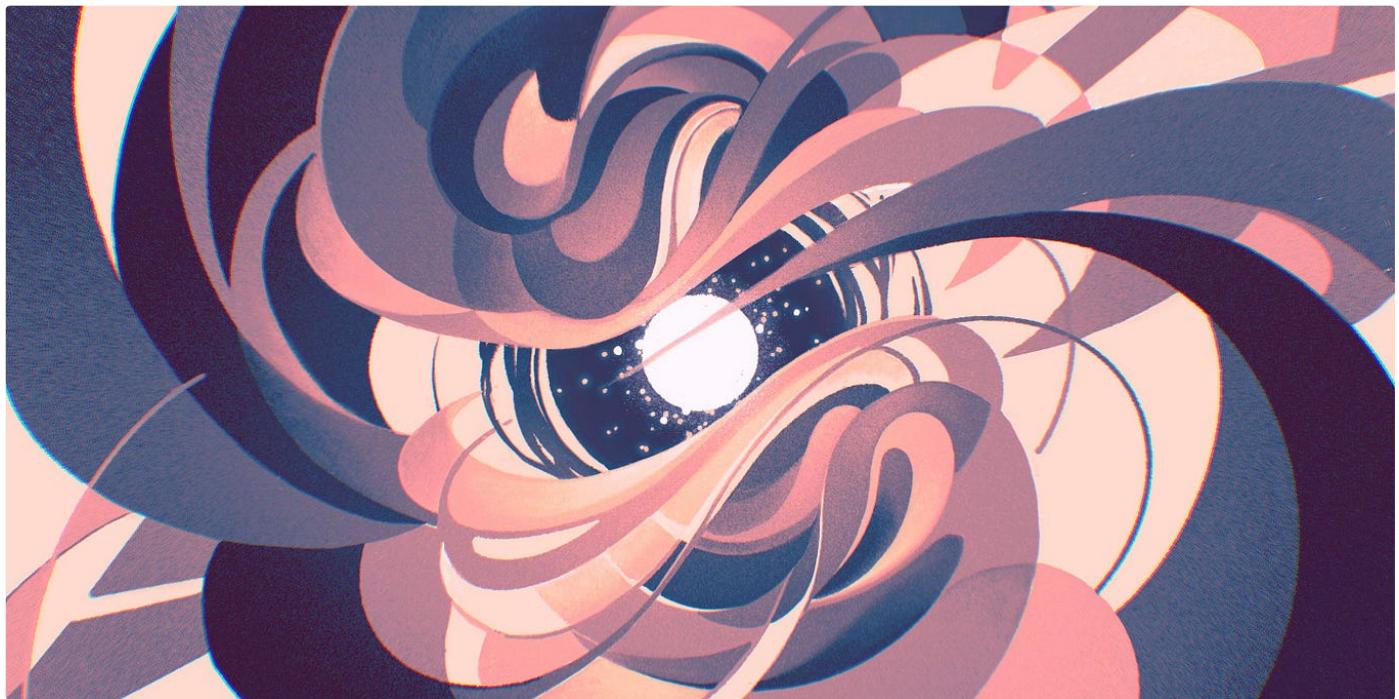
11 min read · Jan 19

 270

 7



...

 Abhijit Mondal

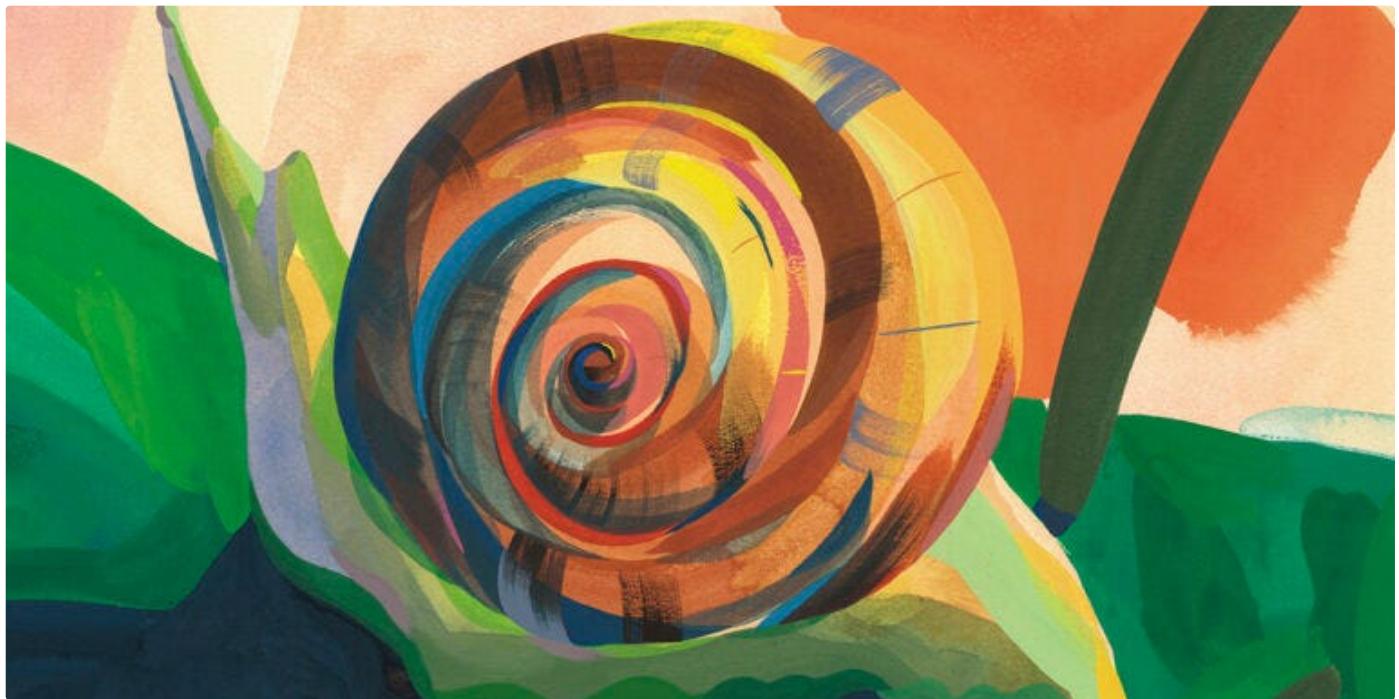
Algo Trading using Machine Learning on Zerodha Kite Connect

In this post I am going to explain how I (un)successfully tried my hands on Algo Trading on the Zerodha Kite Connect platform using Python...

12 min read · Sep 20, 2021

 44 2

...

 Abhijit Mondal

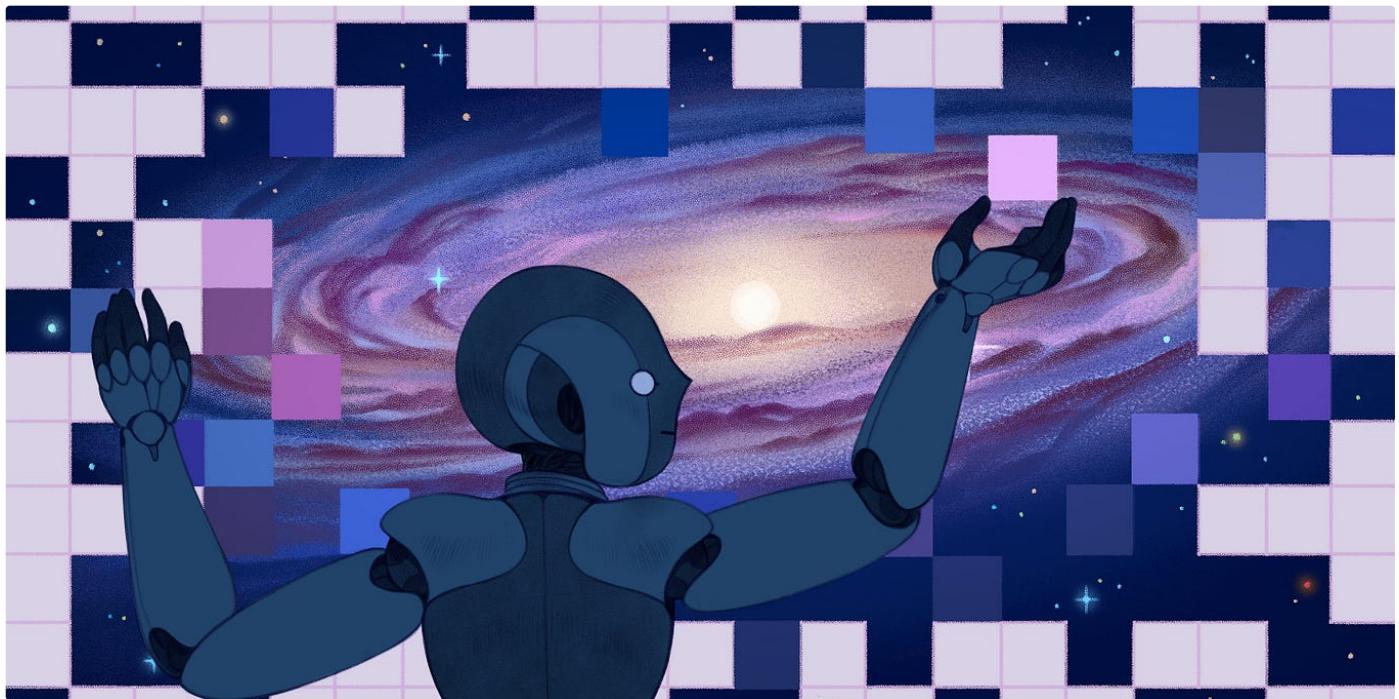
Deep dive into Chord for Distributed Hash Table

In a distributed hash table, a hash table with N keys is partitioned into M partitions or nodes. This is done in order to handle large...

15 min read · May 17

 30

...



 Abhijit Mondal

System Design—Top K Trending Hashtags

Design a system to compute the top K trending hashtags in the last 24 hours for Twitter/Instagram.

9 min read · Jan 29, 2022

 172

 8



...

See all from Abhijit Mondal

Recommended from Medium



Himanshu Joshi in Towards AI

How I cleared AWS Machine Learning Specialty With Three Weeks of Preparation (I will burst some...)

How I prepared for the test, my emotional journey during preparation, and my actual exam experience

9 min read · Jan 18

130

2



...

LeetCode 101: 20 Coding Patterns to the Rescue



Arslan Ahmad in Level Up Coding

Don't Just LeetCode; Follow the Coding Patterns Instead

What if you don't like to practice 100s of coding questions before the interview?

◆ · 5 min read · Sep 15, 2022

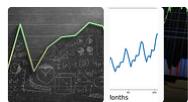
2.9K

10



...

Lists



Predictive Modeling w/ Python

18 stories · 80 saves



Practical Guides to Machine Learning

10 stories · 91 saves



Natural Language Processing

377 stories · 31 saves



The New Chatbots: ChatGPT, Bard, and Beyond

13 stories · 31 saves



 Mehul Gupta in Data Science in your pocket

Contextual Bandits in reinforcement learning explained with example and codes

starting deep reinforcement learning using tensorflow

5 min read · Feb 26

 121



...



Satnam Singh in Python in Plain English

Advanced Python Programming: Writing Efficient Code with Cython and Numba

Python is a popular programming language for its simplicity, ease of use, and flexibility. However, when it comes to performance, Python...

6 min read · Feb 28

👏 46



...



Youssef Hosni in Geek Culture

Crack the Machine Learning Coding Questions

Review these Algorithms and Their Code Before Your Next Machine Learning Interview

◆ · 17 min read · Jan 31



240



...



 Srushti Sonavane

Barclays Interview Experience

In the second week of October 2022, Barclays came to my college for the position of Software Intern (Summer Internship). The recruitment...

4 min read · Jun 7

 41

...

See more recommendations