



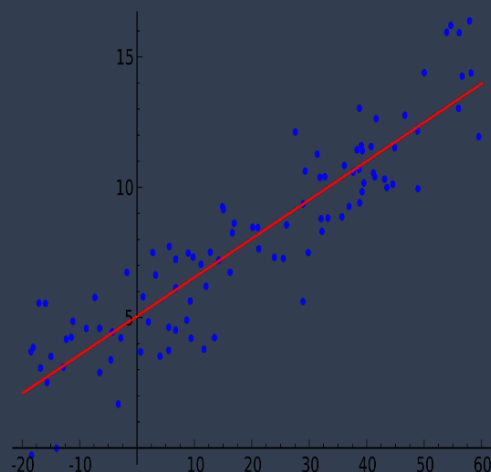
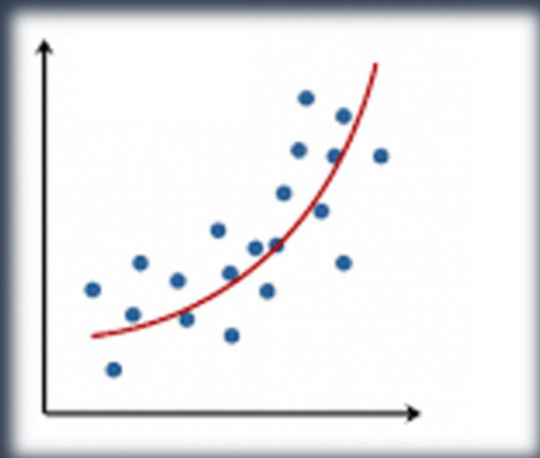
مقدمه و توضیح الگوریتم

"پیش بینی می شود در ابتدای سال آینده، تعداد بیماران شناسایی شده به ۱۰۰ نفر در روز برسد"

احتمالا شما هم در این چند وقت از جملاتی که خبر از پیش بینی موارد ابتلا به کرونا بدهد شنیده اید یا در برخی موارد نمودارهایی مربوط به آن را مشاهده کرده باشید. اما دقیقا مبنای این پیش بینی ها چیست؟

پیش بینی کردن همیشه یکی از دغدغه های مهم انسان بوده است. امروزه نیز بسیاری از فعالیتهای ما نیاز به پیش بینی کردن شرایط دیده می شود. از معاملات روز بازار بورس گرفته تا بررسی شرایط آب و هوایی در روزهای آینده. برای تخمین درست در قدم اول باید روابط بین متغیرها مشخص باشد.

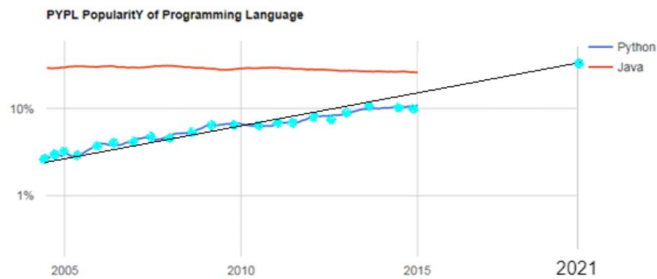
در این مینی پروژه قصد داریم با توجه به دانشی که از جبر خطی کسب کرده ایم، به مدل سازی روابط بین متغیرها بپردازیم.



رگرسیون:

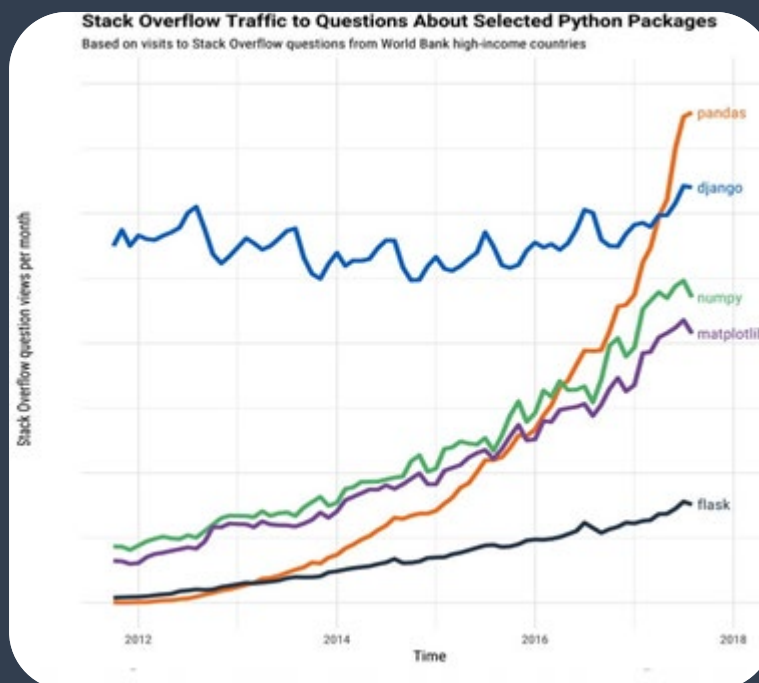
رگرسیون خطی: در نظر بگیرید می خواهیم با داشتن میزان محبوبیت زبان پایتون بین سال های ۲۰۰۵ تا ۲۰۱۷ محبوبیت آن را در انتهای سال ۲۰۲۱ پیش بینی کنیم.

Worldwide, Python is the most popular language, Python grew the most in the last 5 years (17.9%) and Java lost the most (-6.7%)



تصویر فوق یک نمودار از پیش رسم شده برای مقایسه محبوبیت‌ها در سال‌های مختلف می‌باشد. همانطور که از شکل پیداست، نمی‌توان یک رابطه‌ی خطی یافت که داده‌های ما دقیق پیش‌بینی شود اما می‌توان خطی را پیدا کرد که کمترین فاصله را از مجموعه نقاط داشته باشد و در این حالت امیدوار بود که تخمین ما از محبوبیت پایتون دارای کمترین خطا باشد. در اصل امیدواریم رابطه‌ی بین متغیر زمان و محبوبیت پایتون یک رابطه خطی باشد و به طور خطی پیشرفت کند.

رگرسیون غیرخطی:



عکس بالا مقایسه تعداد سوالات پرسیده شده در مورد هر یک از پکیج‌های معروف پایتون است.

فرض کنید می‌خواهیم تعداد سوالات پرسیده شده در مورد پکیج پاندا را در سال ۲۰۲۱ پیش‌بینی کنیم. نمودار نشان می‌دهد که انتخاب کردن یک رابطه‌ی خطی نزدیک به داده‌ها نتیجه نزدیکی به واقعیت نخواهد داشت و هر چه بیشتر از داده‌های واقعی فاصله بگیریم خطای پیش‌بینی بیشتر می‌شود اما به نظر می‌رسد رابطه آن می‌تواند نزدیک به یک چندجمله‌ای از درجه دو باشد و با یافتن هم‌چنین چندجمله‌ای می‌توان تخمین‌های با خطای کمتری داشت.

پس باید چندجمله‌ای زیر را به نحوی پیدا کنیم که به مقادیر واقعی نزدیک‌ترین حالت را داشته باشد:

$$y(t) = b_0 + b_1 t + b_2 t^2$$

و این به این معناست که در حالت ایده‌آل بودن نمودار انتظار داریم که:

$$y(t_1) = b_0 + b_1 t_1 + b_2 t_1^2$$

$$y(t_2) = b_0 + b_1 t_2 + b_2 t_2^2$$

$$\vdots \quad \quad \quad \vdots$$

$$y(t_n) = b_0 + b_1 t_n + b_2 t_n^2$$

و اگر معادلات بالا را به فرم ماتریسی بنویسیم:

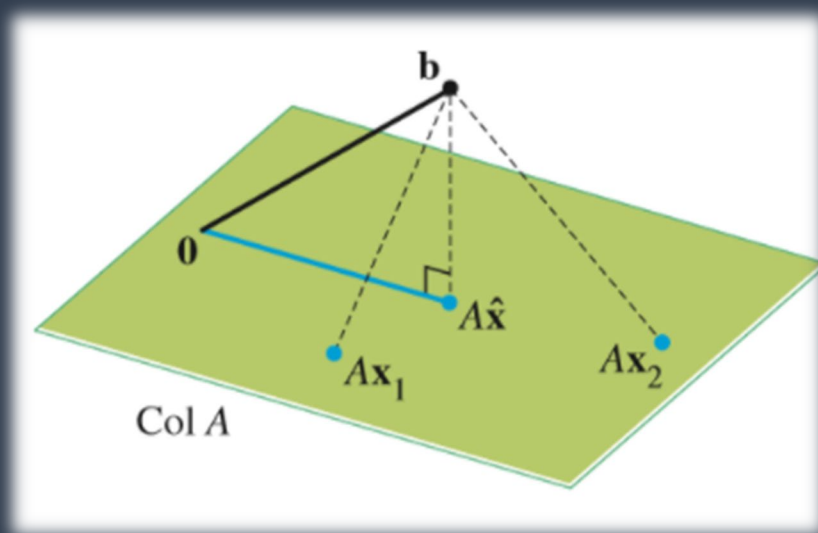
$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & t_1 & t_1^2 \\ 1 & t_2 & t_2^2 \\ \vdots & \vdots & \vdots \\ 1 & t_n & t_n^2 \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix}$$

با جایگذاری کردن زمان و تعداد سوالات، ماتریس را حل کرده و ضرایب دقیق را بدست می‌آوریم. ولی همان‌طور که قبلاً گفته شد، این تنها برای حالت ایده‌آل صدق می‌کند. برای حالت‌های دیگر باید ضرایبی را بدست بیاوریم که نزدیک‌ترین حالت را به نمودار اصلی داشته باشند و برای این کار باید از مسئله کمترین مربعات یا *Least Square* کمک بگیریم.

مسئله کمترین مربعات:

ماتریس $Ax = b$ را در نظر بگیرید. حالتی را در نظر بگیرید که مقدار x برای پاسخ یافت نشود. در این حالت ما به دنبال نزدیک ترین بردار به b هستیم که در این معادله صدق کند.

تصویر b روی A نزدیک ترین بردار می باشد که اگر با $A\hat{x}$ آن را نمایش دهیم خواهیم داشت :



در این مثال $col A$ یک صفحه در فضای سه بعدی است و نزدیک ترین جواب $A\hat{x}$ می باشد.

و واضح است که $b - A\hat{x}$ به تمامی بردارهای $Col A$ عمود است و این یعنی خواهیم داشت:

$$a_j^T (b - A\hat{x}) = 0$$

$$\Rightarrow A^T (b - A\hat{x}) = 0$$

$$\Rightarrow A^T b - A^T A\hat{x} = 0$$

$$\Rightarrow A^T A\hat{x} = A^T b$$

از معادله بالا ما به دنبال \hat{x} هستیم که ضرایب چندجمله ای ما (در اینجا درجه ۲) می باشد.

برای حل دستگاه معادلات بالا می توانید از پروژه اول خود استفاده کنید و یا از این تابع آماده استفاده کنید.

شرح پروژه:

یک فایل از سهام گوگل در روز های متخلف در اختیار شما قرار گرفته است.

این فایل به صورت *CSV* می باشد و شامل هفت ستون می باشد. این هفت ستون مختلف به ترتیب نشان دهنده تاریخ روز، شروع قیمت سهام در ابتدای روز، بالا ترین قیمت سهام در یک روز، کمترین قیمت سهام در یک روز، قیمتی سهام در انتهای روز، حجم معاملات و نام سهام می باشد.

در این پروژه قصد داریم قیمت سهام گوگل را در ابتدای سال ۲۰۰۶ تا پایان سال ۲۰۱۷ بررسی کنیم (چون از *regression* تک متغیره استفاده می کنیم تعدادی از داده ها حذف شدند تا نمودار نهایی برای شما ملموس تر باشد)

نحوه انجام پروژه:

ابتدا فایل *CSV* را دانلود کرده و آن را بخوانید. برای خواندن فایل *CSV* می توانید از این [لینک](#) استفاده کنید (استفاده از کتابخانه و توابع دیگر بلامانع است) ستون دیتای مورد نظر ما *Open* می باشد و ما به بررسی مقادیر این ستون در روزهای مختلف می پردازیم.

به غیر از ده سطر آخر فایل، از تمامی سطرها برای بدست آوردن ضرایب معادلات استفاده کنید. سپس با توجه به ضرایبی که بدست آوردید از ده روز آخر برای بررسی خطای تخمین خود استفاده کنید و آن را نمایش دهید.

در مرحله اول کد شما باید رگرسیون خطی را بررسی کند. برای رگرسیون خطی باید دستگاه معادله زیر حل شود که t شماره روز شماسست و برای هر سطر یک واحد از سطر بالایی خود بیشتر است (روزها متوالی در نظر گرفته شده اند)

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & t_1 \\ 1 & t_2 \\ \vdots & \vdots \\ 1 & t_n \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \end{bmatrix}$$

سپس برای ده روز آخر بررسی کنید که چقدر خطا داشتید و به ازای هر روز، در خروجی سه خط به فرمت زیر نمایش داده شود:

calculated value: 986

actual value: 1060.09

error: - 74.09

در گام بعدی سراغ رگرسیون درجه ۲ می‌رویم و تمامی مراحل بالا (حل دستگاه معادله و محاسبه خطا) را برای این رگرسیون نیز انجام می‌دهیم.

حال با بررسی میزان خطاهای هر کدام تشخیص می‌دهیم که کدام یک از رگرسیون‌ها مناسب داده‌های ما می‌باشد و در آخر نمودار مقادیر تخمینی و مقادیر واقعی مربوط به رگرسیون بهتر را به شکل زیر نمایش می‌دهیم:



نقاط قرمز رنگ مقادیر واقعی و خط آبی رنگ مقادیر تخمین زده شده برای هر روز می‌باشد.

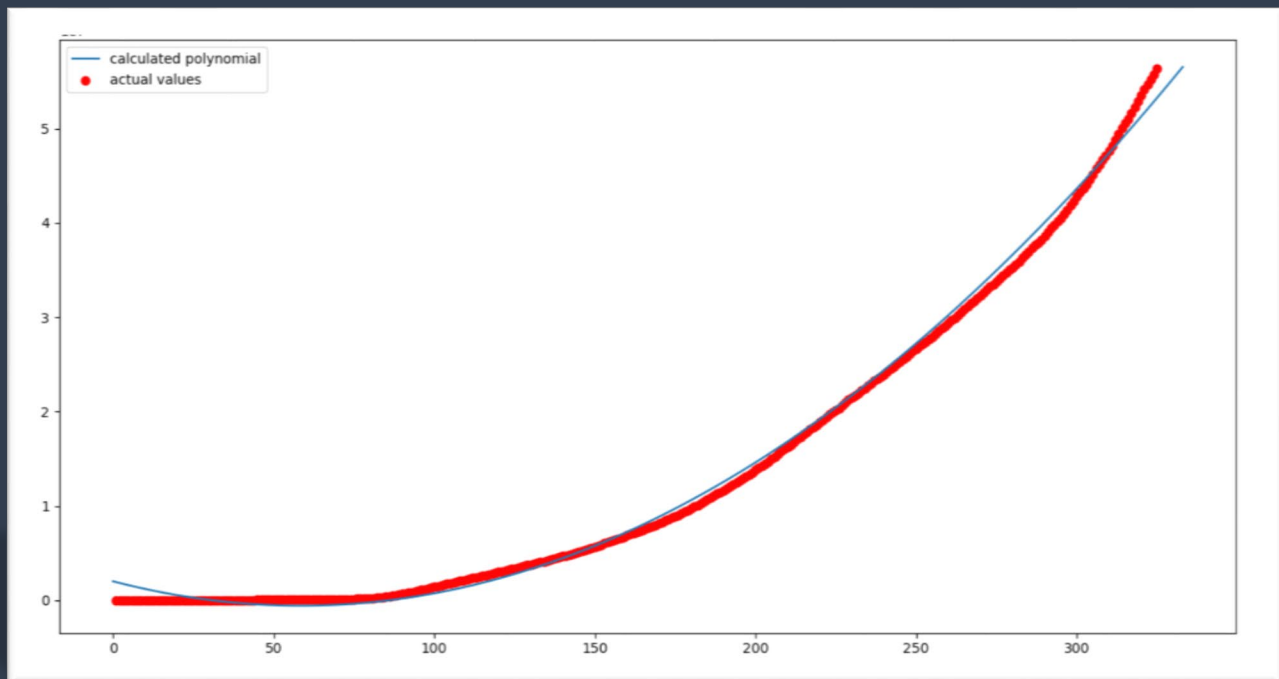
برای رسم نمودار نیز می‌توانید از توضیحات ویدیو استفاده کنید.

قوانین

۱. مینی پروژه به صورت **انفرادی** می‌باشد. تقلب‌ها توسط سامانه کوئرا چک می‌شود اما تصحیح توسط تیم تدریس جاری صورت می‌گیرد.
۲. برای خواندن فایل *CSV* می‌توان از هر کتابخانه‌ای استفاده کرد.
۳. با توجه به مجازی بودن درس و عدم امکان تحویل حضوری، انتظار می‌رود کد پیاده‌سازی شده از توابع مختلف تشکیل شده باشد که فهم مطلب را آسان‌تر کند. وجود مستند برای توضیح کد به شدت استقبال می‌گردد.
۴. برای انجام عملیات ضرب ماتریس‌ها (*dot()*) و ترانپاده (*transpose()*) می‌توانید از توابع آماده استفاده کنید.
۵. فرمت فایل ارسالی بایستی به صورت *miniproject3_STUDENTNUMBER_Name* باشد. به طور مثال:
miniproject3_9628099_SoroushMehraban

برای علاقه مندان

دیتای مبتلایان به کرونا را از روز اول تا به امروز می‌توانید از [اینجا](#) دریافت کنید. فایل `total_cases.csv` در دومین ستون خود حاوی تعداد مبتلایان جهانی کرونا می‌باشد و این آمار به صورت روزانه موجود است. می‌توانید از رگرسیون درجه ۲ خود برای این دیتا استفاده کنید. به طور تقریبی این نمودار حاصل می‌شود :



موفق باشید

تیم تدریس یاری جبر خطی کاربردی

آذر ۹۹