



**دانشگاه صنعتی امیر کبیر**  
(پلی تکنیک تهران)

دانشکده ریاضی و علوم کامپیوتر

تمرین اول درس یادگیری ماشین

عنوان: رگرسیون خطی

موعد تحویل:

## ۱. Linear Regression (رگرسیون خطی)

در این تمرین شما رگرسیون خطی چند متغیره<sup>۱</sup> را پیاده سازی خواهید کرد. این تمرین طولانی به نظر می رسد، اما در عمل شما تنها حدود ۱۰ خط کد را پیاده سازی می کنید. این تمرین برای آشنایی شما با مکانیزم کتابخانه scikit-learn و پیاده سازی الگوریتم های یادگیری ماشین در پایتون طراحی شده است.

### Files and API (فایل ها و توابع)

در این تمرین یک سری فایل و توابع پایتون در اختیار شما قرار دارد:

`test_linreg_univariate.py`: اسکریپتی برای تست کردن رگرسیون خطی تک متغیره<sup>۲</sup>

- `plotData1D`: رسم نمودار نقطه‌ای (پراکندگی)<sup>۳</sup> دیتای یک بعدی
- `plotRegLine1d`: بوسیله تابع `plotData1D` نمودار نقطه‌ای دیتا را رسم می‌کند و بوسیله پارامترهای یادگرفته شده توسط الگوریتم، خط رگرشن را در همان نمودار رسم می‌کند.
- `visualizeObjective`: رسم نمودار `surface` و `contour` تابع هزینه<sup>۴</sup> (شما این تابع را تغییر نمی‌دهید)

`test_linreg_multivariate.py`: اسکریپتی برای تست کردن رگرسیون خطی چند متغیره

`linreg.py`: اسکریپتی شامل کدهای رگرسیون خطی

`LinearRegression`: کلاس رگرسیون خطی چند متغیره

- `__init__`: کانستراکتور کلاس
- `fit`: تابعی برای یادگیری<sup>۵</sup> مدل رگرسیون خطی چند متغیره
- `predict`: تابعی برای پیش بینی ورودی جدید توسط مدل یاد گرفته شده
- `computeCost`: محاسبه مقدار تابع هزینه
- `gradientDescent`: بهینه سازی پارامترهای مدل توسط الگوریتم گرادیان کاهشی<sup>۶</sup>

### Data Sets (مجموعه داده‌ها واقع در فولدر data)

- `univariateData.dat`: مجموعه داده‌ها برای مسئله رگرسیون تک متغیره
- `multivariateData.dat`: مجموعه داده‌ها برای مسئله رگرسیون چند متغیره

### Visualizing the Data (بصرسازی داده‌ها)

بصرسازی داده‌ها بینش ارزشمندی از مسئله به ما ارائه می دهد، اما اغلب به عنوان بخشی از فرآیند یادگیری ماشین نادیده گرفته می شود. ما با رسم مجموعه داده‌های تک متغیره با استفاده از نمودار پراکندگی<sup>۲</sup> بعدی شروع خواهیم کرد. با این حال ما معمولاً با مجموعه داده‌های چند بعدی مواجه هستیم. هنگامی که از دو بعد فراتر می رویم، تجسم

<sup>۱</sup> multivariate linear regression

<sup>۲</sup> univariate linear regression

<sup>۳</sup> scatter plot

<sup>۴</sup> cost function (objective function)

<sup>۵</sup> train

<sup>۶</sup> gradient descent

بسیار دشوارتر می شود. در چنین مواردی یا باید هر بعد را به طور جداگانه رسم کنیم، یا از تکنیک های کاهش ابعاد (مانند PCA) برای کاهش تعداد ویژگی ها<sup>۷</sup> استفاده کنیم. بعداً در کلاس در مورد چنین تکنیک هایی صحبت خواهیم کرد.

با اجرای دستورات زیر در مفسر پایتون از داخل دایرکتوری hw1 می توانید داده های تک متغیره را در متغیرهای  $X$  و  $y$  بصورت ماتریس بارگذاری کنید:

```
import numpy as np
filePath = "data/univariateData.dat"
file = open(filePath, 'r')
allData = np.loadtxt(file, delimiter=',')
X = np.matrix(allData[:, :-1])
y = np.matrix((allData[:, :-1])).T
# get the number of instances (n) and number of features (d)
n, d = X.shape
```

سپس بوسیله ی تابع `plotData1D` نمودار پراکندگی را رسم کنید، خروجی شما باید مانند Figure ۱ باشد:

```
from test_linreg_univariate import plotData1D
plotData1D(X, y)
```

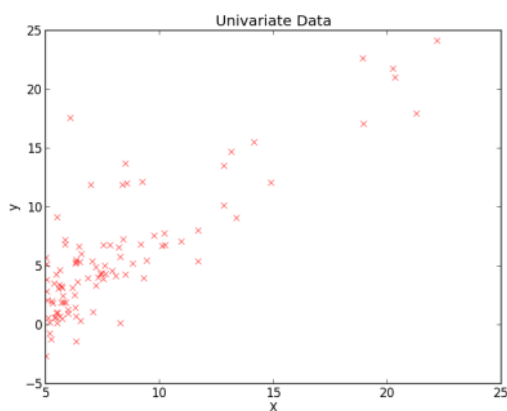


Figure 1: Scatter plot of the 1D data

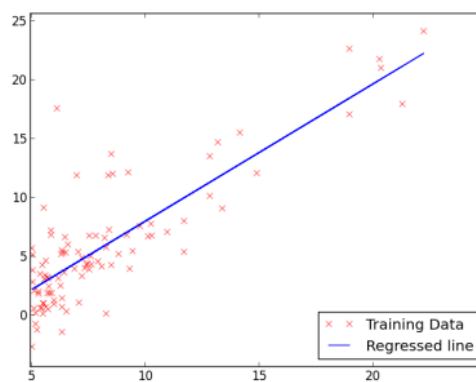


Figure 2: Regressed line of the 1D data

---

<sup>۷</sup> features

## Implementation (پیاده سازی)

با تکمیل کلاس LinearRegression، رگرسیون خطی چند متغیره را از طریق گرادیان کاهشی اجرا کنید. اسم توابع (APIها) را ابداء تغییر ندهید. قسمت‌هایی از کد که باید تغییر دهید با کامنت "TODO" مشخص شده اند.

الگوریتم رگرسیون خطی بعد از یادگیری پارامترهای مدل به کمک دیتا، آن‌ها را بر روی بردار  $\theta^a$  ذخیره می‌کند<sup>۹</sup>. در این تمرین از گرادیان کاهشی برای یافتن جواب بهینه استفاده می‌کنیم. دقت کنید که تابع هزینه رگرسیون خطی  $L_2$ ، محدب است، بنابراین الگوریتم گرادیان کاهشی مینیمم مطلق را پیدا می‌کند.

$$(۱) \quad \hat{\theta} = \min J(\theta)$$

$$(۲) \quad J(\theta) = \frac{1}{2n} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$(۳) \quad h_{\theta}(x) = \theta^T x$$

تابع  $h_{\theta}$  در رگرسیون خطی تک متغیره (ماتریس  $x$  تنها یک ستون دارد) به این شکل  $h_{\theta}(x) = \theta^T x = \theta_0 + \theta_1 x$  است، که  $\theta_0$  همان عرض از مبدا<sup>۱۰</sup> است. برای هندل کردن عرض از مبدا در قالب معادله‌ی (۳) می‌توانیم یک ویژگی جدید به تمام سطرهای داده با مقدار ۱ اضافه کنیم. در واقع  $\theta_0$  را به عنوان ضریب  $x_0$  در نظر می‌گیریم. برای اضافه کردن عرض از مبدا به کل مجموعه داده‌ها می‌توانیم ستونی از یک‌ها را به ماتریس  $X$  اضافه کنیم:

```
X = np.c_[np.ones((n, 1)), X]
```

الگوریتم گرادیان کاهشی برای یافتن مقدار مینیمم تابع  $J(\theta)$  فضای  $\theta$ های ممکن را جست و جو می‌کند. حلقه for اولیه گرادیان کاهشی برای شما پیاده سازی شده است. شما فقط باید معادله را به روز رسانی کنید. در هر مرحله<sup>۱۱</sup> از گرادیان کاهشی باید به کمک معادله زیر همزمان همه پارامترها را به روز رسانی کنید.

$$(۴) \quad \theta_j \leftarrow \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

در هر مرحله از گرادیان کاهشی با به روز رسانی  $\theta$ ، به مقدار مینیمم تابع  $J(\theta)$  نزدیک‌تر می‌شویم. متغیر  $\alpha$  نرخ یادگیری<sup>۱۲</sup> است که آن را بسیار کوچک در نظر می‌گیریم (مثلاً  $\alpha = 0.01$ ). همیشه مقدار اولیه  $\theta$  را یک مقدار کوچک تصادفی در حوالی صفر انتخاب می‌کنیم (می‌توان از توزیع نرمال با واریانسی کوچک و میانگین صفر استفاده کرد).

در نهایت تابع هزینه  $J(\theta)$  (معادله‌ی ۲) را در قالب تابع `computeCost` در کلاس LinearRegression پیاده سازی می‌کنیم.

<sup>۸</sup> vector

<sup>۹</sup> fit

<sup>۱۰</sup> bias term

<sup>۱۱</sup> iteration

<sup>۱۲</sup> learning rate

## مشکلات متداول:

- در هر مرحله از گرادیان کاهش پارامترهای  $\theta$  را بصورت همزمان آپدیت کنید. یعنی در هر مرحله بعد از به روز رسانی یکی از اعضای بردار  $\theta$ ، بردار  $\theta$  را بصورت جزیی به روز رسانی نکنید و مقدار  $h_{\theta}(x^i)$  را محاسبه کنید. در پایان حلقه بردار  $\theta$  را کاملاً به روز رسانی کنید.
- به یاد داشته باشید در هر مرحله از گرادیان کاهش شما تنها فضای  $\theta$ های ممکن را جست و جو می‌کنید و نباید  $X$  و  $y$  را تغییر دهید.

## تست کردن کد نهایی:

یک راه ساده برای تست کردن کد چاپ کردن مقدار  $J(\theta)$  در هر مرحله است. اگر مقدار آن با گذر زمان در هر مرحله بصورت یکنواخت در حال کاهش است یعنی کد شما درست کار می‌کند.

بعد از اتمام پیاده‌سازی کد، مدل خود را بر روی داده‌های univariateData آموزش دهید و سپس تابع plotRegLineID در فایل test\_linreg\_univariate.py را اجرا کنید.

```
from test_linreg_univariate import plotRegLineID
from linreg import LinearRegression
X = np.c_[np.ones((n, 1)), X] # if you didn't this step already
lr_model = LinearRegression(alpha = 0.01, n_iter = 1000)
lr_model.fit(X, y)
plotRegLineID(lr_model, X, y)
```

خروجی باید یک نمودار مانند Figure ۲ باشد.

## Understanding Gradient Descent

در این قسمت نیازی به پیاده‌سازی کد ندارید. برای درک بهتر کد پیاده‌سازی شده در مرحله قبل، تابع هزینه و مسیر رسیدن به مینیمم آن که توسط الگوریتم گرادیان کاهش انتخاب شده را رسم می‌کنیم.

برای مجموعه داده‌های تک متغیره می‌توانیم تغییر تابع هدف را در فضای  $\theta_0$  و  $\theta_1$  به عنوان یک نمودار سطحی<sup>۱۳</sup> و یک نمودار کانتور<sup>۱۴</sup> ترسیم کنیم تا شکل محدب و نزولی آن را نشان دهیم. خط آبی در شکل ۳ مسیر طی شده توسط نزول گرادیان را نشان می‌دهد و نقاط سرخابی نقاط را در هر تکرار نشان می‌دهد.

<sup>۱۳</sup> Surface plot

<sup>۱۴</sup> Contour plot

موعد تحویل:

معیار ارزیابی شما:

گزارش کار:

۱. نکته مهم در گزارش نویسی روشن بودن پاسخ ها می باشد، اگر فرضی برای حل سوال استفاده میکنید حتما آن را ذکر کنید، اگر جواب نهایی عددی است به صورت واضح آن را بیان کنید.
۲. هرگونه شباهت در گزارش به منزله تقلب می باشد و کل نمره تمرین صفر می باشد.
- ۳.