



## تمرین کامپیوتری سوم



سیستم‌های عامل - پاییز ۱۳۹۹

گزارش کار

دانشکده مهندسی برق و کامپیوتر

نام و نام خانوادگی:

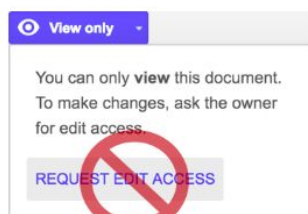
مهار کریمی

تاریخ:

۲۴ آذر ماه ۱۳۹۹

استاد:

دکتر مهدی کارگهی

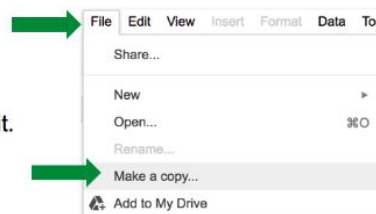


### How to use this template:

This is a view-only file and cannot be edited.

**Create your own copy** of this template to edit.

In the menu, click **File > Make a copy...**



2

مقدمه

3

پیاده‌سازی سری

3

سوال اول

3

سوال دوم

3

جدول اول

3

پیاده‌سازی چندریسه‌ای

3

سوال سوم

4

سوال چهارم

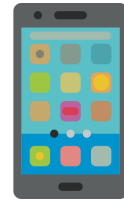
4

سوال پنجم

4

جدول دوم

## مقدمه



در این تمرین به تحلیل داده‌هایی که از مشخصات و قیمت فروش گوشی‌های موبایل جمع‌آوری شده‌است پرداخته شده است. در ابتدا برنامه اقدام به خواندن و تجزیه مجموعه داده<sup>۱</sup>ی ارائه شده می‌کند و آنها را در حافظه خود ذخیره می‌کند. پس از استخراج داده‌ها و ویژگی‌های آنها، برنامه اقدام به نرمال‌سازی<sup>۲</sup> داده‌ها و در نهایت اقدام به تعیین طبقه قیمتی گوشی‌ها می‌کند. این تمرین به دو روش این مسئله پیاده‌سازی شده است که در ادامه گزارش، نتایج حاصل آمده است.



---

<sup>۱</sup> Dataset

<sup>۲</sup> Data Normalization

## پیاده‌سازی سری

### سوال اول

چرا برای پیاده‌سازی یک برنامه بصورت چندریسه‌ای، بهتر است ابتدا این برنامه بصورت سری پیاده‌سازی شود؟  
جواب: پیاده‌سازی و اشکال‌زدایی (Debugging) برنامه بصورت سری آسانتر از پیاده‌سازی چندریسه‌ای است و برای بررسی درستی الگوریتم مورد نظر، روش بهتری نسبت به پیاده‌سازی چندریسه‌ای می‌باشد. هم‌چنین پس از پیاده‌سازی سری، تشخیص بخش‌هایی از کد که می‌توانند بصورت موازی اجرا شوند سریعتر انجام می‌شود.

### سوال دوم

با بررسی زمان اجرای بخش‌های مختلف برنامه، <sup>3</sup>Hotspot های برنامه را مشخص کنید.

جواب: می‌دانیم برای یک پردازش تک‌ریسه‌ای، در خروجی برنامه time رابطه زیر برقرار است:

$$T_{real} = T_{user} + T_{sys}$$

با بررسی خروجی time در چند اجرای متوالی، سهم  $T_{user}$  بیشتر از  $T_{sys}$  است؛ پس می‌توانیم نتیجه بگیریم که برنامه ما CPU-Bound می‌باشد و بیشتر زمان اجرای برنامه صرف انجام محاسبات و تجزیه فایل‌های ارائه شده می‌گردد.  
در پیاده‌سازی ما، توابع زیر بیشترین زمان اجرا را به خود اختصاص می‌دهند:

وظیفه تابع	نام تابع
خواندن و تجزیه نمونه‌های ورودی در فایل train.csv	Evaluator::read_data
نرمال‌سازی مقدار featureها	Evaluator::normalize_data
دسته‌بندی هر نمونه ورودی و بررسی عملکرد مدل به ازای ورودی مورد نظر	Evaluator::evaluate_model

<sup>3</sup> توابعی که در برنامه‌تان بیشترین زمان اجرا را به خود اختصاص می‌دهند.

## جدول اول

زمان‌های اجرای ۶ اجرای متوالی از برنامه و میانگین آن‌ها را بازای ورودی نمونه‌ای که در شرح تمرین آمده است، در جدول زیر بیاورید.

میانگین	اجرای ششم	اجرای پنجم	اجرای چهارم	اجرای سوم	اجرای دوم	اجرای اول
0.0398	0.040	0.039	0.040	0.042	0.039	0.039

توضیح: مقادیر بالا در مقیاس ثانیه هستند.

## پیاده‌سازی چندریسه‌ای

### سوال سوم

اگر هنگام موازی‌سازی برنامه به زمان اجرای بیشتری نسبت به حالت سری برخورد کنید، چه رویکردهایی را برای کاهش زمان اجرا و استفاده حداکثری از موازی‌سازی پیش می‌گیرید؟

جواب: برای اینکه ریشه‌ها بتوانند به صورت همزمان از منابع سیستم استفاده کنند، تعداد آن‌ها را تا حد امکان کم نگه می‌داریم (افزایش تعداد ریشه‌ها باعث بوجود آمدن سربار اجرا برای تعویض متن بین ریشه‌ها و زمان‌بندی آنها می‌شود، زیرا تعداد ریشه‌هایی که بطور همزمان اجرا می‌شوند حداکثر برابر تعداد هسته‌های پردازشی منطقی سیستم است). همچنین، تنها بخشهایی به ریشه‌ها اختصاص یافته‌اند که قابل موازی‌سازی باشند و محاسبات روی مجموعه‌های کوچک داده به صورت سری انجام شده‌است.

### سوال چهارم

در هنگام پیاده‌سازی این بخش، به چه چالش‌هایی برخورد کردید و بیان کنید که به چه صورت آن‌ها را رفع کردید.

جواب: می‌دانیم که همه دسته‌بندها از مجموعه یکتایی از بردارهای وزن برای دسته‌بندی استفاده می‌کنند؛ به همین دلیل، پس از یک بار خواندن فایل بردارهای وزن، برای ارزیابی مدل به ازای هر ریشه، اطلاعات مربوط به این بردار نیز با ارجاع پاس داده می‌شوند. همچنین، برای پیشگیری از استفاده متغیرهای global، استراکچری با عنوان Parallel\_Arg ساخته‌ایم که همه اعضای آن بصورت رفرنس هستند. تعریف این استراکچر به صورت زیر است:

```
struct Parallel_Arg {
    const std::vector<std::vector<double> >& weights;
    std::pair<int, int>& total;
    std::string dataset_dir;
    int idx;
    pthread_mutex_t& mutex;
    Evaluator pe;
};
```

## سوال پنجم

با توجه به تجربه‌ای که در پیاده‌سازی این تمرین بدست آوردید، به نظر شما در چه مواقعی از قفل<sup>4</sup> در یک طراحی چندریسه‌ای ضروری است؟ تاثیر استفاده از قفل‌ها را بر روی کارایی<sup>5</sup> سامانه بیان کنید.

جواب: استفاده از قفل تنها زمانی ضروری است که بخواهیم مقدار متغیری را توسط چند ریشه تغییر دهیم؛ برای مثال در پیاده‌سازی ما، زوج مرتب total برای نگه داشتن تعداد کل نمونه‌ها و تعداد حدس‌های درست به کار رفته است. هر ریشه پس از پایان محاسبه روی فایل متناظر خود، این مقدار را بروز رسانی می‌کند. برای جلوگیری از ایجاد اشتباه در محاسبات، از یک قفل برای مشخص کردن بخش مربوط به بروز رسانی این متغیر استفاده کرده‌ایم. با توجه به این که سایر بخش‌های داده یا فقط توسط هر ریشه خوانده می‌شوند و یا بطور محلی برای همان ریشه هستند، استفاده از قفل در سایر مراحل محاسبه ضروری نیست.

با قفل کردن قسمتی از کد، همه ریشه‌ها (بجز ریشه‌ای که در حال حاضر در حال اجرا می‌باشد یا قفل را فعال کرده است) برای اجرای کد قفل شده باید منتظر بمانند تا ریشه قفل کننده اجرای محدوده بحرانی (Critical) را به پایان برساند. به همین دلیل، استفاده از قفل سبب اضافه شدن به سرپار اجرای برنامه ما شود و بایستی از قفل‌ها بطور محدود و بهینه استفاده کنیم.

## جدول دوم

زمان‌های اجرای ۶ اجرای متوالی از برنامه و میانگین آن‌ها را بازای ورودی نمونه‌ای که در شرح تمرین آمده است، در جدول زیر بیاورید.

<sup>4</sup> Lock

<sup>5</sup> Performance

میانگین	اجرای ششم	اجرای پنجم	اجرای چهارم	اجرای سوم	اجرای دوم	اجرای اول
0.0275	0.27	0.26	0.27	0.31	0.25	0.29

میزان تسریع (  $\frac{Serial\ Time}{Parallel\ Time}$  ) برنامه نسبت به حالت سری را در زیر بیاورید.

میزان تسریع	میانگین زمان اجرای موازی	میانگین زمان اجرای سری
1.45	0.0275	0.0398