# به نام خالق یکتا



تمرین دوم مبانی یادگیریماشین

پاییز 1402

#### رگرسیون منطقی

1- فرض کنید به شما یک مجموعه آموزشی D شامل n نمونه به صورت  $(\mathbf{x}^{(n)},t^{(n)}),...,(\mathbf{x}^{(n)},t^{(n)})$  نمونه به صورت  $\mathbf{x}^{(i)}$  نمونه به صورت دودویی باشد، یعنی هر که هر ورودی  $\mathbf{x}^{(i)}$  یک بردار  $\mathbf{x}^{(i)}$  بعدی  $\mathbf{x}^{(i)}$  باشد. مدلی را در نظر بگیرید که برای هر نمونه  $\mathbf{p}(t=1|\mathbf{x}^{(i)},\mathbf{w},b)$  شکل یک تابع منطقی را به خود می گیرد:

$$p(t=1|\mathbf{x}^{(i)}, \mathbf{w}, b) = \sigma(\mathbf{w}^{T}\mathbf{x}^{(i)} + b) = \frac{1}{1 + \exp(-\sum_{i=1}^{d} \mathbf{w}_{d}\mathbf{x}_{d}^{(i)} - b)}$$

. تابع درستنمایی یا شباهت $^1$  به صورت  $p(t^{(1)},t^{(2)},...,t^{(N)}|\mathbf{x}^{(1)},\mathbf{x}^{(2)},...,\mathbf{x}^{(N)},\mathbf{w},b)$  تعریف شده است

علاوه بر این، فرض کنید برای منظم سازی $^2$ ، یک توزیع پیشین گاسی بـر روی وزنهـا  $\mathbf{w} = \{\mathbf{w}_1, \dots, \mathbf{w}_d\}$  قـرار داده شـده بـه طور یکه  $\mathbf{p}(\mathbf{w}) = N(\mathbf{w}|0, \alpha^{-1}\mathbf{I})$  است.

الف) تابع زیان  $^{3}$  را طوری تعریف کنید که منفی احتمال پسین خطا $^{4}$  نسبت به وزنها باشد. تابع زیان را تا جایی که میشود ساده بنویسید (تمام مشتقها را نشان دهید). برای این که این کار را کنید، فرض کنید دادهها  $^{5}$  هستند (یعنی، همه دادهها از هم مستقل هستند و از یک توزیع یکسان تولید شدهاند).

ب) نشان دهید که مشتق تابع زیان نسبت به وزنهای b,  $\mathbf{W_i}$  چه خواهد بود، یعنی  $\frac{\partial \ \mathrm{Loss}}{\partial \ \mathbf{w_i}}$  و حساب کنید.

پ) در نهایت، شبه کد برای نزول گرادیان را با استفاده از مشتقهای محاسبه شده در قسمت قبل بنویسید.

2- فرض کنید ما مدلهای رگرسیون منطقی زیر را برای دستهبندی دودویی با تابع  $\frac{1}{1+e^{-z}}$  در نظر داریم:

$$\mathsf{P}(y=1|\mathbf{x},w_1,w_2) = \mathsf{sigmoid}(w_1x_1+w_2x_2)$$
 عدل اول:  $\mathsf{P}(y=1|\mathbf{x},w_1,w_2) = \mathsf{sigmoid}(w_0+w_1x_1+w_2x_2)$  عدل دوم:  $\mathsf{O}$ 

ما نمونه های آموزشی زیر را داریم:

$$\mathbf{x}^{(1)} = [1,1]^{\mathrm{T}}$$
  $\mathbf{x}^{(2)} = [1,0]^{\mathrm{T}}$   $\mathbf{x}^{(3)} = [0,0]^{\mathrm{T}}$   $y^{(1)} = 1$   $y^{(2)} = -1$   $y^{(3)} = 1$ 

الف) اگر از مدل اول استفاده کنیم فرقی می کند که برچسب نمونه سوم چه عددی باشد؟

ب) اگر برچسب سومین نمونه آموزشی را به 1- تغییر دهیم، مقدار وزنهای  $w_2$  و  $w_1$  بعد از آموزش تغییری می کند؟

<sup>&</sup>lt;sup>1</sup> Likelihood Function

<sup>&</sup>lt;sup>2</sup> Regularization

<sup>&</sup>lt;sup>3</sup> Loss Function

<sup>&</sup>lt;sup>4</sup> Loss

<sup>&</sup>lt;sup>5</sup> Independent and Identically distributed

ج) اگر از مدل اول استفاده کنید چه تغییری برای مقادیر وزنها اتفاق میافتد؟

د) اگر از مدل دوم استفاده کنید چه تغییری برای مقادیر وزنها اتفاق میافتد؟

 $y^{(1)}, y^{(2)}, \dots, y^{(n)}$  با برچسبهای  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}$  با برچسبهای دوم را برای n داده n داده از بیشینه کردن احتمال لگاریتم شباهت یا درستنمایی n برچسبها آموزش دهیم.

$$\sum_i \log P(y^{(i)}|\mathbf{x}^{(i)}, \mathbf{w}) - \frac{\lambda}{2} ||\mathbf{w}||^2 = \sum_i \log \operatorname{sigmoid}(y^{(i)}\mathbf{w}^T\mathbf{x}^{(i)}) - \frac{\lambda}{2} ||\mathbf{w}||^2$$
 (1) برای  $\lambda$ های بزرگ (منظمسازی قوی<sup>7</sup>)، لگاریتم شباهت به صورت تابع خطی از  $\mathbf{w}$  عمل خواهد کرد:

log sigmoid
$$(y^{(i)}\mathbf{w}^{\mathrm{T}}\mathbf{x}^{(i)}) \approx \frac{1}{2} y^{(i)}\mathbf{w}^{\mathrm{T}}\mathbf{x}^{(i)}$$
 (2)

لگاریتم شباهت نهایی را با استفاده از این تقریب بیان کنید (با مدل 1) و از عبارت بیشینه شباهت نسبت به  $\mathbf{w}$  مشتق بگیرید و جواب را بر حسب  $\lambda$  و مجموعه آموزشی  $\{\mathbf{x}^{(i)}, y^{(i)}\}$  بیان کنید. هم چنین بگویید رفتار  $\mathbf{w}$  با افزایش  $\lambda$  چگونه تغییر می کنید (فرض کنید هر  $\mathbf{x}^{(i)}, \mathbf{x}^{(i)} = (\mathbf{x}_1^{(i)}, \mathbf{x}_2^{(i)})$  باشد و هر  $\mathbf{x}^{(i)}, \mathbf{y}^{(i)}$  باشد).

#### رگرسیون خطی

3- دو متغیر حقیقی x و y که x مشروط به y به شکل زیر تولید می شود را در نظر بگیرید:

$$\varepsilon \sim N(0, \sigma^2)$$

$$y = ax^2 + bx + c + \varepsilon$$

که x یک متغیر است که یک نویز گاسی با توزیع نرمال با میانگین a و واریانس a را نمایندگی می کند. یک رگرسیون خطی دارای و علی متغیر است که یک نویز گاسی با توزیع نرمال با میانگین a دارای رابط و می متغیر است. احتمال شرطی a دارای رابط و پارامترهای a و دارای رابط و و دارای رابط

است که به شکل زیر نوشته می شود:

$$P(y|x, a, b, c) = \frac{1}{\sigma \sqrt{2\pi}} exp(-\frac{1}{2\sigma^2} (y - ax^2 + bx + c)^2)$$

الف) فرض کنید n زوج آموزشی  $(x^{(i)}, y^{(i)})$  داریم و واریانس  $\sigma$  مجهول است. با استفاده از بیشینه شباهت هـر یـک از پـارامتر های این مدل را تخمین بزنید.

ب) اگر مجموعه داده ما چهار زوج آموزشی به شکل  $D=\{(5,8),(1,4),(3,-2),(0,-8)\}$  داشته باشد، با استفاده از بیشینه شباهت هر یک از پارامترهای این مدل را محاسبه کنید.

\_

<sup>&</sup>lt;sup>6</sup> Log-Likelihood

Strong regularization

باشد و مدل زیر را داشته باشیم:  $\mathbf{x}_i\in\mathbb{R}^n$  و  $\mathbf{x}_i\in\mathbb{R}^m$  ،  $\mathbf{D}=\{(\mathbf{x}_i,y_i)\}_{i=1}^n$  -4 -4  $y_i\sim N(\mathbf{w}^{\mathrm{T}}\mathbf{x}_i,\sigma^2)$ 

شباهت به صورت زیر تعریف خواهد شد:

$$P(y|\mathbf{x},\mathbf{w}) = N(\mathbf{x}\mathbf{w},\sigma^2\mathbf{I}) = \prod_{i=1}^{|D|} \frac{1}{\sigma\sqrt{2\pi}} \exp(-\frac{(y_i - \mathbf{w}^T\mathbf{x}_i)^2}{2\sigma^2})$$

توزیع پیشین برای وزنها را به صورت زیر در نظر بگیرید:

$$P(\mathbf{w}) = N(0, \sigma_0^2 \mathbf{I}) = \prod_{j=1}^{|\mathbf{w}|} \frac{1}{\sigma \sqrt{2\pi}} \exp(-\frac{(\mathbf{w}_j)^2}{2\sigma_0^2})$$

تخمین MAP را برای وزنها به دست آورید.

5- درست یا غلط بودن عبارات زیر را با دلیل مشخص کنید.

الف) بیشینه شباهت پارامتر a در مدل میتواند با استفاده از رگرسیون خطی یاد گرفته شود اگر مدل به این شکل باشد:  $\varepsilon_i \sim N(0, \sigma^2)$  و  $y_i = \log(x_1^{a_1} e^{a_2}) + \varepsilon_i$ 

ب) بیشینه شباهت پارامتر a در مدل میتواند با استفاده از رگرسیون خطی یاد گرفته شود اگر مدل به این شکل باشد:  $\varepsilon_i \sim N(0, \sigma^2)$  و  $y_i = x_1^{a_1} e^{a_2} + \varepsilon_i$ 

6- میخواهیم پارامترهای یک رگرسیون را برای یک مجموعه داده که میدانیم توسط یک تابع چند جملهای تولید شده است اما درجه این چند جملهای را نمیدانیم، یاد بگیریم. فرض کنید داده واقعا توسط یک چند جملهای درجه با یک نـویز گاسی جمع  $y = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + \varepsilon$  بین بــه صــورت  $y = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + \varepsilon$  بین بین بین بین بین بین بین برای آموزش و تعداد 100 زوج  $\{x,y\}$  برای آموزش و تعداد 100 زوج  $\{x,y\}$  دیگر برای آزمایش موجود است.

B چون درجه چندجملهای را نمی دانیم، پارامترهای دو مدل را برای داده یاد می گیریم، مدل A یک چندجملهای درجه 2 و مـدل 2 و مـدل یک چندجملهای درجه 3 است. کدام یک از این دو مدل بهتر بر داده آزمایشی منطبق می شوند؟ چرا؟

### بخش پیاده سازی

## سوال اول)

#### 1. مقدمه

در این سوال، هدف پیادهسازی رگرسیون لاجستیک برای پیش بینی نجات یافتن یا نیافتن مسافران کشتی تایتانیک است. در این تمرین، با پلتفرم کگل نیز آنشا می شوید. کگل به عنوان منبع عظیمی از دادهها و پروژههای داده معتبر شناخته می شود و جامعه بزرگی از داده کاوانها و تحلیل گران داده را به خود جذب کرده است. این پلتفرم به افراد امکان می دهد تا مهارتهای خود را در حوزه داده و یادگیری ماشین بهبود بخشیده و در پروژههای واقعی مشارکت کنند و همچنبن مسابقات بزرگی در حوزه یادگیری ماشین برگزار میکند.

### 2. معرفي مجموعهداده

این مجموعهداده، اطلاعات مربوط به مسافران کشتی تایتانیک را دارد. نام و توضیحات هر یک از ستونهای این مجموعهداده به شرح زیر است:

- Survival: زنده ماندن یا نماندن مسافر که با 0 یا 1 نمایش داده میشود
  - pclass: کلاس بلیط مسافر که میتواند کلاس 1، 2 و یا 3 باشد
    - sex: جنسیت
    - Age: سن مسافر
    - sibsp: تعداد همسران یا خواهران و برادران همراه مسافر
      - parch: تعداد فرزندان یا والدین همراه مسافر
        - ticket: شماره بليط مسافر
        - cabin: شماره کابین مسافر
- embarked و یا S بندر مربوطه که می تواند C = Cherbourg, Q = Queenstown, S = Southampton کد بندر مربوطه که می تواند C = Cherbourg, Q = Queenstown, S = Southampton

## 3. ثبت پاسخ

برای ثبت پاسخ مساله از طریق لینک زیر پاسخ خود را در سایت کگل آپلود نمایید و اسکرینشات نتیجهی آنرا به مستندات تحویلی خود اضافه نمایید.

https://www.kaggle.com/competitions/titanic/data

# سوال دوم)

#### 1. مقدمه

در این سوال، هدف پیادهسازی الگوریتم k نزدیکترین همسایه برای پیشبینی کلاس گلها است. این پیشبینی به کمک اطلاعات مجموعهدادهی Iris صورت خواهد گرفت.

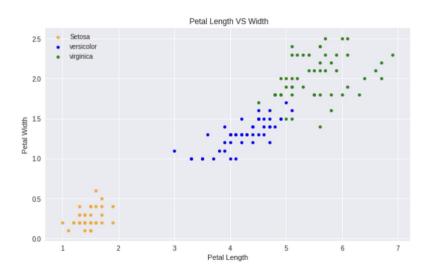
#### 2. معرفی مجموعه داده

مجموعه داده ی انواع گل iris، اطلاعات مربوط به کاسبرگ و گلبرگ گل و نوع آنها را دارد. نام و توضیحات مربوط به هر یک از ستونهای این مجموعه داده به شرح زیر است:

- ID: آی دی منسوب به هر گل
- SepalLengthCm: طول کاسبرگ گل به سانتی متر
- SepalWidthCm: عرض کاسبرگ گل به سانتی متر
- PetalLengthCm: طول گلبرگ گل به سانتی متر
- PetalWidthCm: عرض گلبرگ گل به سانتی متر
  - Species: نوع گل

# 3. مصورسازي

برای این سوال لازم است که نموداری مانند نمودار زیر رسم کنید و در آن نتایج پیشبینی مدل خود را نشان دهید.



# 4. ارزیابی

معیار accuracy را برای راه حل خود گزارش کنید.

### 5. آپلود

در انتها نتایج مدل خود را در فایلی به فرمت .CSV ذخیره کنید. این فایل باید به تعداد سطرهای داده ی تست پیشبینی داشته باشد که در یک ستون با نام Species ذخیره شده است. فایل نتایج را به علاوه ی کد در پوشه ی مربوط به سوال در ریپازیتوری گیتهاب قرار دهید.

## سوال سوم)

#### 1. مقدمه

در این تمرین، هدف پیادهسازی رگرسیون لاجستیک برای پیش بینی اسپم بودن یا نبودن ایمیل ها است بر روی دیتاست Scikit learn است. در این سوال شما اجازه دارید تا از کتابخانهی Scikit learn استفاده کنید.

### 2. معرفي مجموعهداده

دادههای مجموعهداده شامل اطلاعاتی عمدتا آماری از کلماتی است که در متن ایمیل وجود دارند. همچنین اسپم بودن یا نبودن ایمیل نیز در این ویژگیها نمایش داده شده است. 58 ویژگی در این مجموعهداده وجود دارد که توضیحات مربوط به آنها به شرح زیر است:

- 48 ستون با مقادیر پیوسته با نامهایی به فرمت word\_freq\_WORD که درصد کلاماتی که با لغت WORD مطابقت دارند را نمایش می دهد.
- 6 ستون با مقادیر پیوسته با نامهایی به فرمت char\_freq\_CHAR که درصد کرکترهایی که با کرکتر CHAR تطابق دارند را نمایش می دهد.
  - 1 ستون با نام capital\_run\_length\_average که میانگین طول رشتههای متشکل از حروف بزرگ (Capital) را نمایش می دهد.
- 1 ستون به نام capital\_run\_length\_longest که طول بلندترین رشتهی متشکل از حروف بزرگ را نشان می دهد.
- 1 ستون به نام capital\_run\_length\_total که تعداد کل حروف بزرگ به کاررفته در متن ایمیل را نمایش می دهد.
  - 1 ستون به نام Class که در آن اسپم نبودن یا بودن ایمیل، به ترتیب با مقدار 0 یا 1 نمایش داده شده است.

این مجموعه داده یک مجموعهدادهی نامتعادل (imbalanced) است (تقریبا 39.4 درصد) که برای مدیریت آن باید از سه تکنیک undersampling، و تابع loss وزن دار استفاده نمایید.

همچنین در سند نهایی شما باید سه معیار f1, recall precission وaccuracy را برای هر سه تکنیک گزارش کنید.

# 3. مصورسازي

برای این تمرین لازم است تا ماتریس درهمریختگی (Confusion Matrix) مربوطه را رسم کنید و در سند نهایی تمرین قرار دهید.

### 4. آيلود

در انتها نتایج مدل خود را در فایلی به فرمت .CSV ذخیره کنید. این فایل باید به تعداد سطرهای دادهی تست پیشبینی داشته باشد که در یک ستون با نام Class ذخیره شده است. فایل نتایج را به علاوهی کد در پوشهی مربوط به سوال در ریپازیتوری گیتهاب قرار دهید.

### آنچه باید تحویل دهید:

لازم است که موارد زیر را به عنوان موارد مورد تحویل هر سوال در پوشهی مربوط به آن سوال داخل ریپازیتوری گیتهای قرار دهید:

- فایل ژوپیتر نوت بوک به فرمت .ipynb که شامل کد پاسخ و توضیحات مربوط به آن است.
- فایل پیشبینیهای خروجی مدل که به فرمت Q#.CSV خواهد بود و # شمارهی سوال میباشد.
- برای سوال اول، اسکرین شات نتیجه بارگذاری پیشبینی مدل در سایت کگل به شکلی که نام اکانت مشخص باشد.
  - تحویل همه موارد فقط و فقط از گیت هاب صورت میگیرد. لینک تمرین:

•

https://classroom.github.com/a/pBm5fJ4H

• مهلت تمرین تا سه شنبه 30 آبان، ساعت 3 بامداد می باشد.

موفق باشيد