

نام و نام خانوادگی: پیمان هاشمی

شماره دانشجویی: 400131032

تمرین: سری اول

بخش اول

سوال (1)

الف: نادرست - زیرا ممکن است دیتای مربوط به آموزش دارای همه انواع مقادیر کلاس داده نباشد و در این صورت خطای مربوط به آموزش کاهش پیدا میکند ولی در صورتی که خطای آزمون به شدت بالا رفته به دلیل اینکه دیتا تنها بر روی یک نوع خاصی دیتا train شده است.

ب: نادرست - با کاهش تعداد داده های آموزش الگوریتم under train میشود.

ج: نادرست - در مدل هایی که کم بر ارزش شده اند به دلیل اینکه اندازه کافی واریانس ندارد باعث میشود خطای آزمون بالا باشد در این حالت یعنی مدل پیچیدگی لازم را ندارد. در این حالت باید پیچیدگی مدل را افزایش داد تا زمانی که خطا کم شده و خطا های مربوط به واریانس از بین بروند.

د: درست - به دلیل اینکه بر خلاف MSE یا RMSE به توان 2 نمیرسد و در قدر مطلق قرار دارد.

سوال (2)

به این معنی میباشد که پیچیدگی مدل استفاده شده بسیار زیاد بوده و به اصطلاح مدل overfit شده است.

برای حل این موضوع میتوان در مرحله اول پیچیدگی مدل را کاهش داد.

باید توازن بین واریانس و بایاس برقرار باشد.

سوال (3)

(الف)

Handwritten mathematical derivation for linear regression:

$$\begin{aligned} J(w) &= \sum_{i=1}^n (y_i - w^T x_i)^2 \Rightarrow J(\theta_0, \theta_1) = \frac{1}{2n} \sum_{i=1}^n (y_i - w^T x_i)^2 \Rightarrow J(\theta) = \frac{1}{2n} (Xw - y)^T (Xw - y) \\ y &= Xw + \epsilon \Rightarrow \epsilon = y - Xw \Rightarrow \|\epsilon\|_2 = \|y - Xw\|_2 \\ \|\epsilon\|_2^2 &= \|y - Xw\|_2^2 = (y - Xw)^T (y - Xw) = (y - Xw)^T (y - Xw) = (Xw - y)^T (Xw - y) \\ &= (Xw)^T Xw - (Xw)^T y - y^T Xw + y^T y = (Xw)^T Xw - 2(Xw)^T y + y^T y \\ \frac{\partial J}{\partial w} &= 2X^T Xw - 2X^T y = 0 \Rightarrow w = (X^T X)^{-1} X^T y \end{aligned}$$

(ب)

یکی از مشکلات علامت است که آن را با ضرب داخلی حل میکنیم
در ضرب $X^T * X$ زمانی که ستون ها بهم وابسته باشند مشکل ساز میشود که باید از PCA استفاده کنیم
زمانی تعداد ستون از تعداد داده ها بسیار بیشتر باشد که نتواند به طور خوب پوشش بدهد
(د) جمله منظم ساز به نویز حساس تر است

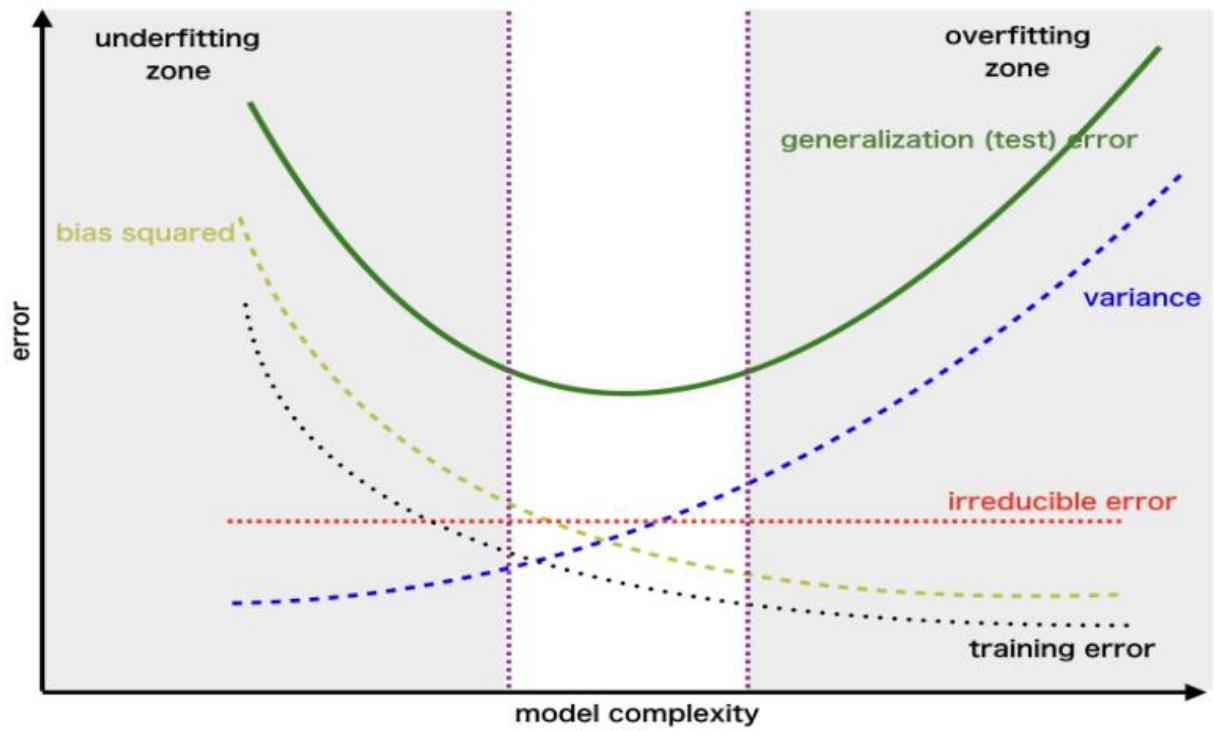
سوال 4

الف: از گرادیان نزولی mini batch استفاده میکنیم زیرا در داده های بزرگ یا بیگ دیتا استفاده از معادله نرمال زمان بسیار زیادی را به خود اختصاص میدهد.
ب: در معادله نرمال تاثیری نمیگذارد و میتوان در حالت عادی از الگوریتم استفاده کرد ولی در استفاده از گرادیان نزولی سرعت بسیار کم میشود و در این حالت میتوانیم از feature scalling استفاده کنیم تا تعداد تکرار ها را کمتر کنیم.

سوال 5

در موارد با واریانس بالا الگوریتم به نوعی درگیر بیش برآزش میشود که حاصل تمرکز آن بر روی داده های نویز است ولی با اضافه کردن داده های بیشتر نیز تعداد داده های دارای نویز بالا میرود و الگوریتم قادر بیش برآزش نمیشود بنابر این تمرکز اصلی را بر روی داده های دیگر میگذارد و به نوعی از بیش برآزش خارج میشود پس اضافه کردن به داده های آموزش برای واریانس بالا مناسب است.
در موارد با بایاس بالا اضافه کردن به داده های آموزش معمولاً یا تغییری در بایاس ایجاد نمیکند یا باعث افزایش بایاس میشود.

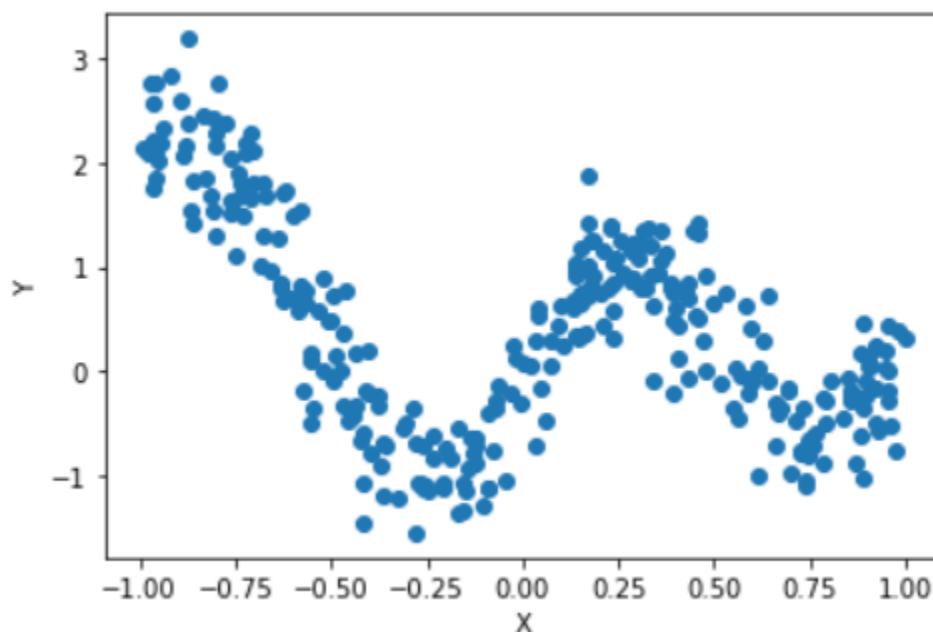
سوال (6)



بخش دو

سوال 1

الف:

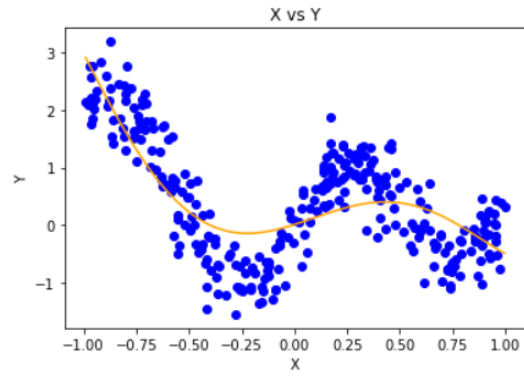


ب:

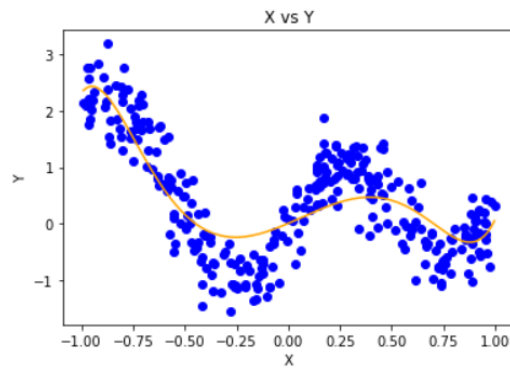
شافل کردن دیتا زمانی انجام میشود دیتاست شما بر مبنای یک ستون خاص مرتب شده است و یا زمانی که دارای high bias شده ایم از شافل کردن دیتا استفاده میکنیم تا بتوان یک توزیع مناسب برای داده ها ایجاد کنیم تا مجموعه دیتا آموزش دارای انواع دیتا باشد.

ج:

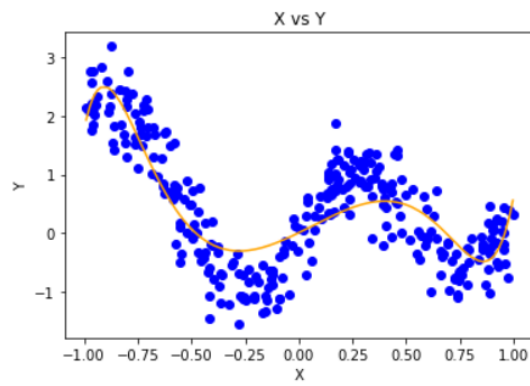
بررسی داده ها با 5 هزار بار تکرار و learning rate: 0.01



درجه 5

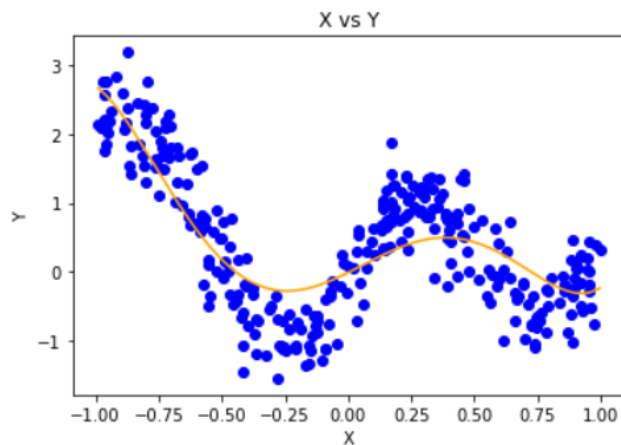


درجه 8

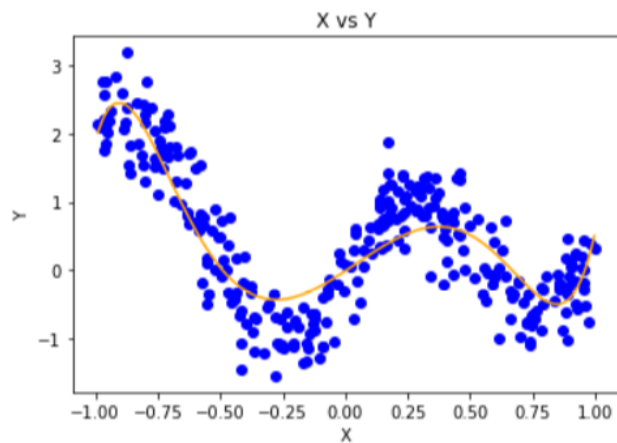


درجه 10

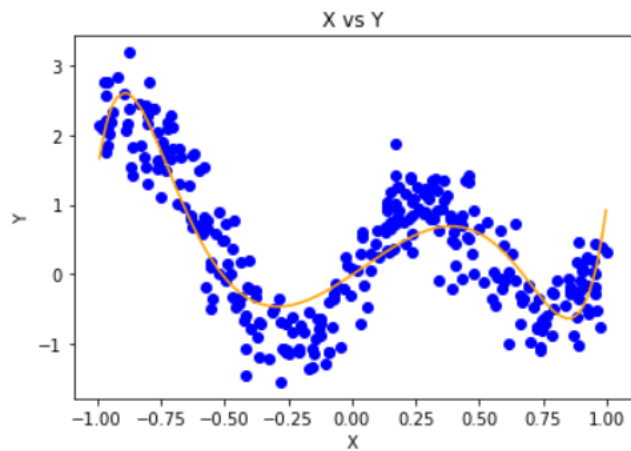
در نمودار های زیر الگوریتم بالا با تعداد تکرار 10000 بار انجام شده است



درجه 5



درجه 8



درجه 10

در درجه 10 با تعداد تکرار 10000 بار بیش برآزش اتفاق افتاده است.

در این الگوریتم معادله اصلی که با آن کار میشود برابر زیر است:

$$Y = ax^2 + bx + c$$

که در آن a, b, c متناظر برابر میباشند. هدف این است که ضرایب چند جمله ای پیدا شود. در ابتدا همه ضرایب به صورت رندوم اختصاص داده میشود و بنابر منحنی رسم شده نسبت آن را با x و y میسنجند و

حال با گرادیانت دیسنت سعی میکنیم تا لاس فانکشن را به کمترین میزان ممکن برسانیم. در این حالت نتایج بهتری بدست می آید و آن ها را در معادله میگذاریم. همین روند تا زمانی که منحنی به داده ها نزدیک شود ادامه میدهد.

سوال (2)

الف:

برای پر کردن داده های گمشده در یک فایل یک راه حل مشخص وجود ندارد و بنا بر مسئله و نوع داده، راهکار های حل این مشکل میتواند متفاوت باشد. راحت ترین راه برای مقابله با داده های گمشده حذف ردیف هایی میباشد که داده ها در آن قرار دارند. در این راه حل در صورتی که ردیف های دارای داده گمشده دارای اطلاعات حساسی باشند، این اطلاعات از بین میروند که داده آموزش را با مشکل مواجه میکند.

راه حل دیگر برای مقابله با داده های گمشده استفاده از دستور زیر میباشد. روش کار این دستور به این نحو است که از داده های بالا و پایین داده گمشده میانگین گرفته و آن را جایگزین داده گمشده میکند. این راه حل، مناسب دیتاست هایی با موارد گمشده زیاد نیست و همین طور اگر داده گمشده از ضریب همبستگی بالایی برخوردار باشد ممکن است باعث شود برآزش به صورت اشتباه انجام شود.

`Dataset. Interpolate()`

یک روش دیگر پر کردن داده های گمشده با مقدار صفر است که برای این روش از دستور زیر استفاده میشود.

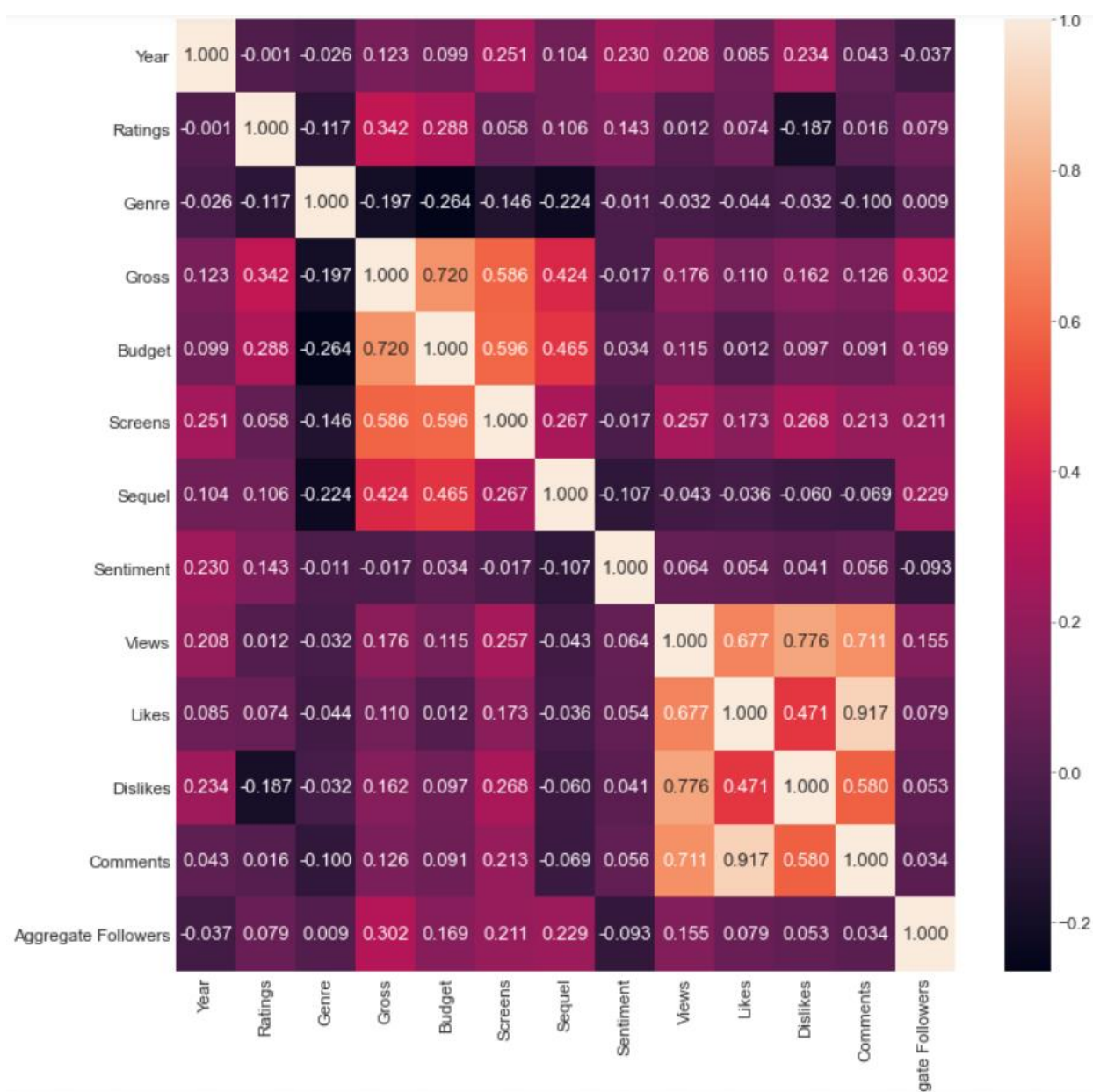
`Dataset.fillna(0)`

روشی دیگر بر کردن داده های گمشده با مقدار قبلی یا بعدی است که این روش نزی مانند روش های دیگر دارای نکات منفی از جمله آموزش اشتباه الگوریتم است.

یک روش دیگر برای مقابله با داده های گمشده استفاده از کتابخانه sklearn میباشد. در این کتابخانه دستوری به شکل زیر وجود دارد. با استفاده از دستور گفته شده میتوان مقادیر گمشده را با مقادیر نظیر، میانگین، بیشترین، کمترین، پر تکرار ترین یا عدد رندوم در آن ستون جایگزین کرد.

SimpleImputer

ب:



ج:

با توجه به نمودار و اینکه کلاس مورد هدف ما rating میباشد میتوان آن دسته از ویژگی هایی که دارای ضریب همبستگی بسیار نزدیک به عدد صفر دارند را حذف کرد که آن بسته به دیدگاه نسبت سوال دارد ولی با توجه به این که ستون های year, views, comments بین 0.02 و -0.02 است میتوان آن ها حذف کرد.

به این دلیل این ستون ها کمترین تاثیر را نسبت به rating دارند پس میشود از آن ها صرفه نظر کرد.

سوال (3)

الف: عملکرد آن به جز در مواردی که متغیر دسته ای مقادیر خیلی زیادی بگیرد، بسیار خوب است (معمولا از این روش برای متغیرهایی که بیش تر از ۱۵ مقدار متفاوت بگیرند، مناسب نیست. در برخی از مواردی که تعداد متغیرها کمتر است نیز امکان دارد گزینه مناسبی نباشد). کدبندی One-Hot ستون های دودویی (binary) جدیدی می سازد که هر یک مربوط به یکی از مقادیری هستند که متغیر به خود می گیرد. برای درک بهتر موضوع، مثالی در ادامه ارائه شده است.

فرض می شود که یک متغیر با عنوان Color در مجموعه داده وجود دارد که مقادیر آن Red ، Yellow و Green هستند. اکنون، کلیه مقادیر این متغیر به سه ستون جدا با عنوان های Red ، Yellow و Green تبدیل می شوند. فرض کنید پنج «نمونه» در مجموعه داده وجود دارد که مقدار متغیر Color برای آن ها به ترتیب برابر با Red ، Yellow ، Red ، Red و Yellow است. اکنون، برای نمونه اول که رنگ آن قرمز است در ستون Red عدد ۱ و در ستون های Yellow و Green عدد صفر وارد می شوند. برای نمونه دوم نیز که رنگ آن قرمز است به همین صورت عمل می شود. برای نمونه سوم که رنگ آن زرد است در ستون Red و Green مقدار صفر و در ستون Yellow مقدار ۱ وارد می شود. برای متغیر چهارم که رنگ آن Green است در خانه های Red و Yellow مقدار صفر و در خانه Green مقدار ۱ وارد می شود. برای سایر نمونه ها نیز به همین صورت عمل می شود.

Interger encoding :

جایی که هر برجسب منحصر به فرد به یک عدد صحیح نگاشت می شود.

One Hot Encoding : جایی که هر برجسب به یک بردار باینری نگاشت می شود. جاسازی آموخته شده: جایی که یک نمایش توزیع شده از دسته ها آموخته میشود