

One-hot Vector (One-Hot Encoding)

1. Dữ liệu mẫu

```
crawl > phongvu > 9. Màn hình LCD Acer 21.5Inch R221QB.txt
1  Mô tả sản phẩmMàn hình LCD Acer 21.5 inch R221QB
2
3  THIẾT KẾ SANG TRỌNG, CỰC MỎNG
4  Màn hình LCD Acer 21.5 inch R221QB là màn hình doanh nghiệp có thiết kế hiện đại cực kỳ mỏng chỉ 6.6mm, và tràn viền, mang lại diện tích hi
5  Chân đế cách điệu thời thượng, mang lại không gian làm việc hiện đại.
6
7  Hình ảnh sắc nét chất lượng
8  Độ phân giải Full HD trên tấm nền IPS mang lại hình ảnh cực kỳ sắc nét và sống động và chân thực ở mọi góc nhìn
9  Độ tương phản động cực cao 100,000,000:1 càng làm cho độ hiển thị thêm sắc nét, và có chiều sâu hơn trên màn hình doanh nghiệp Acer 21.5 in
10
11  CHUYỂN ĐỘNG MƯỢT MÀ
12  Cùng trong series còn có sự góp mặt của màn hình doanh nghiệp Acer R221Q mang tần số quét 60Hz, và thời gian phản hồi 4ms cho hình ảnh chuy
13  Riêng màn hình doanh nghiệp R221QB được cải tiến lên tần số quét 75Hz, cùng 1ms phản hồi, tích hợp thêm công nghệ Freesync, loại bỏ hiện t
14
15  Bảo vệ thị lực suốt cả ngày
16  Màn hình LCD Acer 21.5 inch R221QB siêu mỏng được tích hợp công nghệ Bluelight Shield giúp loại bỏ 80% lượng ánh sáng xanh có hại cho mắt
17  Công nghệ khử nhấp Flickerless giúp bảo vệ thị giác người dùng xuyên suốt thời gian dài sử dụng, tránh tình trạng mỏi mắt, khó chịu
18  Công nghệ chống chói ComfyView giúp hình ảnh vẫn hiển thị tốt trong môi trường sáng, cho góc nhìn thoải mái suốt cả ngày.
19
20  TIẾT KIỆM ĐIỆN NĂNG TỐI ĐA
21  Màn hình LCD Acer 21.5 inch R221QB được trang bị công nghệ tiết kiệm năng lượng với tiêu chuẩn Energy Star 6.0 giúp tiết kiệm điện năng lên
22
23
```

2. Code

```
text_representation > OneHotEncoding.py > OneHotEncoding > run
1 class OneHotEncoding:
2     def __init__(self):
3         self.documents = ""
4
5     def run(self):
6         self.getData()
7         self.buildTheVocabulary()
8         print(self.getOnehotVector("Mô TIẾT").lower())
9
10    def getData(self):
11        file = open(
12            "D:/E23.1/NaturalLanguageProcessing/Ex/crawl/phongvu/9. Màn hình LCD Acer 21.5Inch R221QB.txt", "r", encoding="utf8")
13        self.documents = file.read()
14        self.processed_docs = self.documents.replace(
15            ".", " ").replace(",", " ").replace("-", " ").lower()
16
17    def buildTheVocabulary(self):
18        self.vocab = {}
19        self.count = 0
20        for word in self.processed_docs.split():
21            if word not in self.vocab:
22                self.count = self.count + 1
23                self.vocab[word] = self.count
24
25    # Get one hot representation for any string based on this vocabulary.
26    # If the word exists in the vocabulary, its representation is returned.
27    # If not, a list of zeroes is returned for that word.
28    def getOnehotVector(self, somestring):
29        onehot_encoded = []
30        for word in somestring.split():
31            temp = [0]*len(self.vocab)
32            if word in self.vocab:
33                # -1 is to take care of the fact indexing in array starts from 0 and not 1
34                temp[self.vocab[word]-1] = 1
35            onehot_encoded.append(temp)
36        return onehot_encoded
37
```

3. Kết quả

[illegible]

1. Dữ liệu mẫu: Như bài 1
2. Code

```

1  from sklearn.feature_extraction.text import CountVectorizer
2
3
4  class BagOfWords:
5      def __init__(self):
6          self.documents = ""
7
8      def run(self):
9          self.getData()
10         self.buildTheVocabulary()
11         print(self.getBagOfWords(20))
12
13     def getData(self):
14         file = open(
15             "D:/E23.1/NaturallanguageProcessing/Ex/crawl/phongvu/9. Màn hình LCD Acer 21.5Inch R221QB.txt", "r", encoding="utf8")
16         self.documents = file.read()
17         self.processed_docs = self.documents.replace(
18             ".", " ").replace(",", " ").replace("-", " ").lower().split()
19
20     def buildTheVocabulary(self):
21         self.count_vect = CountVectorizer()
22         # Build a BOW representation for the corpus
23         self.bow_rep = self.count_vect.fit_transform(self.processed_docs)
24
25         # Look at the vocabulary mapping
26         print("Our vocabulary: ", self.count_vect.vocabulary_)
27
28     def getBagOfWords(self, position):
29         return self.bow_rep[position].toarray()
30
31
32 if __name__ == '__main__':
33     BagOfWords().run()
34

```

3. Kết quả

Bag of n-grams

1. Dữ liệu mẫu: Như bài 1
2. Code

3. Kết quả

TF-IDF

2. Code

3. Kết quả

[illegible]