

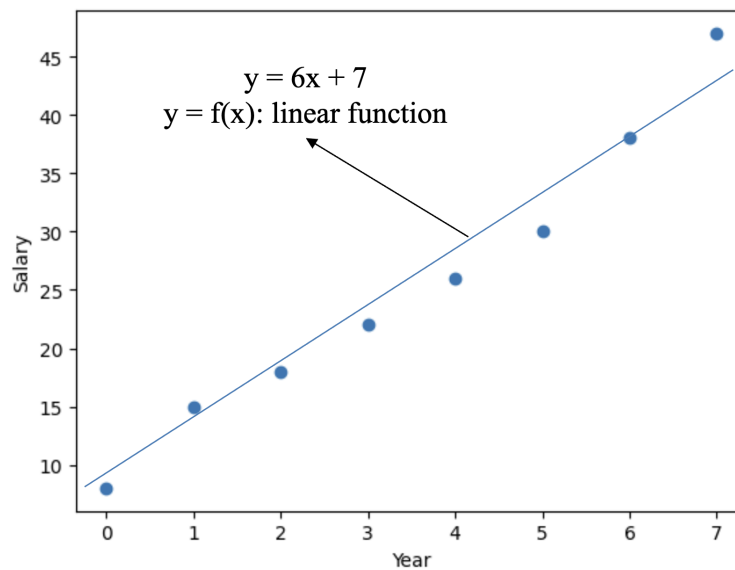
Project: Linear and non-linear Regressions

Ngày 9 tháng 7 năm 2023

Hồi quy tuyến tính - Linear Regression là mô hình phân tích mối quan hệ phụ thuộc của Y với một hay nhiều biến X sử dụng hàm tuyến tính (hàm đa thức bậc 1). Các tham số của mô hình (hay hàm số) được ước lượng từ dữ liệu huấn luyện.

Mô hình hồi quy tuyến tính cho bài toán dự đoán mức lương dựa vào số năm kinh nghiệm được mô tả như sau: $salary = year * w + b$. Phương trình tuyến tính xác định bởi 2 tham số w và b .

Visualization



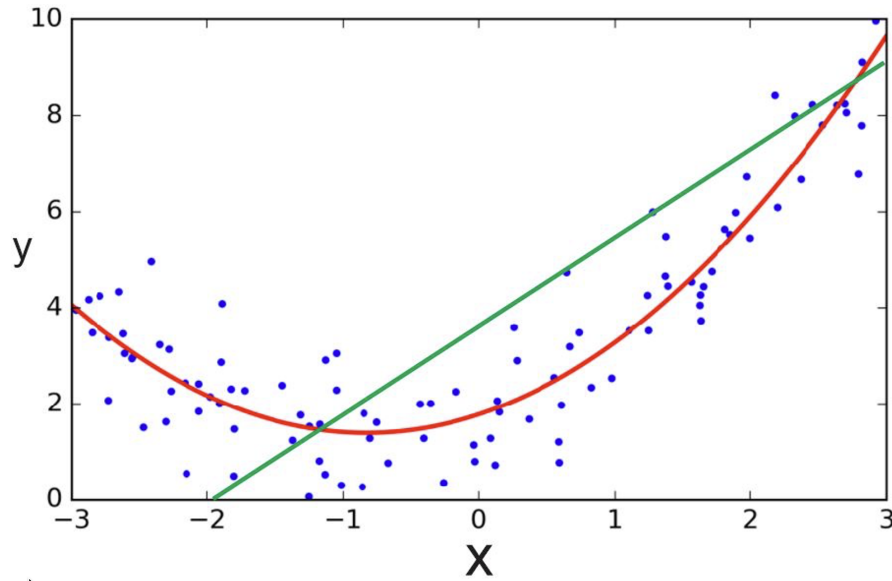
Hình 1: Mô tả dữ liệu dự đoán mức lương dựa vào số năm kinh nghiệm sử dụng mô hình hồi quy tuyến tính.

Dựa vào dữ liệu huấn luyện thông qua phương pháp học gradient descent để tìm bộ tham số tối ưu cho phương trình tuyến tính. Trong mỗi lần học, từ dữ liệu huấn luyện tính toán giá trị dự đoán, sau đó kết hợp với giá trị cần dự đoán trong thực tế để tính toán hàm mất mát (loss), rồi cập nhật lại tham số mô hình dựa vào tính đạo hàm. Quá trình huấn luyện cần chọn lựa các tham số ảnh hưởng đến kết quả của mô hình như tốc độ học (learning rate), giá trị khởi tạo của tham số và số vòng lặp thực hiện cập nhật trọng số).

Tuy nhiên, một trong những nhược điểm của mô hình hồi quy tuyến tính là giá trị dự đoán sự phụ thuộc tuyến tính của các giá trị đầu vào và tham số mô hình. Điều này dẫn tới khả năng tổng quát hoá của mô hình kém và không tối ưu cho dữ liệu trong thực tế. Vì vậy, mở rộng mối quan hệ phụ thuộc tuyến tính thông qua xây dựng hàm tuyến tính bằng các hàm phi tuyến sẽ giúp mô hình có tính tổng quát hơn. Các mô hình sử dụng các hàm phi tuyến với các giá trị đặc trưng đầu vào được gọi là hồi quy

phi tuyến (Nonlinear Regression). Một số phi tuyến điển hình như hàm đa thức (polynomial), hàm mũ hoặc hàm logarit,...

Ví dụ, hàm đa thức bậc 2 cho bài toán dự đoán lương dựa vào năm kinh nghiệm được mô tả như hình sau:



Hình 2: Mô tả dữ liệu dự đoán mức lương dựa vào năm kinh nghiệm sử dụng mô hình hồi quy đa thức bậc 2.

1. Linear regression

Để củng cố kiến thức về hồi quy tuyến tính, dựa vào file colab và đoạn code sau đây, để trả lời các câu hỏi.

```

1 class LinearRegression:
2     def __init__(self, X_data, y_target, learning_rate=0.01, num_epochs=10000):
3         self.X_data = X_data # Shape: (num_samples, num_features)
4         self.y_target = y_target # Shape: (num_samples, )
5         self.learning_rate = learning_rate
6         self.num_epochs = num_epochs
7         self.num_samples = self.X_data.shape[0]
8
9         # Initial Coefficients
10        self.theta = np.random.randn(self.X_data.shape[1])
11        self.losses = []
12
13    def compute_loss(self, y_pred, y_target):
14        loss = (y_pred-y_target)*(y_pred-y_target)
15        loss = np.mean(loss)
16        return loss
17
18    def predict(self, X_data):
19
20        *** Your code here ***
21
22        return y_pred
23
24    def fit(self):
25        for epoch in range(self.num_epochs):
26
27            *** Your code here ***
28
29            print(f'Epoch: {epoch} - Loss: {loss}')
30
31        return {
32            'loss': sum(self.losses)/len(self.losses),
33            'weight': self.theta
34        }

```

Câu hỏi 1 Dòng code nào sau đây giúp hoàn thiện hàm *predict()*:

- a) `y_pred = X_data.dot(self.theta)`
- a) `y_pred = X_data.dot(theta)`
- c) `y_pred = X_data.dot(self.theta)**2`
- d) `y_pred = X_data**2.dot(self.theta)`

Câu hỏi 2 Sắp xếp nào sau đây là đúng để hoàn thiện code trong hàm *fit()*:

Đoạn A:

```

1 # compute gradient
2 k = 2*(y_pred-self.y_target)
3 gradients = self.X_data.T.dot(k)/self.num_samples

```

Đoạn B:

```

1 # compute loss
2 loss = self.compute_loss(y_pred, self.y_target)
3 self.losses.append(loss)

```

Đoạn C:

```
1 # predict
2 y_pred = self.predict(self.X_data)
```

Đoạn D:

```
1 # update weight
2 self.theta = self.theta - self.learning_rate*gradients
```

Thứ tự đúng của các đoạn trên là:

- a) D - C - B - A
- b) C - D - B - A
- c) C - B - A - D
- d) C - B - D - A

Câu hỏi 3 Định nghĩa hàm tính hệ số xác định (Coefficient of determination) để đánh giá tính hiệu quả của mô hình hồi quy tuyến tính

Hệ số xác định (R^2) được xác định dựa vào công thức sau:

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

Chọn đáp án đúng trong các đáp án sau đây:

```
1 (a) def r2score(y_pred, y):
2     rss = np.sum((y_pred - y) ** 2)
3     tss = np.sum((y-y.mean()))
4     r2 = 1 - (rss / tss)
5     return r2
```

```
1 (b) def r2score(y_pred, y):
2     rss = np.sum((y_pred - y))
3     tss = np.sum((y-y.mean()) ** 2)
4     r2 = 1 - (rss / tss)
5     return r2
```

```
1 (c) def r2score(y_pred, y):
2     rss = np.sum((y_pred - y) ** 2)
3     tss = np.sum((y-y_pred.mean()) ** 2)
4     r2 = 1 - (rss / tss)
5     return r2
```

```
1 (d) def r2score(y_pred, y):
2     rss = np.sum((y_pred - y) ** 2)
3     tss = np.sum((y-y.mean()) ** 2)
4     r2 = 1 - (rss / tss)
5     return r2
```

Câu hỏi 4 Cho giá trị y_pred và y như sau dựa vào hàm $r2score$ được định nghĩa ở câu 3, tính các giá trị $r2score$ tương ứng sau đây:

```
1 y_pred = np.array([1, 2, 3, 4, 5])
2 y = np.array([1, 2, 3, 4, 5])
3 r2score(y_pred, y)
```

```
1 y_pred = np.array([1, 2, 3, 4, 5])
2 y = np.array([3, 5, 5, 2, 4])
3 r2score(y_pred, y)
```

- Giá trị $r2score(y_pred, y)$ lần lượt là:
- a) -2.2 và 1.0
 - b) 1.0 và -2.2
 - c) 1.0 và -4.4
 - d) -4.4 và 1.0

2. Polynomial Regression

Câu hỏi 5 Hàm số nào sau đây không phải là hàm phi tuyến:

- a) $y = 5x^2 + 3x$
- b) $y = 6x + 7$
- c) $y = x^3 + 2x^2 + 3$
- d) $y = \ln(x)$

Câu hỏi 6 Phát biểu nào sau đây là đúng nhất về polynomial regression

- a) Giá trị y phụ thuộc tuyến tính vào giá trị đầu vào X
- b) Giá trị y phụ thuộc phi tuyến vào giá trị đầu vào X và tuyến tính vào giá trị trọng số θ
- c) Giá trị y phụ thuộc tuyến tính vào giá trị trọng số θ
- d) Giá trị y phụ thuộc phi tuyến vào giá trị đầu vào X và phi tuyến vào giá trị trọng số θ

Polynomial Feature

Các đặc trưng phi tuyến cho hàm đa thức (polynomial feature) được tính toán như hình dưới. Với giá trị đầu vào $x = 2$, hàm đa thức bậc 3 sẽ tạo ra các feature tương ứng là $x^1 = 2, x^2 = 4, x^3 = 8$. Từ đó ta chuyển từ bài toán đơn biến chỉ có một giá trị x sang bài toán hồi quy tuyến tính đa biến.

$\psi(\varphi(i))$

Input	$\varphi(i)$	$\varphi(i)^2$
0	0	0
1	1	1
2	2	4
3	3	9
4	4	16
5	5	25

Hình 3: Ví dụ mô tả quá trình tạo các polynomial feature từ một dữ liệu đầu vào.

Examples

Input:

```
X = np.array([[1],
               [2],
               [3]])
```

```
degree = 2
```

Output:

```
X_poly = np.array([[1, 1],
                   [2, 4],
                   [3, 9]])
```

Hoàn thành đoạn code sau đây để tạo các polynomial feature từ một dữ liệu đầu vào:

```
1 def create_polynomial_features(X, degree=2):
2     """Creates the polynomial features
3     Args:
4         X: A array tensor for the data.
5         degree: A intege for the degree of
6             the generated polynomial function.
7     """
8     ** Your code here **
9     return X_new
10
11 X = np.array([[1], [2], [3]])
```

Câu hỏi 7 Chọn đáp án đúng trong các đáp án sau đây:

```
1 (a) X_new = X
2     for d in range(2, degree+1):
3         X_new = np.c_[X_new, np.power(X, d)]
```

```
1 (b) X_new = X
2     for d in range(2, degree+1):
3         X_new = np.c_[X_new, np.add(X, d)]
```

```
1 (c) X_new = X
2     for d in range(2, degree+1):
3         X_new = np.c_[X_new, np.square(X, d)]
```

```
1 (d) X_new = X
2     for d in range(1, degree+1):
3         X_new = np.c_[X_new, np.power(X, d)]
```

Xây dựng các đặc trưng phi tuyến cho mô hình đa biến thông qua việc xem xét các biến này phi tuyến riêng lẻ với nhau. Quá trình này được mô tả như hình sau:

$$(a+b)^2 \Rightarrow a^2 + b^2 + ab + a + b + 1$$

Multivariable

X		$\varphi(i)$		$\varphi(i)^2$		
0	2	0	2	0	4	0
1	1	1	1	1	1	1
2	2	2	2	4	4	4
3	1	3	1	9	1	3
4	2	4	2	16	4	8
5	1	5	1	25	1	5

Hình 4: Ví dụ mô tả quá trình tạo các polynomial feature từ nhiều dữ liệu đầu vào và xem xét các dữ liệu đầu vào phi tuyến độc lập với nhau.

Câu hỏi 8 Hàm nào trong các hàm sau đây phù hợp để xây dựng các đặc trưng phi tuyến được mô tả như hình trên:

Examples

Input:

```
X = np.array([[1, 2],
               [2, 3],
               [3, 4]])
```

degree = 2

Output:

```
X_poly = np.array([[1, 1, 2, 4],
                   [2, 4, 3, 9],
                   [3, 9, 4, 16]])
```

```
1 a) def create_polynomial_features(X, degree=2):
2     """Creates the polynomial features
3     Args:
4         X: A array for the data.
5         degree: A intege for the degree of
6             the generated polynomial function.
7     """
8     X_mem = []
9     for X_sub in X.T:
10        X_new = X_sub
11        for d in range(2, degree+1):
12            X_new = np.c_[X_new, np.power(X_sub, d)]
13        X_mem.extend(X_new.T)
14    return np.c_[X_mem].T
```

```
1 b) def create_polynomial_features(X, degree=2):
2     """Creates the polynomial features
3     Args:
```



```

4         X: A array for the data.
5         degree: A intege for the degree of
6         the generated polynomial function.
7         """
8         X_mem = []
9         for X_sub in X.T:
10             X_sub = X_sub.T
11             X_new = X_sub
12             for d in range(2, degree+1):
13                 X_new = np.c_[X_new, np.power(X_sub, d)]
14             X_mem.append(X_new.T)
15         return np.c_[X_mem].T

```

```

1 c) def create_polynomial_features(X, degree=2):
2     """Creates the polynomial features
3     Args:
4         X: A array for the data.
5         degree: A intege for the degree of
6         the generated polynomial function.
7     """
8     X_mem = []
9     for X_sub in X.T:
10         X_sub = X_sub.T
11         X_new = X_sub
12         for d in range(2, degree+1):
13             X_new = np.c_[X_sub, np.power(X_sub, d)]
14         X_mem.extend(X_new.T)
15     return np.c_[X_mem].T

```

```

1 d) def create_polynomial_features(X, degree=2):
2     """Creates the polynomial features
3     Args:
4         X: A array for the data.
5         degree: A intege for the degree of
6         the generated polynomial function.
7     """
8     X_mem = []
9     for X_sub in X.T:
10         X_sub = X_sub.T
11         X_new = X_sub
12         for d in range(2, degree+1):
13             X_new = np.c_[X_new, np.power(X_sub, d)]
14         X_mem.extend(X_new.T)
15     return np.c_[X_mem].T

```

Bộ dữ liệu dự đoán khối lượng của cá (file Fish.csv) được mô tả như hình sau:

	Species	Weight	VerticalLen	DiagonalLen	CrossLen	Height	Width
0	Bream	242.0	23.2	25.4	30.0	11.5200	4.0200
1	Bream	290.0	24.0	26.3	31.2	12.4800	4.3056
2	Bream	340.0	23.9	26.5	31.1	12.3778	4.6961
3	Bream	363.0	26.3	29.0	33.5	12.7300	4.4555
4	Bream	430.0	26.5	29.0	34.0	12.4440	5.1340
...
154	Smelt	12.2	11.5	12.2	13.4	2.0904	1.3936
155	Smelt	13.4	11.7	12.4	13.5	2.4300	1.2690
156	Smelt	12.2	12.1	13.0	13.8	2.2770	1.2558
157	Smelt	19.7	13.2	14.3	15.2	2.8728	2.0672
158	Smelt	19.9	13.8	15.0	16.2	2.9322	1.8792

159 rows x 7 columns

Hình 5: Dữ liệu bài toán dự đoán khối lượng của cá.

Trong đó trường thuộc tính *Species* thuộc kiểu dữ liệu text (các nhãn) - category. Vì vậy, các phương pháp mã hoá dữ liệu dạng nhãn thành số hoặc vector để xử lý trường thuộc tính này. Trong đó có hai phương pháp cơ bản điển hình là One Hot Encoding và Label Encoding.

Câu hỏi 9 Đoạn code nào sau đây để chuyển trường dữ liệu *Species* theo phương pháp One Hot Encoding

a) `encode_species = pd.get_dummies(df.Species)`

b) `encode_species = pd.get_dummies(~df.Species)`

c) `encode_species = pd.groupby(df.Species)`

d) `encode_species = pd.groupby(~df.Species)`

Câu hỏi 10 Đoạn code nào sau đây để chuyển trường dữ liệu *Species* theo phương pháp Label Encoding

a) `df["Species"] = df["Species"].astype('category')`
`label_encoding_species = df["Species"].cat.codes`

b) `label_encoding_species = df["Species"].cat.codes`

c) `df["Species"] = df["Species"].astype('numerical')`
`label_encoding_species = df["Species"].cat.codes`

d) `label_encoding_species = df["Species"].cat.numerical`

- Hết -