

Session 05

F1-Score & One-Hot Encoding

Machine Learning | Zahra Amini



Telegram: @zahraamini_ai & Instagram: @zahraamini_ai & LinkedIn: @zahraamini-ai

<https://zil.ink/zahraamini>

چ: اگر در مسائل classification فقط از accuracy برای ارزیابی مدل استفاده کنیم، آیا می‌توانیم تصمیم بگیریم

Cancer $\frac{\text{تعداد درست‌ها}}{\text{تعداد کل}}$ که عملکرد مدل خوب است؟

→ Accuracy: 95% → چ عملکرد مدل

مدل ما اگر همیشه "سالم" را پیش‌بینی کند با
دقت 95٪ دارد درست عمل می‌کند. ∴

100 نفر $\begin{cases} 95 \text{ سالم} \\ 5 \text{ بیمار} \end{cases}$

برای رفع این مشکل باید از Confusion Matrix استفاده کنیم.

y $\begin{cases} \text{True} \\ \text{False} \end{cases}$

y_Pred $\begin{cases} \text{Positive} \\ \text{Negative} \end{cases}$

زهرا امینی

@zahraamini_ai

True Positive \rightsquigarrow Predict Positive , True

True Negative \rightsquigarrow Predict Negative , True

False Positive \rightsquigarrow Predict Positive , False

False Negative \rightsquigarrow Predict Negative , False

TP TN

FP FN

| | y | | thr = 0.6 | Recall | Precision | Accuracy |
|------|---|-----|-----------|---------------|---------------|---------------|
| TN ← | 0 | 0.5 | 0 | | | |
| TP ← | 1 | 0.9 | 1 | | | |
| FP ← | 0 | 0.7 | 1 | $\frac{1}{2}$ | $\frac{2}{3}$ | $\frac{4}{7}$ |
| TP ← | 1 | 0.7 | 1 | | | |
| FN ← | 1 | 0.3 | 0 | | | |
| TN ← | 0 | 0.4 | 0 | | | |
| FN ← | 1 | 0.5 | 0 | | | |

| | | | |
|---------|---|--------|----|
| | | Actual | |
| Predict | 1 | TP | FP |
| | 0 | FN | TN |
| | | 1 | 0 |

$$\text{Recall} = \frac{TP}{TP + FN}$$

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

$$\text{Accuracy} = \frac{TP + TN}{\text{Total}}$$

$$\text{F1-Score} = \frac{2}{\frac{1}{\text{Recall}} + \frac{1}{\text{Precision}}}$$

بیشتر مدل های یادگیری ماشین با داده ها عددی کار می کنند.

؟ اگر بخواهیم بررسی داده های Categorical کار کنیم، چه طور باید این مشکل را حل کنیم ؟

خب ساده است باید داده هایمان را به عدد تبدیل کنیم، اما چه طور ؟

1. Integer Encoding

2. One-Hot Encoding

زهرا امینی
@zahraamini_ai

"Red" : 1

"Yellow" : 2

"Green" : 3

| Color |
|--------|
| Red |
| Red |
| Yellow |
| Green |
| Yellow |

Integer Encoding
→

| Color |
|-------|
| 1 |
| 1 |
| 2 |
| 3 |
| 2 |

ظاهر آنکه همه چیز به نظر خوب می‌رسد، اما واقعاً اینجوری نیست؟ **ordered relationship**

اگر ما به هر **Category** یک عدد بدهیم و مدل ما ارزش آن عدد را در حسابات دخیل کند در حالی که

عدد بزرگتر نشان دهنده‌ی برتری یک **Category** نسبت به دیگری نیست.

؟: برای رفع این مشکل چه کنیم؟ می‌توانیم از **One-Hot Encoding** استفاده کنیم.

| Color | | Red | Yellow | Green |
|--------|--------------------|-----|--------|-------|
| Red | | 1 | 0 | 0 |
| Red | One-Hot Encoding → | 1 | 0 | 0 |
| Yellow | | 0 | 1 | 0 |
| Green | | 0 | 0 | 1 |
| Yellow | | 0 | 1 | 0 |

۲: مشکلات One-Hot Encoding چیست؟

۱. افزایش تعداد ویژگی‌ها

۲. Dummy Variable Trap

| Red | Yellow | Green | | Red | Yellow |
|-----|--------|-------|---|-----|--------|
| 1 | 0 | 0 | → | 1 | 0 |
| 1 | 0 | 0 | | 1 | 0 |
| 0 | 1 | 0 | | 0 | 1 |
| 0 | 0 | 1 | | 0 | 0 |
| 0 | 1 | 0 | | 0 | 1 |