

# تمرین‌های سری یک Machine Learning

فرهاد دلیرانی

۹۶۱۳۱۱۲۵

[dalirani@aut.ac.ir](mailto:dalirani@aut.ac.ir)

[dalirani.1373@gmail.com](mailto:dalirani.1373@gmail.com)

**یادگیری با نظارت:** هرگاه داده های آموزش شامل نمونه هایی باشد که به وضوح مشخص شده باشد که خروجی چه باید باشد آنگاه می توان گفت در محدوده ی یادگیری با نظارت هستیم. مانند تشخیص رقم ها، اسپم بودن/نبودن ایمیل. در این دسته از روش های یادگیری، مجموعه ی آموزش به صورت زیر است:

$$(X_1, Y_1) (X_2, Y_2) \dots (X_n, Y_n)$$

که  $X_i$  داده ی آموزش است و  $Y_i$  برچسب آن است.

**یادگیری نیمه نظارتی:** یادگیری نیمه نظارتی نوعی یادگیری بین دسته ی یادگیری بانظارت و بدون نظارت است. در این نوع یادگیری برای آموزش هم از داده های با برچسب و هم بدون برچسب استفاده می شود.

**یادگیری بدون نظارت:** در یادگیری بدون نظارت مجموعه ی یادگیری هیچگونه برچسبی ندارد. یادگیری بدون نظارت را می توان به عنوان کاری دید که به دنبال ساختار و روابط میان داده های ورودی است. مانند دسته بندی تعدادی کتاب در موضوعات مختلف به طوری که کتاب های هر دسته از نظر موضوعی شباهت داشته باشند. همین طور از یادگیری بدون نظارت می توان به عنوان یک روش برای نمایش و آرایه ی داده ها به صورت قابل درک و معنا دار استفاده کرد.

**یادگیری تقویتی:** در یادگیری تقویتی برخلاف یادگیری بانظارت که هر نمونه در مجموعه ی آموزش برچسبی دارد، تابع  $target$  مشخصی که یک برچسب را به یک نمونه نسبت بدهد را نداریم. بلکه در این نوع یادگیری برای هر اکشن یک مقیاس داریم که مشخص می کند که اکشن چه میزان خوب یا بد بوده است. در یادگیری با نظارت مجموعه ی یادگیری به صورت (input, correct output) است در حالی که در یادگیری با تقویتی مجموعه ی یادگیری به صورت زیر است:

$$(Input, Some Output, Grade for this output)$$

این نوع یادگیری کاربردهای زیادی در یادگیری انجام بازی دارد.

**یادگیری برخط:** در این نوع یادگیری مجموعه‌ای که برای یادگیری استفاده می‌شود به صورت یکجا در اختیار نیست بلکه هر زمان یک نمونه در اختیار آن قرار می‌گیرد. به عبارت دیگر دیتاهای آموزش یک دنباله هستند هر عضو دنباله در زمان متفاوتی در اختیار الگوریتم قرار می‌گیرد.

در صورتی که مقدار تابع هزینه (خطا) بر روی ۸۰ درصد داده آموزش کم باشد ولی مقدار تابع هزینه برای ۲۰ درصد باقی مانده زیاد باشد متوجه می شویم مدل مورد استفاده دچار بیش برازش شده است زیرا خطای آن بر روی داده هایی که با آن ها آموزش دیده است کم است ولی توانایی پیشبینی درست برای داده هایی را که ندیده است را ندارد.

خطای MSE : خطای Mean Square Error خطایی است که میزان میانگین مربع خطاها را محاسبه می- کند. منظور از خطا اختلاف بین تخمین مدل و مقدار واقعی است. که از رابطه‌ی زیر محاسبه می‌شود:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (\hat{Y}_i - Y_i)^2$$

خطای RMSE : خطای Root-Mean Square Error یک مقیاس پرکاربرد برای سنجش خطای تخمین- های مدل و مقدار واقعی است. که میزان جذر میانگین مربع خطاها است که از رابطه‌ی زیر قابل محاسبه است:

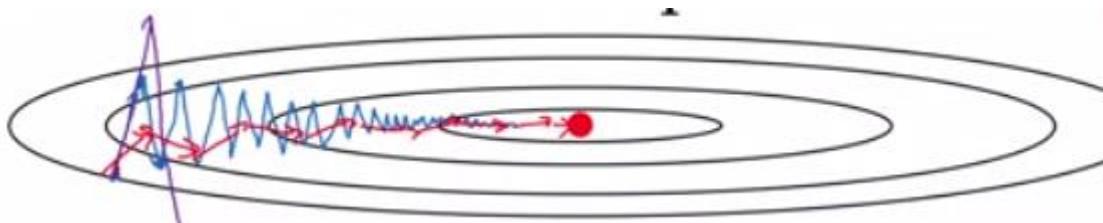
$$\text{RMSE}(\theta_1, \theta_2) = \sqrt{\text{MSE}(\theta_1, \theta_2)} = \sqrt{E((\theta_1 - \theta_2)^2)} = \sqrt{\frac{\sum_{i=1}^n (x_{1,i} - x_{2,i})^2}{n}}.$$

برای دیتاست با داده‌های پرت و ناهنجار استفاده از RMSE بهتر است زیرا در مقابل داده‌های پرت با توجه به جذری که در آن در مقایسه با MSE هست Robust تر است و مقدارش به صورت در مقابله با آن داده‌ها کم تر تغییر می‌کند.

ایده‌ی اصلی Gradient Descent With Momentum استفاده از Exponentially Weighted

Average گرادیان ها است. این روش همیشه سریع تر از روش گرادیان نزولی عادی عمل می کند.

در تصویر کانتورهای یک تابع هزینه را مشاهده می کنید که در راستای افق کشیده تر از راستای عمودی است. بردارهای آبی گرادیان نزولی عادی را نشان می دهد. همان طور که مشاهده می کنید تعداد گام هایی که گرادیان طی می کند بسیار زیاد است و همین طور بردارهای آبی زاویه ی زیادی با نقطه ی آپتیمال دارند و همین باعث می شود که نتوانیم ضریب یادگیری بزرگی انتخاب کنیم. بردارهای قرمز به عنوان مثال گرادیان نزولی به همراه تکانه را نشان می دهد، این روش همان طور که در تصویر زیر مشاهده می کنید باعث شده است که مولفه ی عمودی هر گام کوچک باشد و مولفه ی افقی آن بلند و این در مثال زیر باعث می شود بسیار سریع به نقطه ی آپتیمال برسیم زیرا هر گام نسبت به گرادیان نزولی بیشتر در جهت نقطه ی آپتیمال است و گرادیان نزولی همراه تکانه مسیر مستقیم تری را طی می کند.



تتای صفر ضریب  $X_0$  است.  $X_0$  ویژگی است که ما به sample ها اضافه می‌کنیم و مقدارش همیشه برابر یک است. از آنجایی که مقدار  $X_0$  همواره برابر یک است ضریب معادل آن را وارد Regularization نمی‌کنیم زیرا تاثیر چندانی بر نتیجه‌ی نهایی ندارد. تتای صفر نقش بایاس (عرض از مبدا) را دارد.

$$\frac{\partial J(\theta)}{\partial \theta_j} = \left( \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \right) + \frac{\lambda}{m} \theta_j \quad \text{for } j \geq 1$$

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

$$\min_{\theta} J(\theta)$$

افزایش داده‌های آموزش، احتمال رخ دادن  $\text{over fitting}$  را کم می‌کند.

زیرا اضافه کردن تعداد بیشتری دیتا  $\text{diversity}$  را بیشتر می‌کند و هر چه  $\text{diversity}$  بیشتر باشد مدل که پیدا می‌شود جنرال‌تر است و در مدل جنرال‌تر احتمال  $\text{over fitting}$  کم‌تر است.



- افزایش داده‌های مجموعه‌ی آموزش: هرچه تعداد داده‌های مجموعه‌ی آموزش بیشتر باشد diversity داده‌ها بیشتر می‌شود و diversity بالا منجر به می‌شود مدلی که بر اثر آموزش ایجاد می‌شود جنرال‌تر باشد و در مدل جنرال احتمال رخ داد بیش برآزش کم می‌شود.
- کاهش تعداد features: هر چه تعداد features بیشتر باشد امکان پیدا کردن مدل‌های پیچیده‌تری فراهم می‌شود و بنابراین ممکن است تعداد فیچرهای زیاد منجر به این شده باشند که مدلی که پیدا کرده ایم برای داده‌ها پیچیده باشد به طوریکه خطای مجموعه‌ی آموزش بسیار کم شده باشد در حالی که خطا برای داده‌هایی که تاکنون مشاهده نکرده‌ایم زیاد باشد. با کم کردن فیچرها می‌توان از پیچیدگی مدل کاست و از بیش برآزش دوری کرد.
- رگولاریزیشن: رگولاریزیشن تکنیکی در یادگیری ماشین است که با استفاده از آن می‌توان از بیش-برآزش دوری کرد. ایده‌ی اصلی پشت آن افزودن پناستی به تابع هزینه است که باعث می‌شود مقدار تابع هزینه در صورت پیچیده‌شدن و over fitting زیاد شود و از آنجایی که به دنبال مینیموم کردن هزینه‌ها هستیم باعث ایجاد تعادل بین پیچیدگی زیاد و رسیدن به مینیموم تابع هزینه‌ها می‌شود یکی از راه‌های رگولاریزیشن در MSE را در تصویر زیر مشاهده می‌کنید:

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

$$\min_{\theta} J(\theta)$$

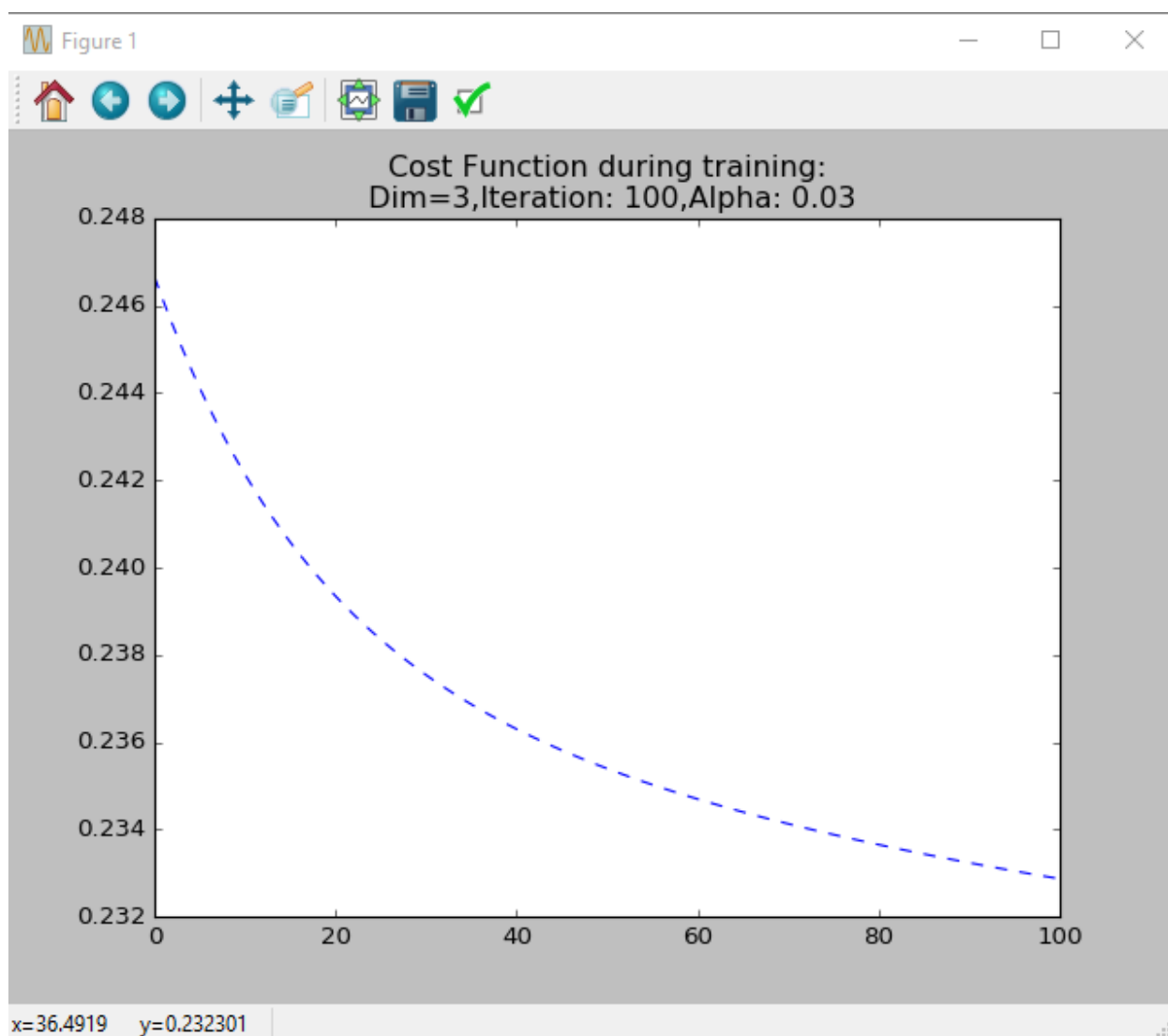
- تغییر الگوریتم یادگیری: الگوریتم‌های مختلف یادگیری ویژگی‌های مختلفی دارند با تغییر الگوریتم‌های یادگیری و انتخاب مدلی مناسب می‌توان از بیش برآزش جلوگیری کرد.

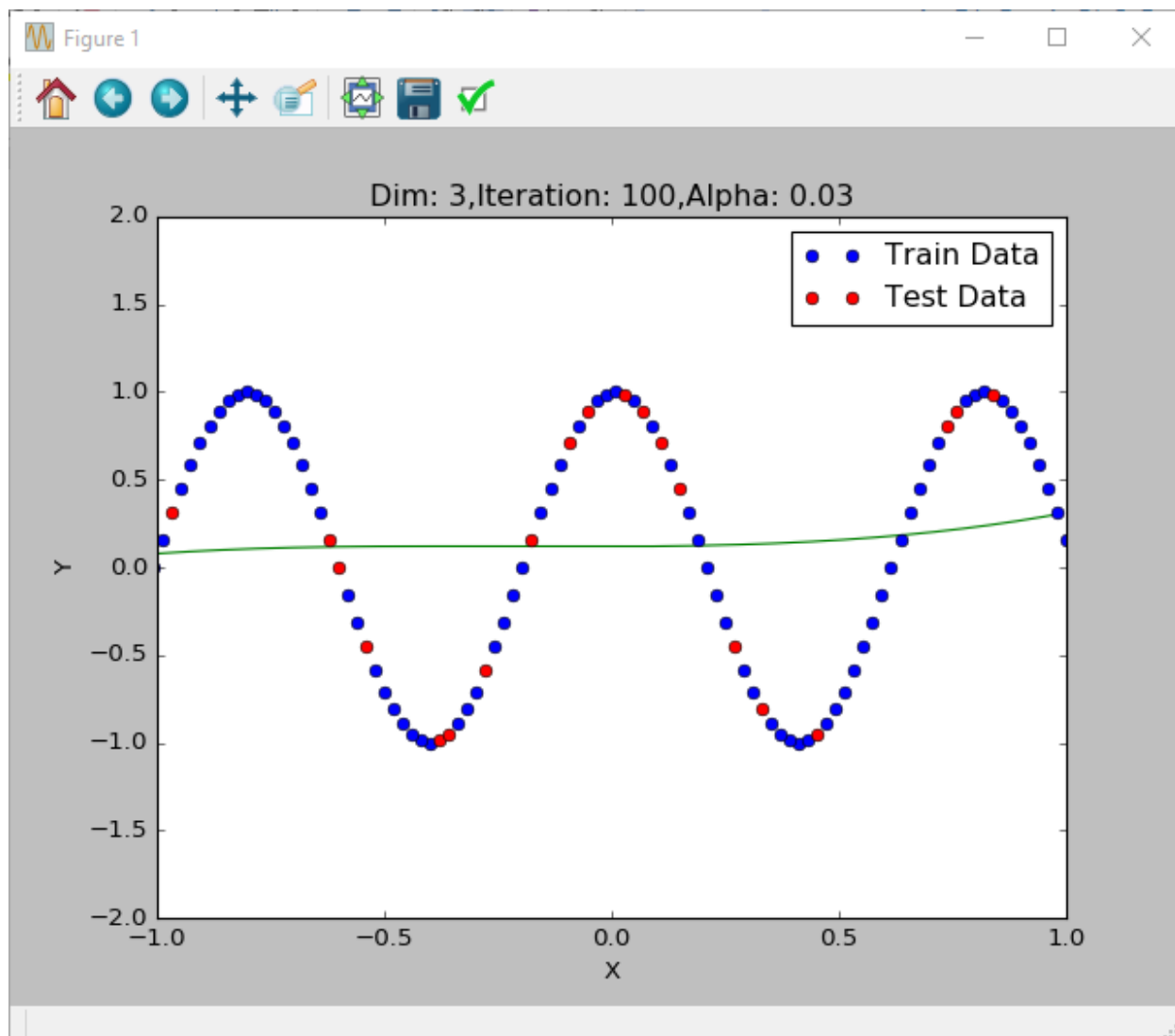
(الف)

کد این بخش از سوال هشت در فایل `gradientDescent.py` موجود است که به خوبی کامنت گذاری شده است.

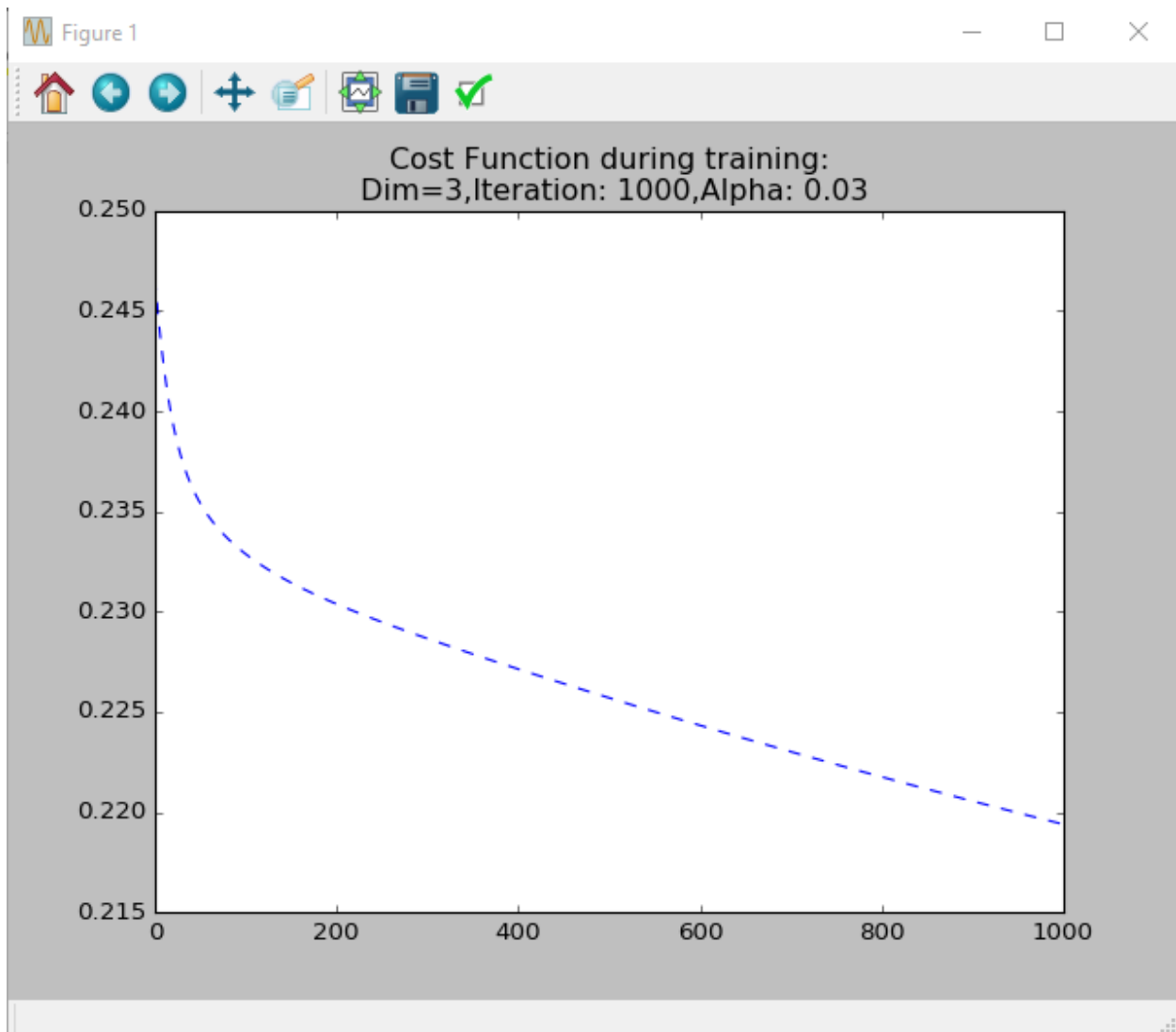
نتایج حاصل را در تصویرهای زیر می بینید:

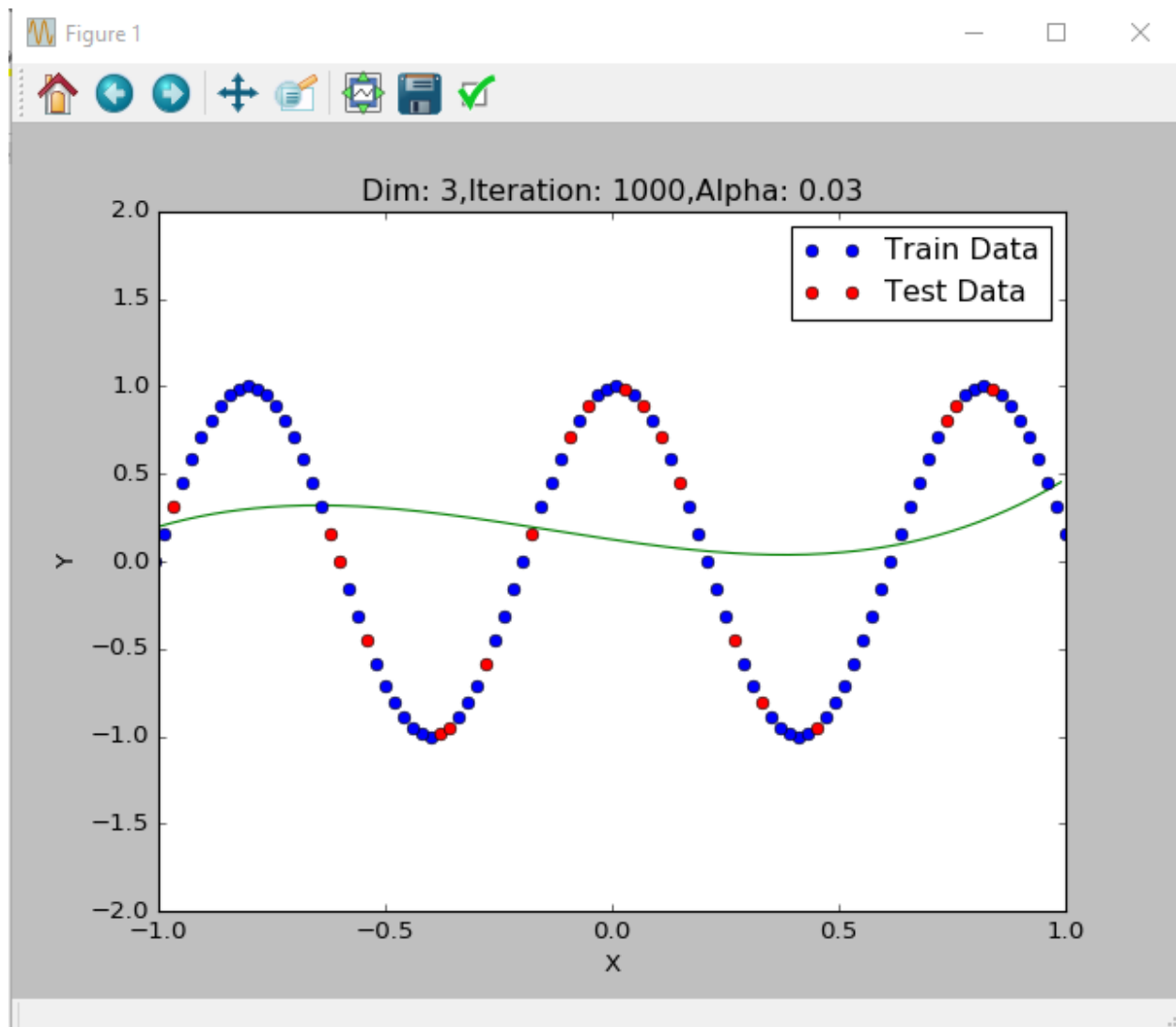
درجه ی ۳ – گرادیان نزولی در ۱۰۰ گام – میزان تابع خطا در طول ۱۰۰ گام





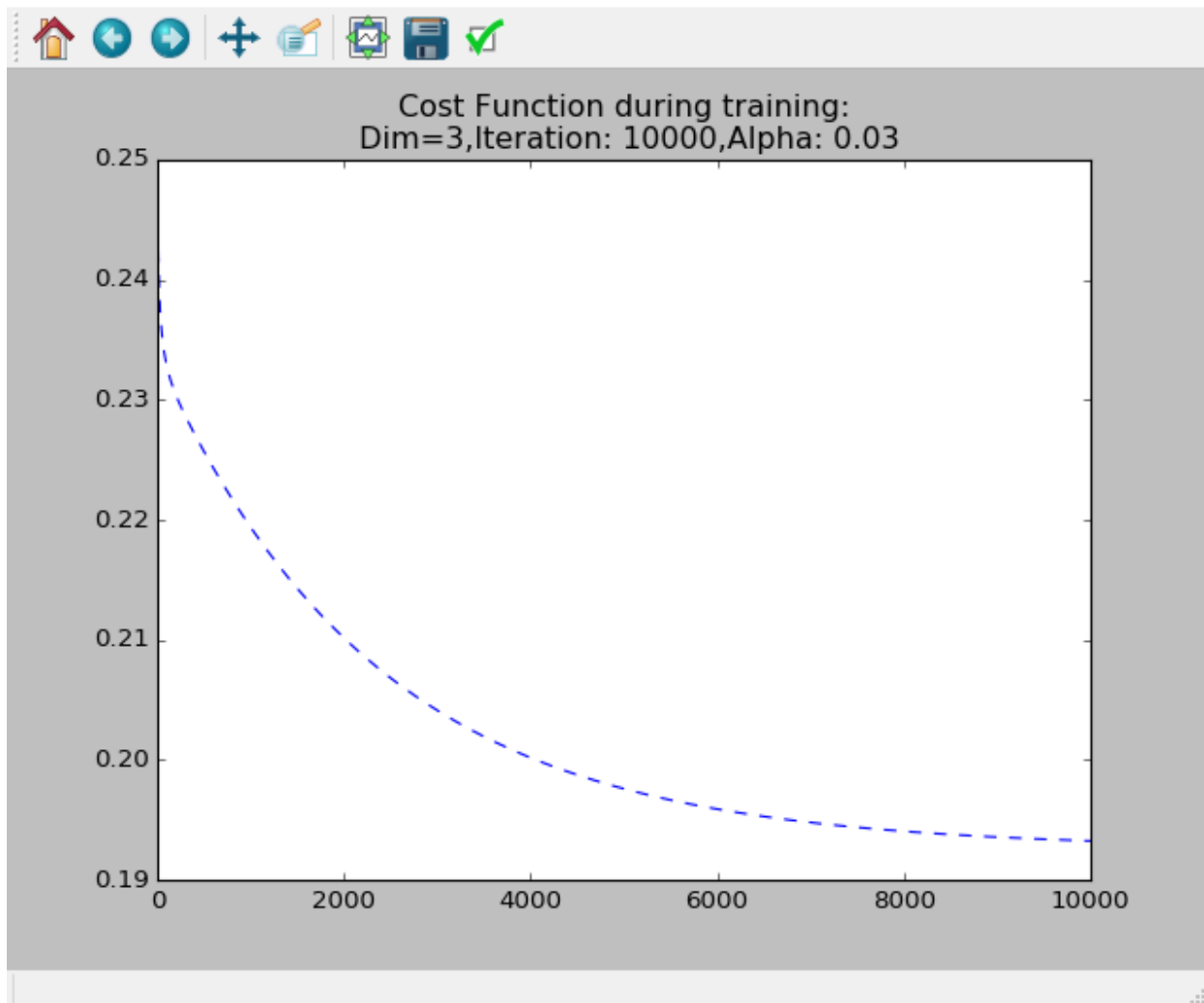
درجه ی ۳ - گرادیان نزولی در ۱۰۰۰ گام - میزان تابع خطا در طول ۱۰۰۰ گام آموزش:

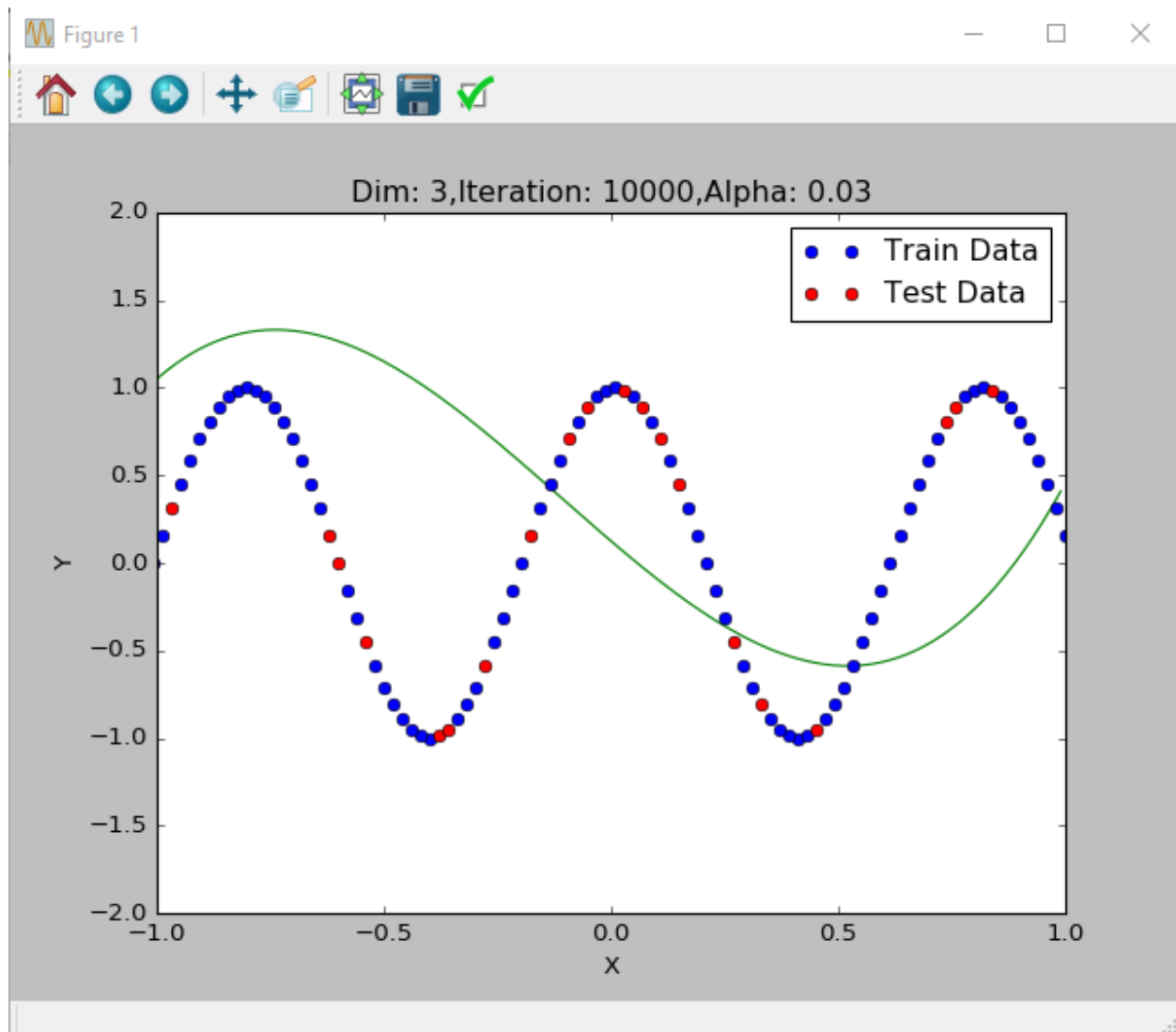




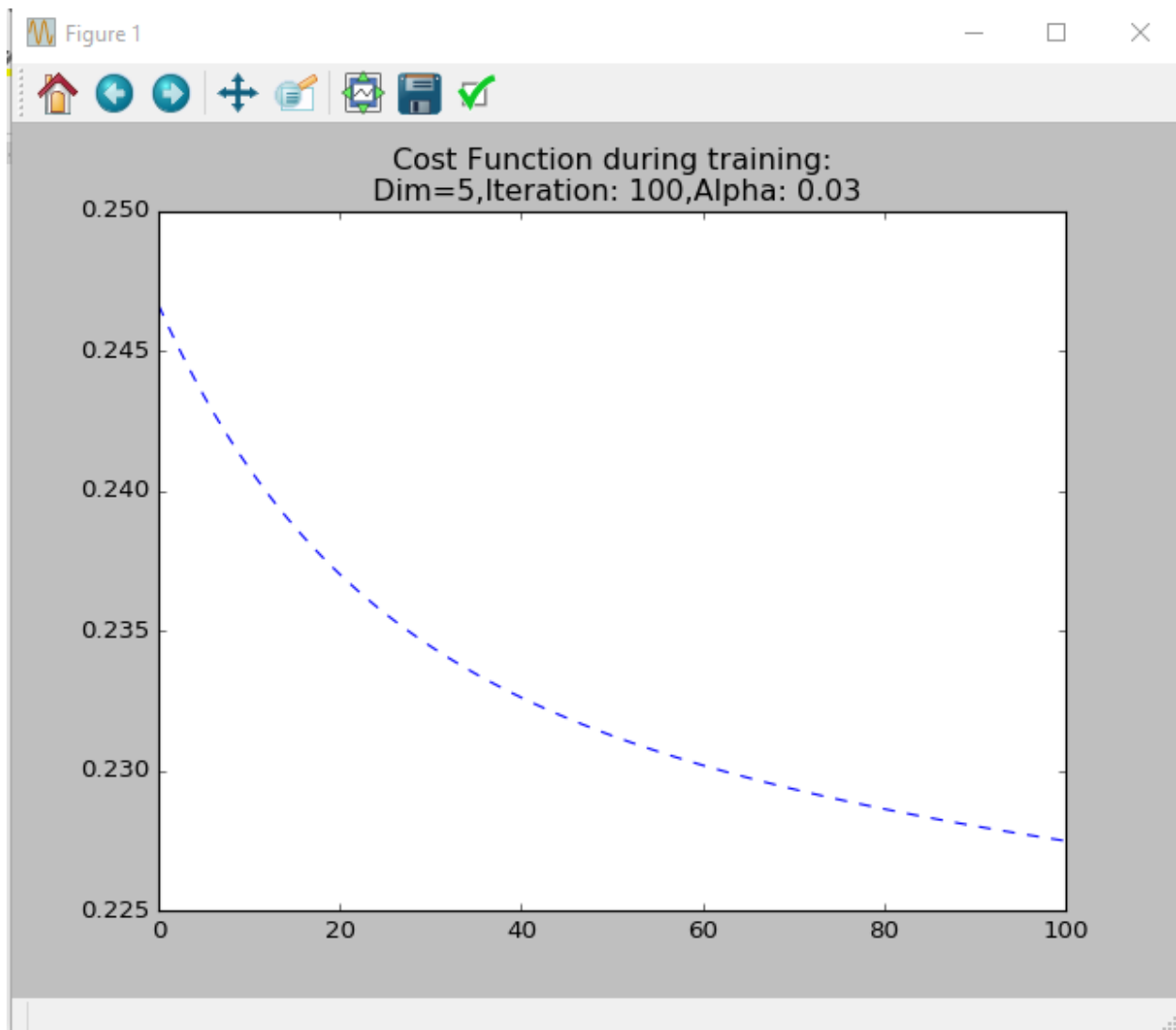
درجه‌ی ۳ – گرادیان نزولی در ۱۰۰۰۰ گام – میزان تابع خطا در طول ۱۰۰۰۰ گام آموزش:

Figure 1

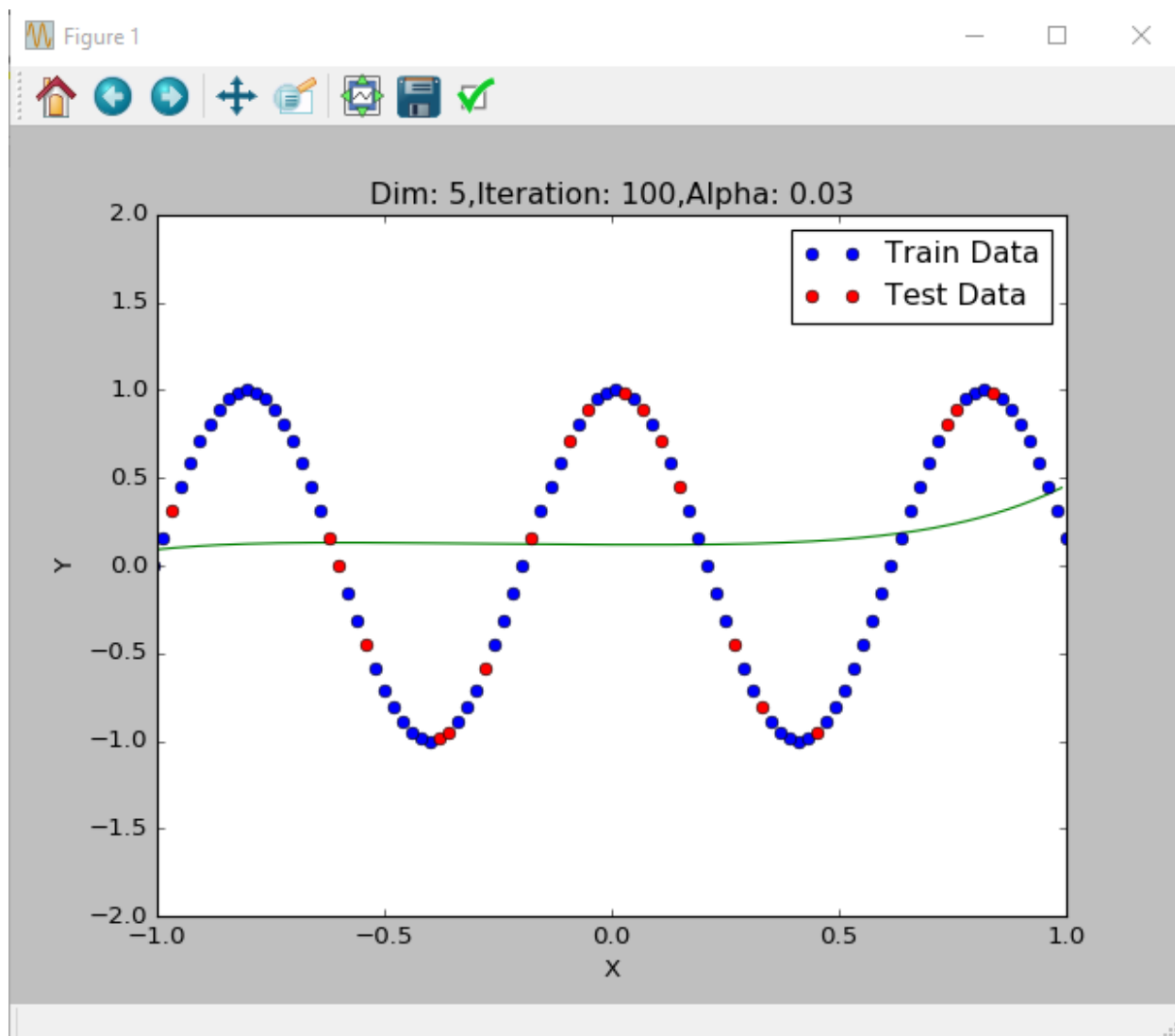




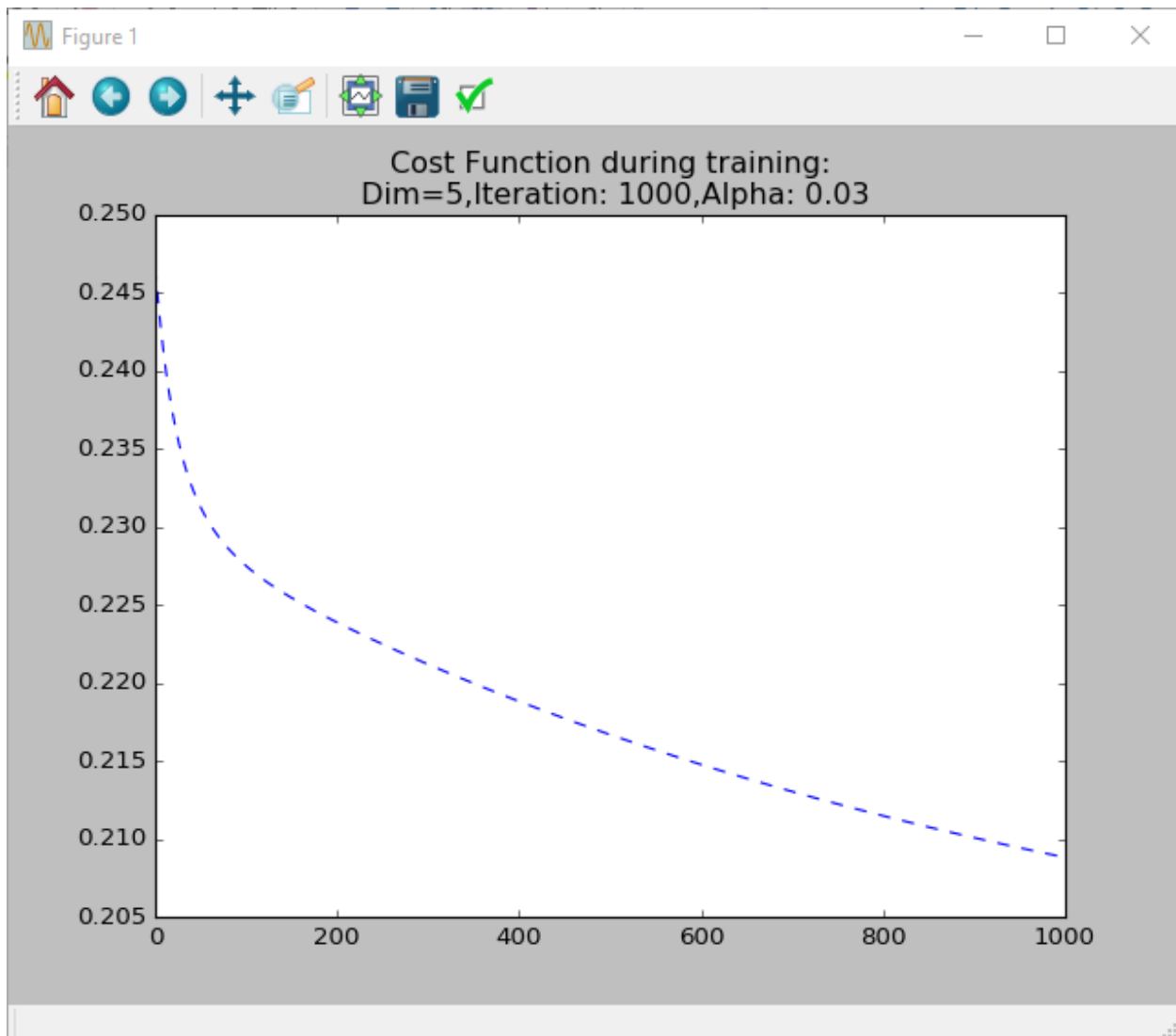
درجه‌ی ۵ - گرادیان نزولی در ۱۰۰ گام - میزان تابع خطا در طول ۱۰۰ گام آموزش:

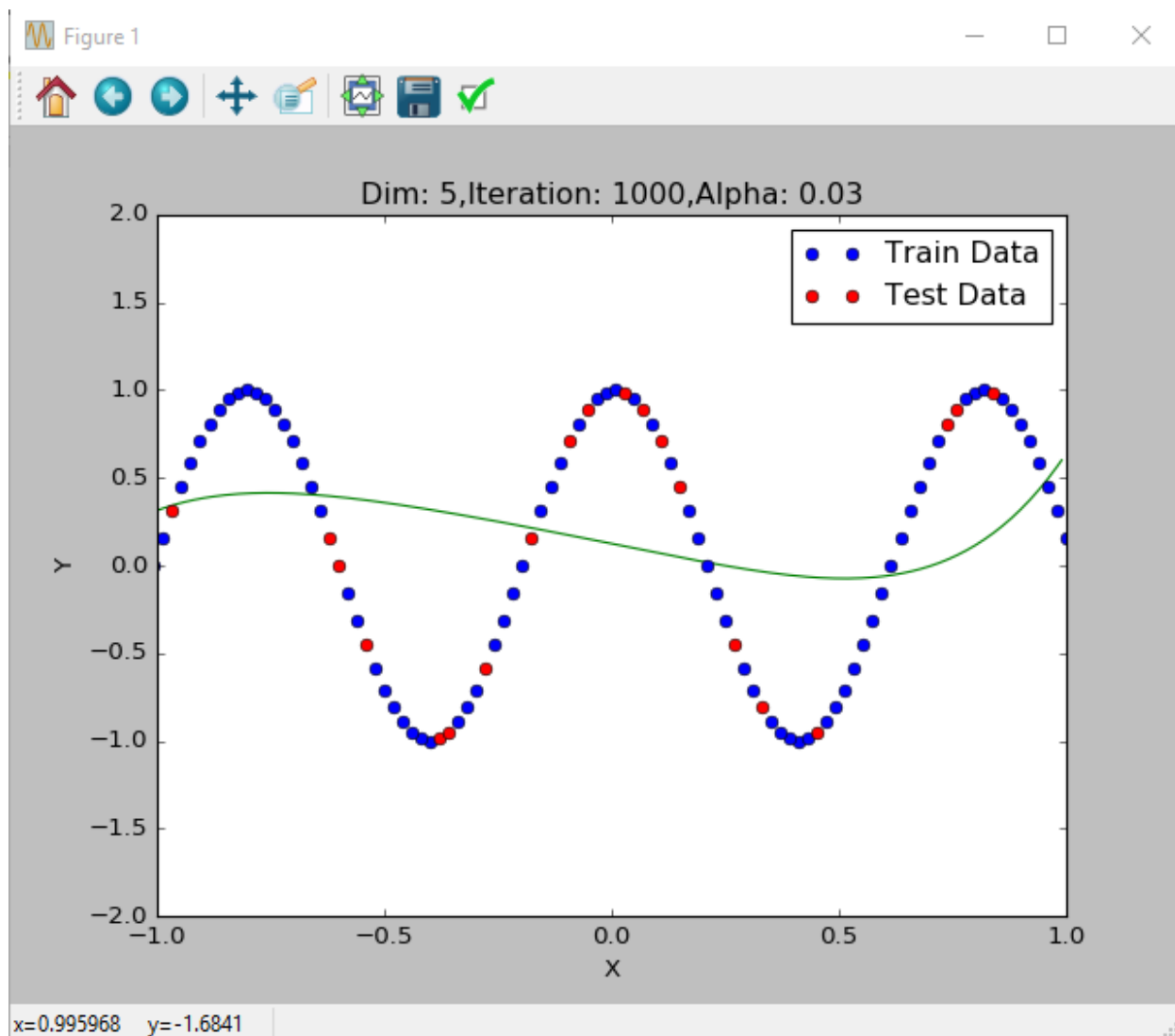




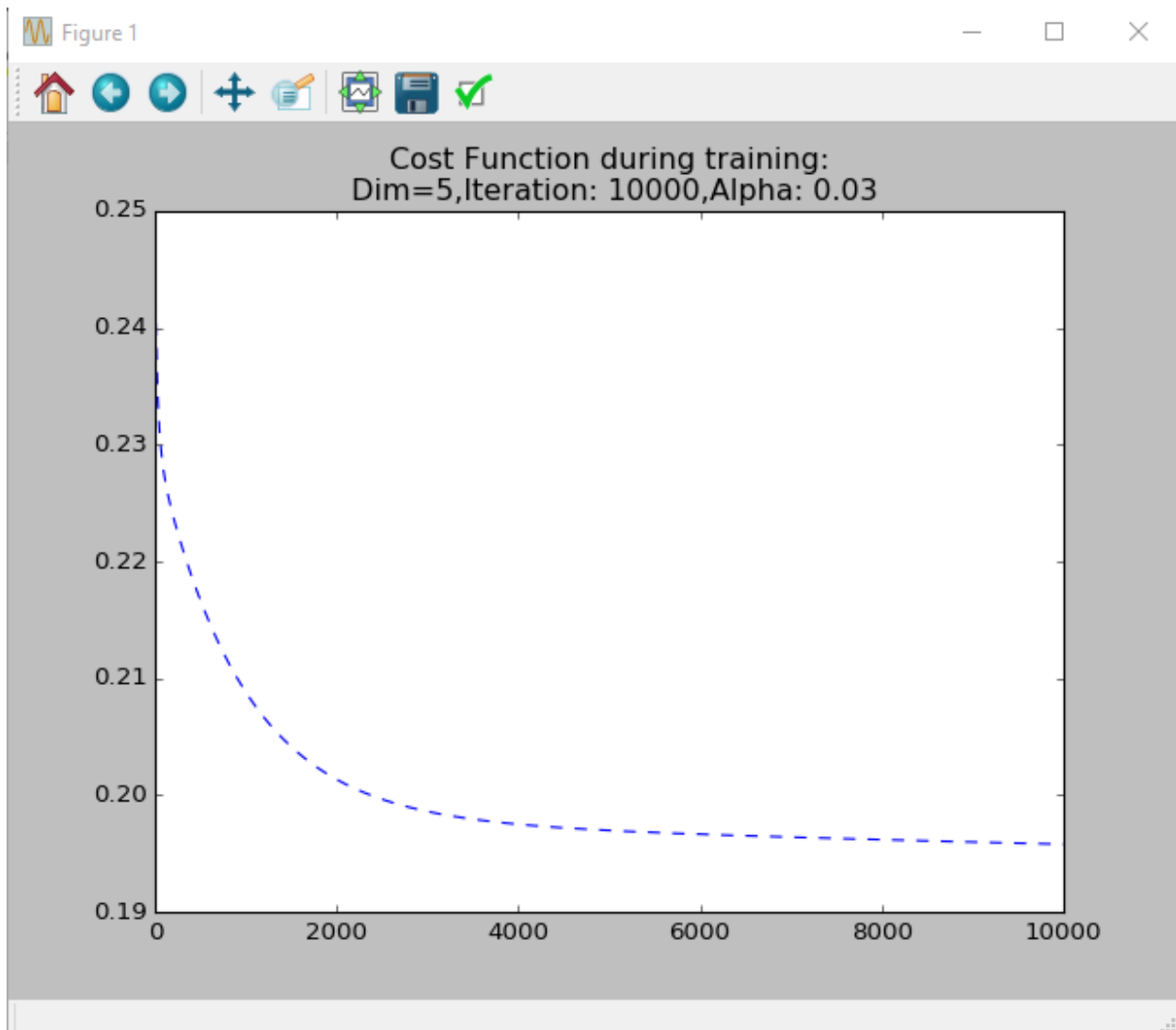


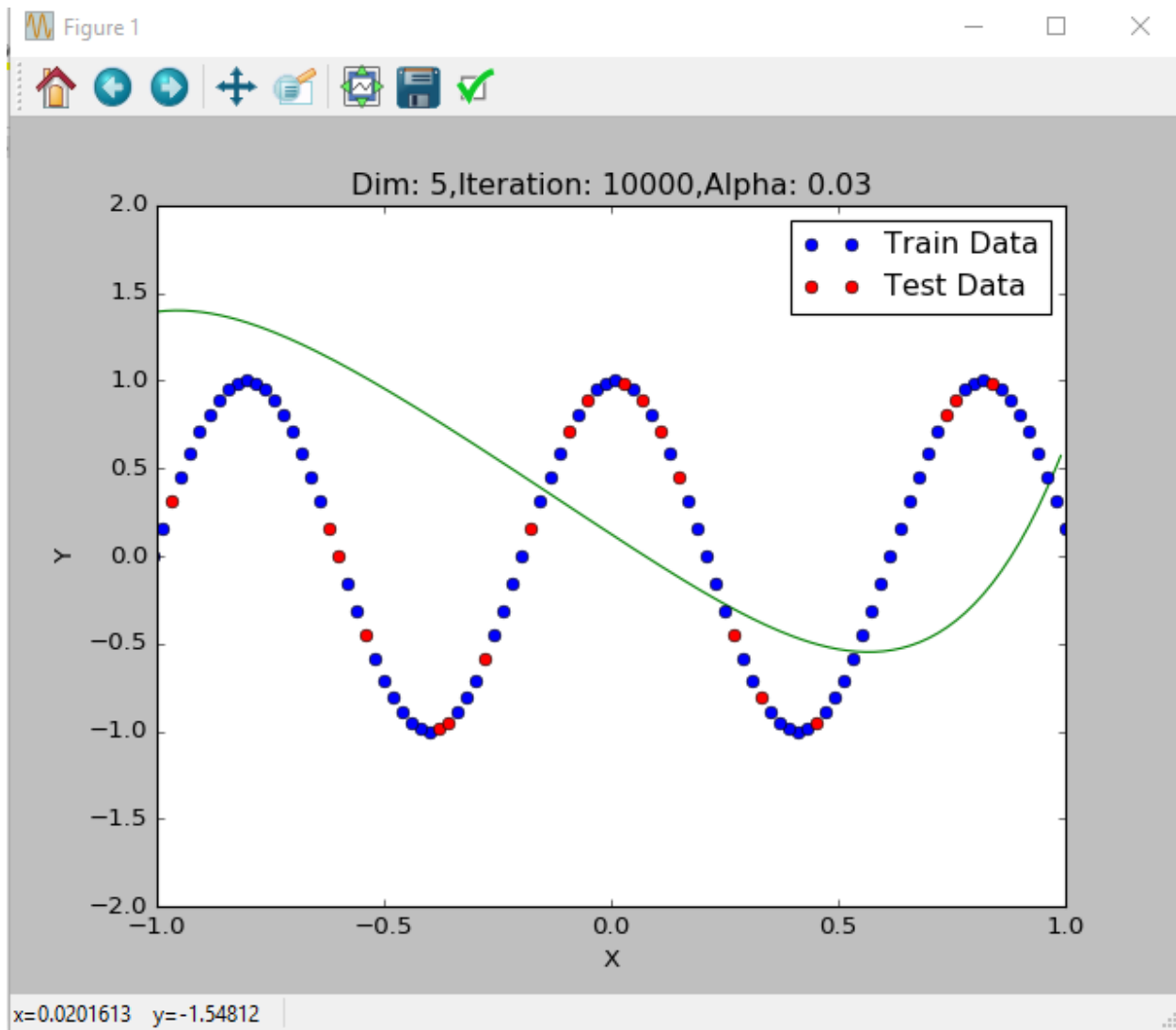
درجه ی ۵ - گرادیان نزولی در ۱۰۰۰ گام - میزان تابع خطا در طول ۱۰۰۰ گام آموزش:



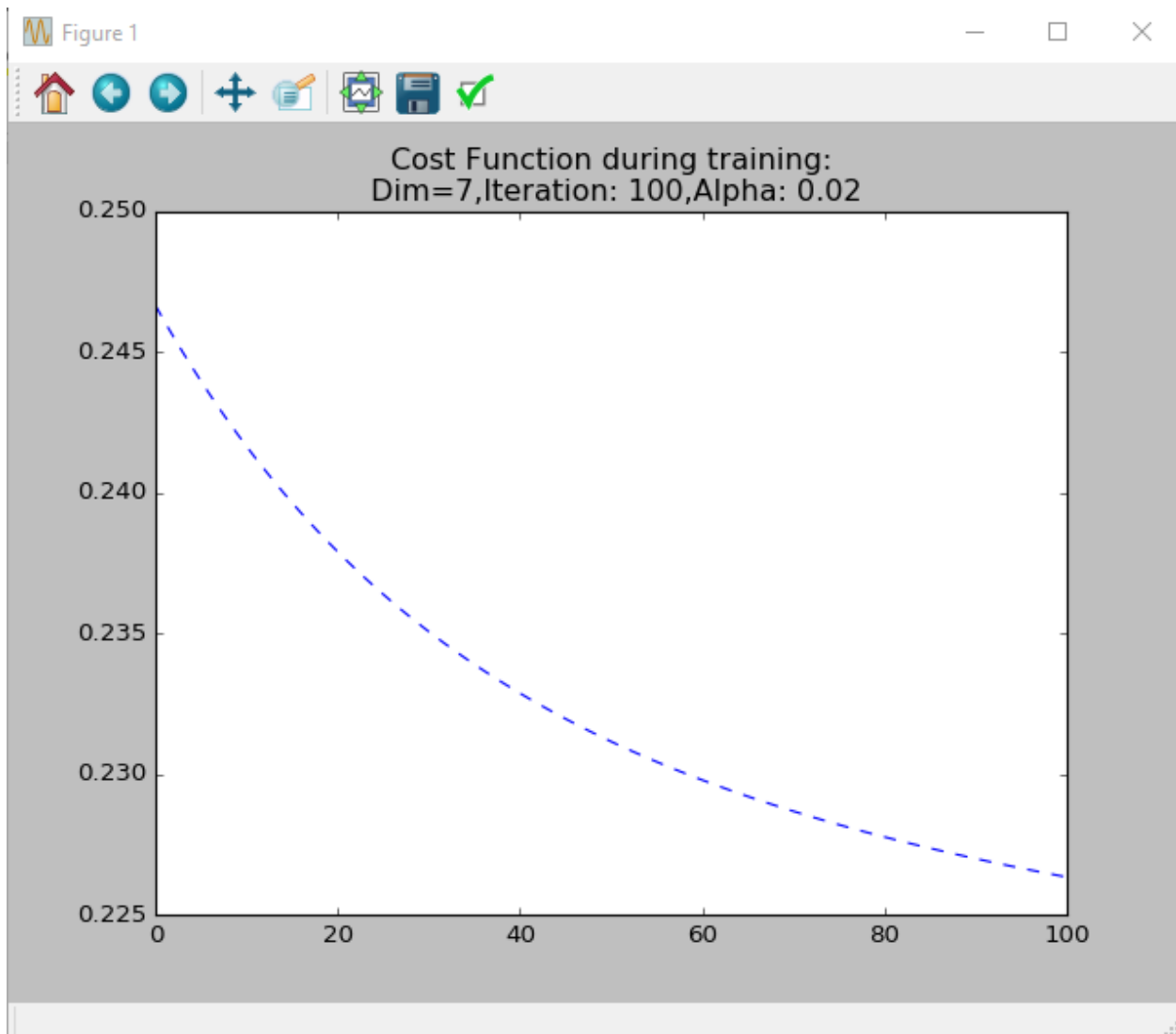


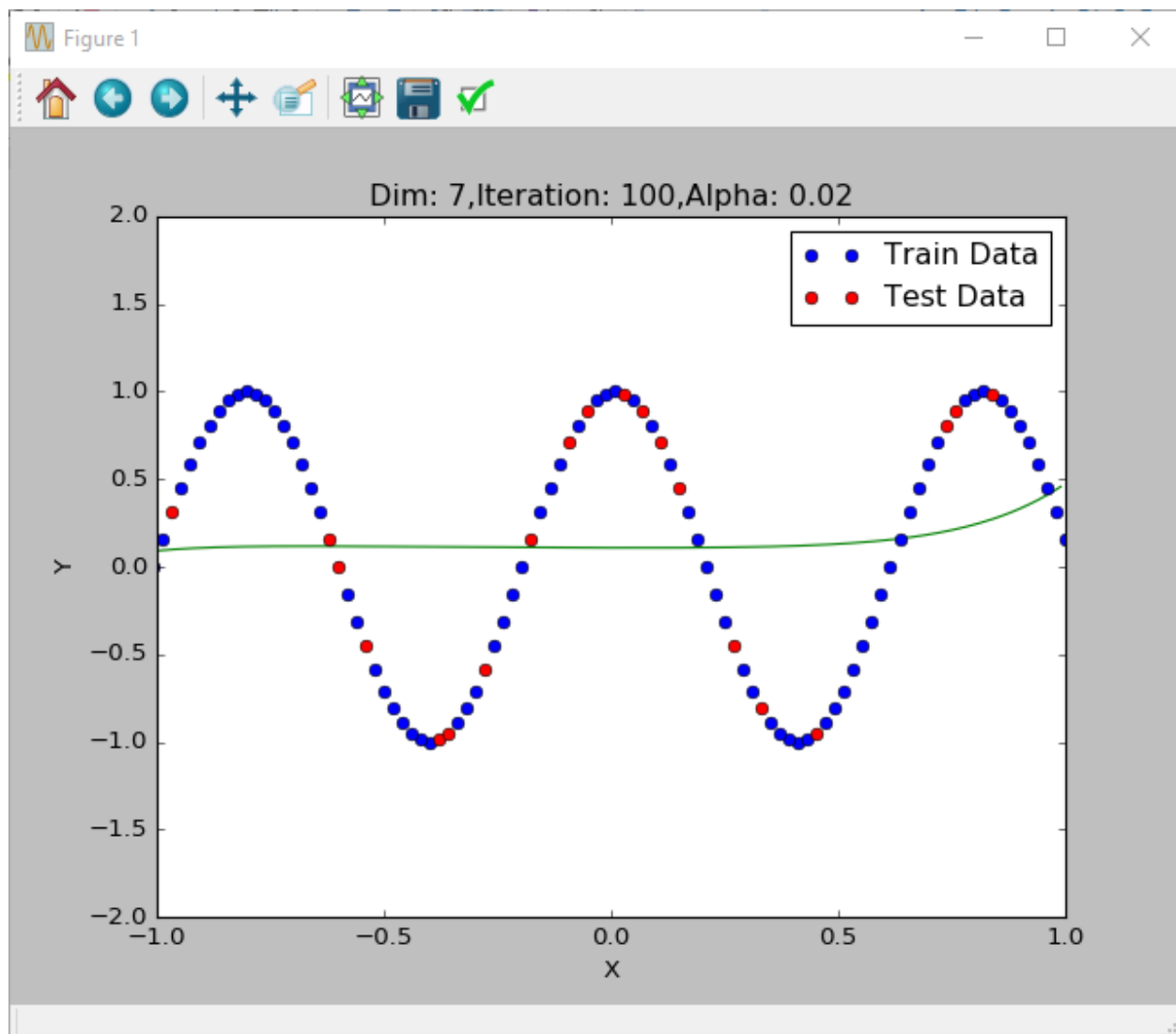
درجه ی ۵ – گرادیان نزولی در ۱۰۰۰۰ گام – میزان تابع خطا در طول ۱۰۰۰۰ گام آموزش:



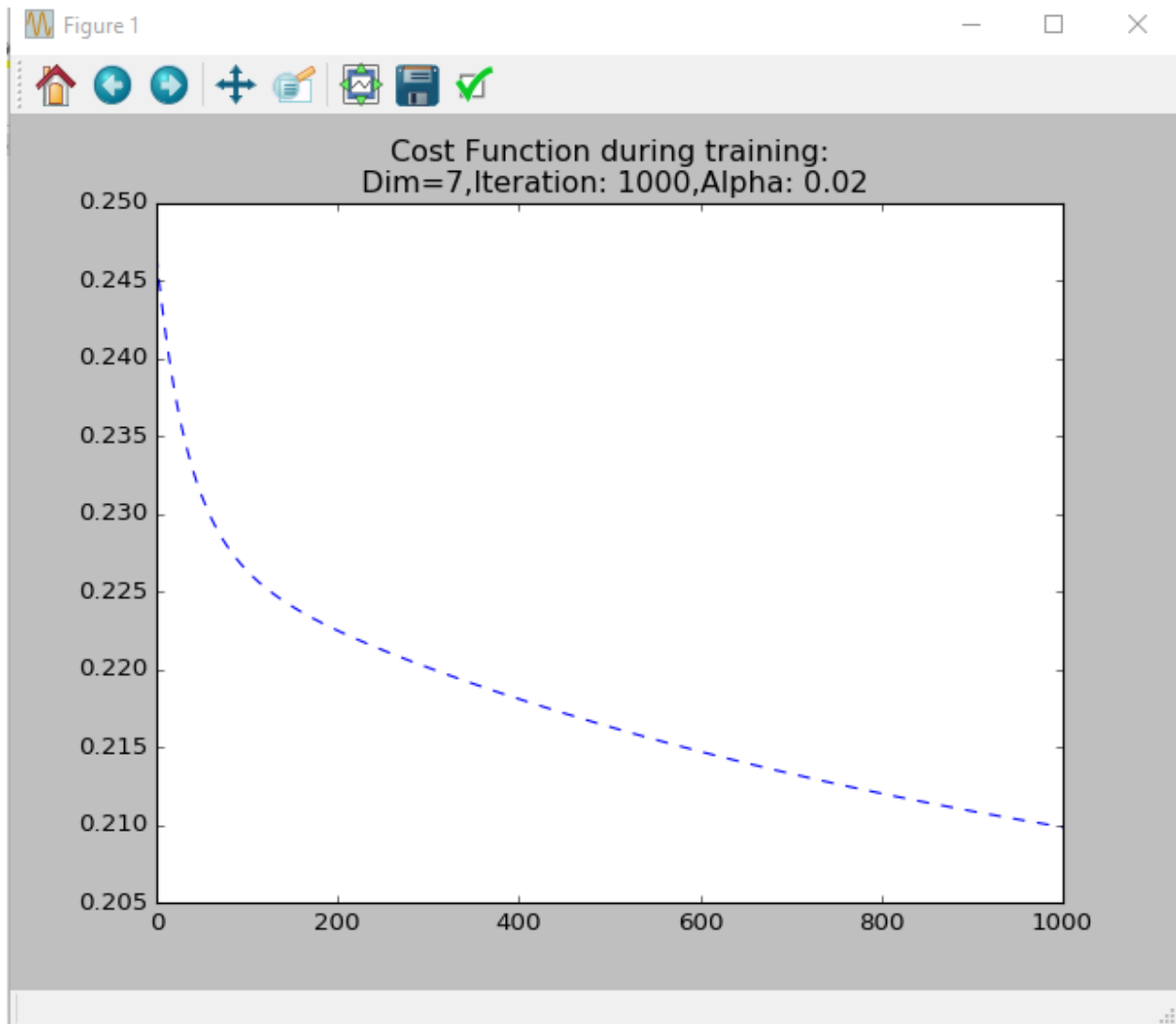


درجه‌ی ۷ – گرادیان نزولی در ۱۰۰ گام – میزان تابع خطا در طول ۱۰۰ گام آموزش:

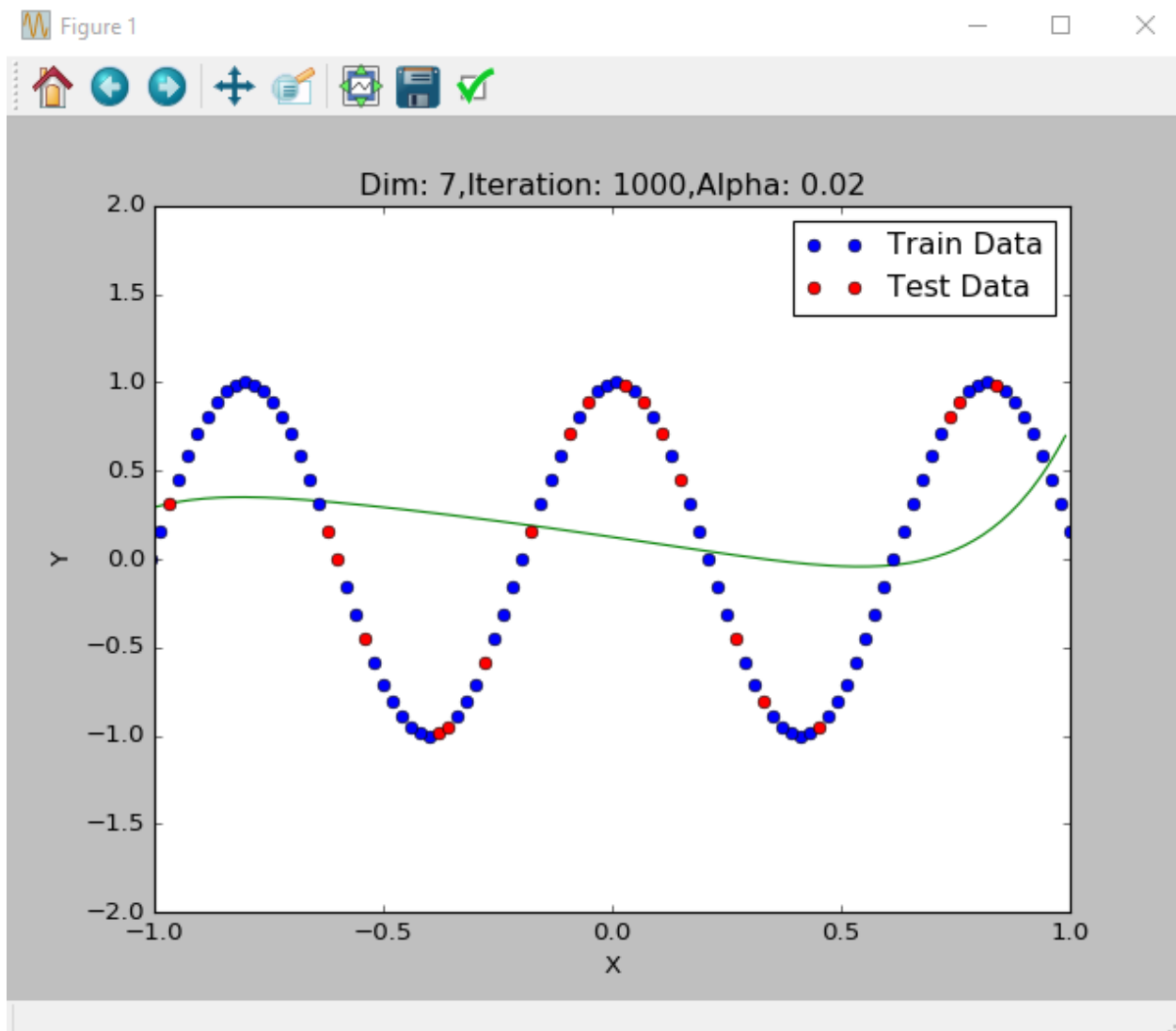




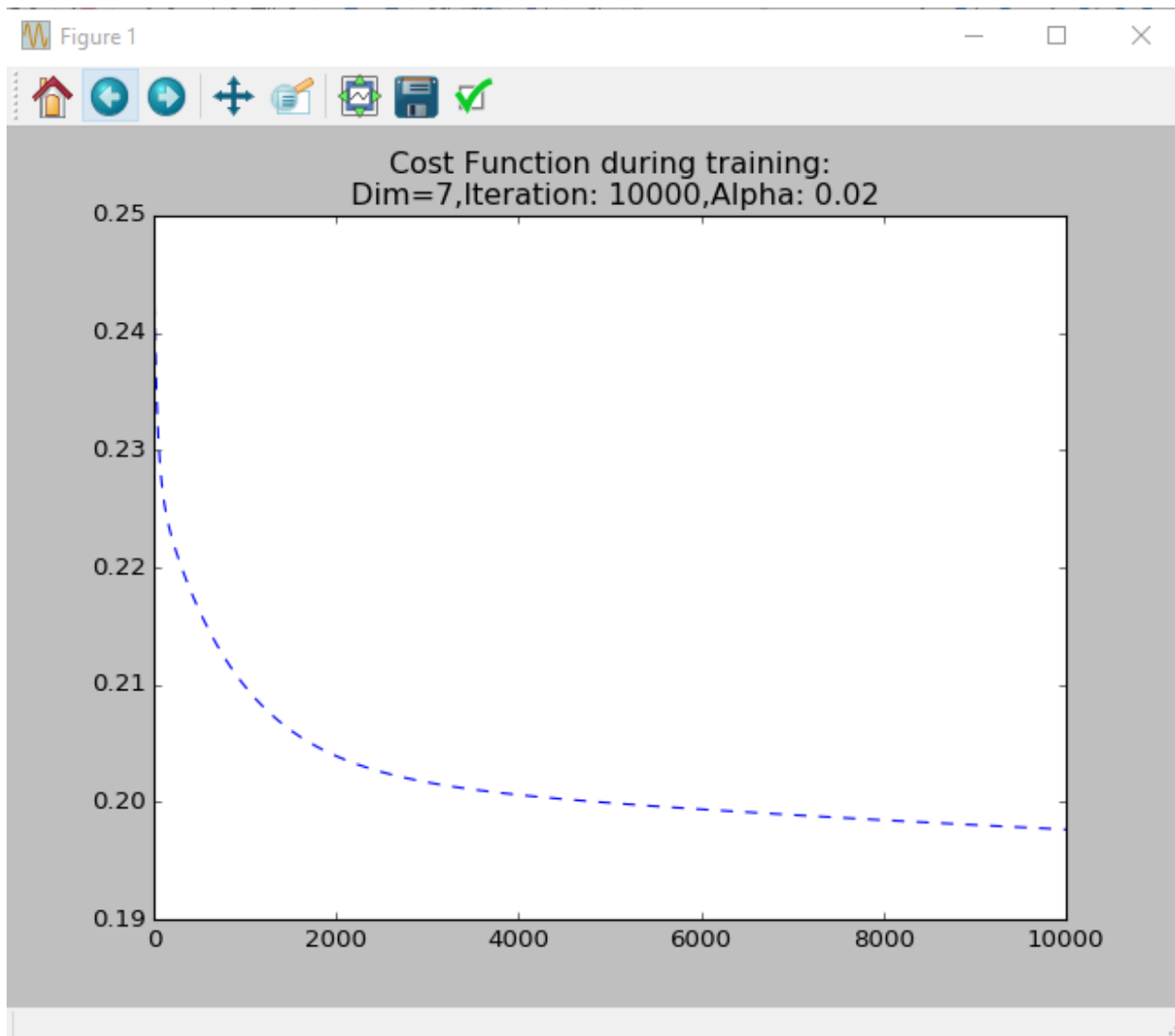
درجه‌ی ۷ - گرادیان نزولی در ۱۰۰۰ گام - میزان تابع خطا در طول ۱۰۰۰ گام آموزش:

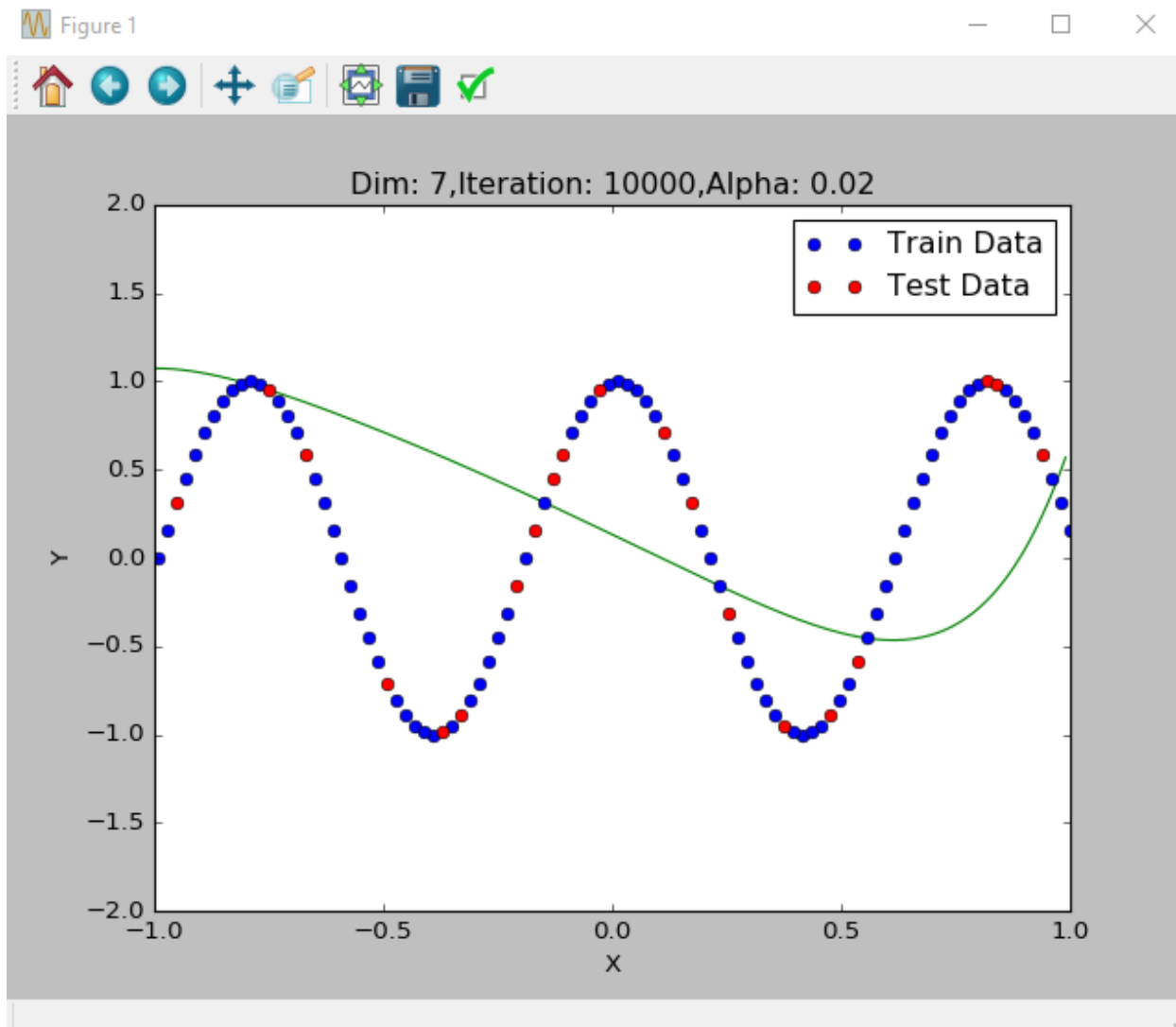




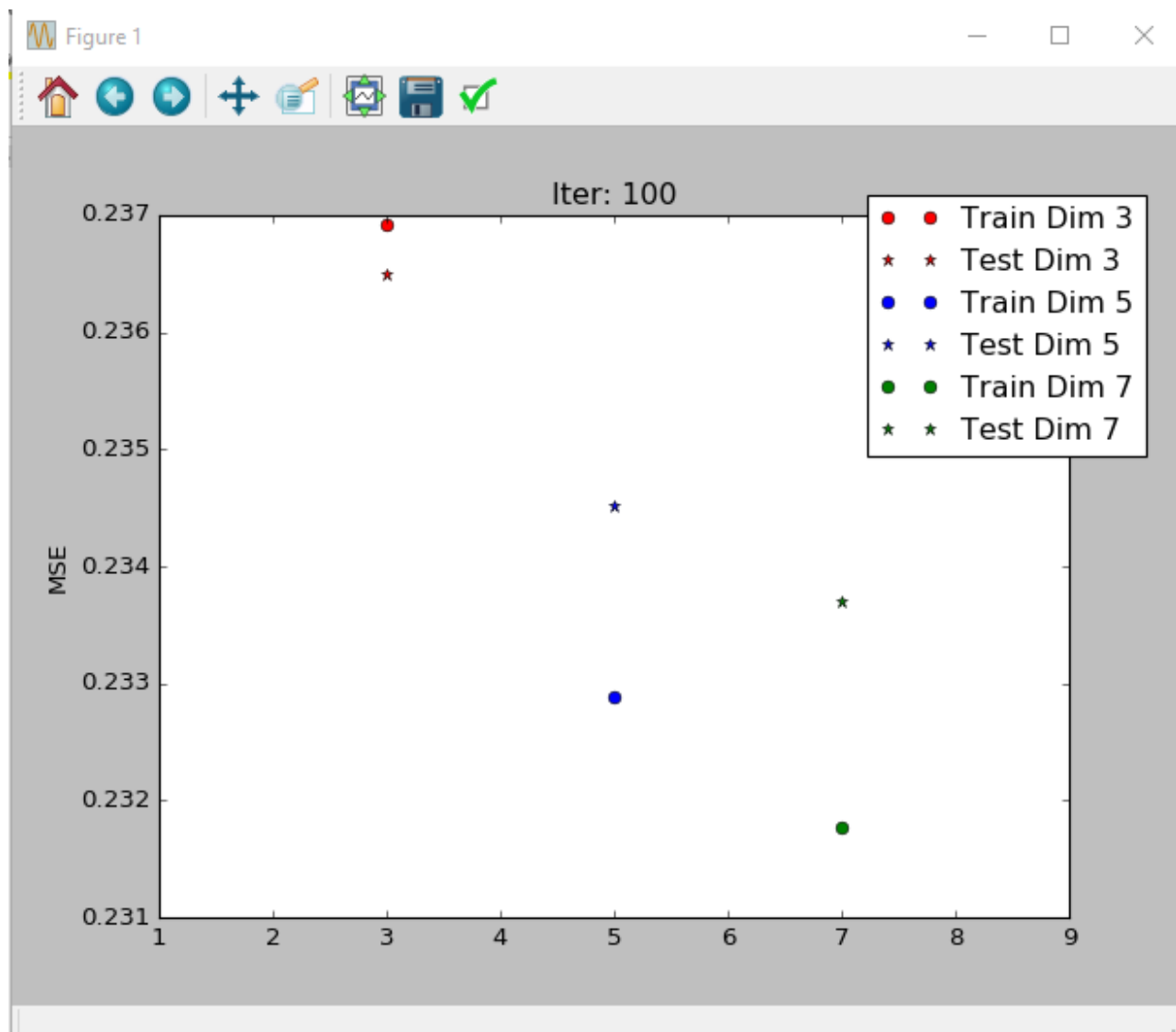


درجه‌ی ۷ - گرادیان نزولی در ۱۰۰۰۰ گام - میزان تابع خطا در طول ۱۰۰۰۰ گام آموزش:

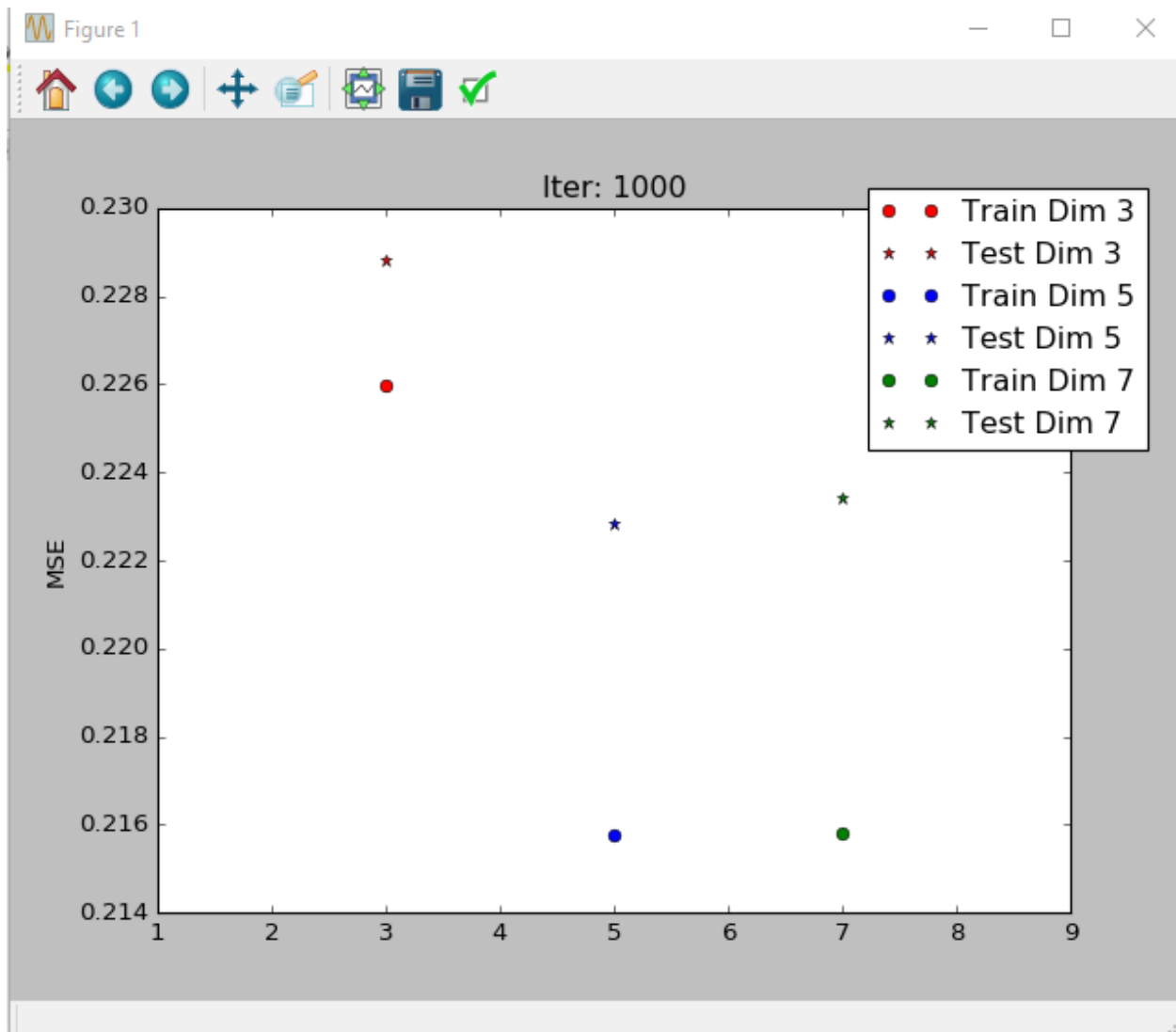




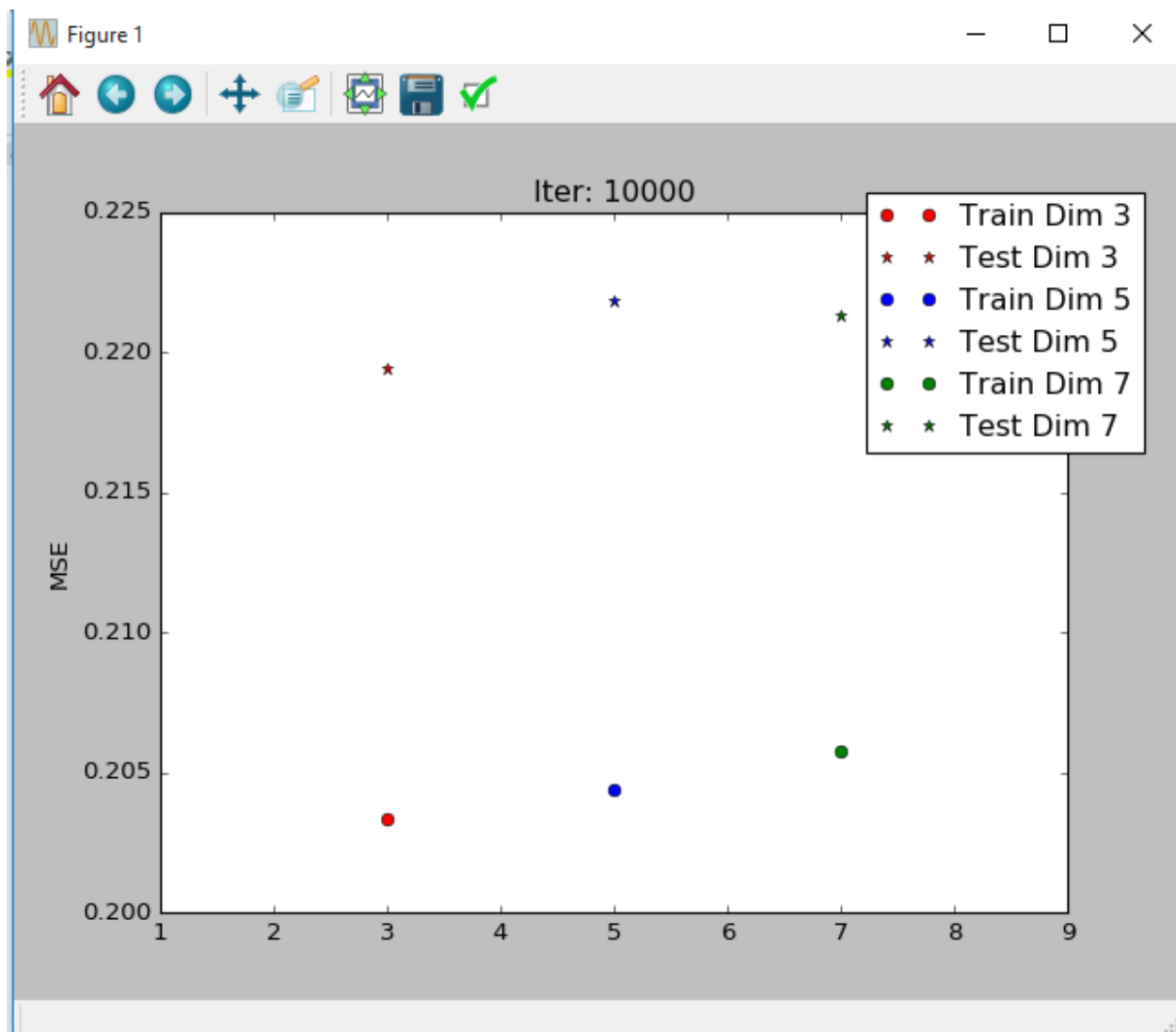
میزان خطای مجموعه‌ی آموزش و تست برای درجه‌های مختلف ۳، ۵ و ۷ با استفاده از گرادیان نزولی با ۱۰۰ گام:



میزان خطای مجموعه‌ی آموزش و تست برای درجه‌های مختلف ۳، ۵ و ۷ با استفاده از گرادینان نزولی با ۱۰۰۰ گام:

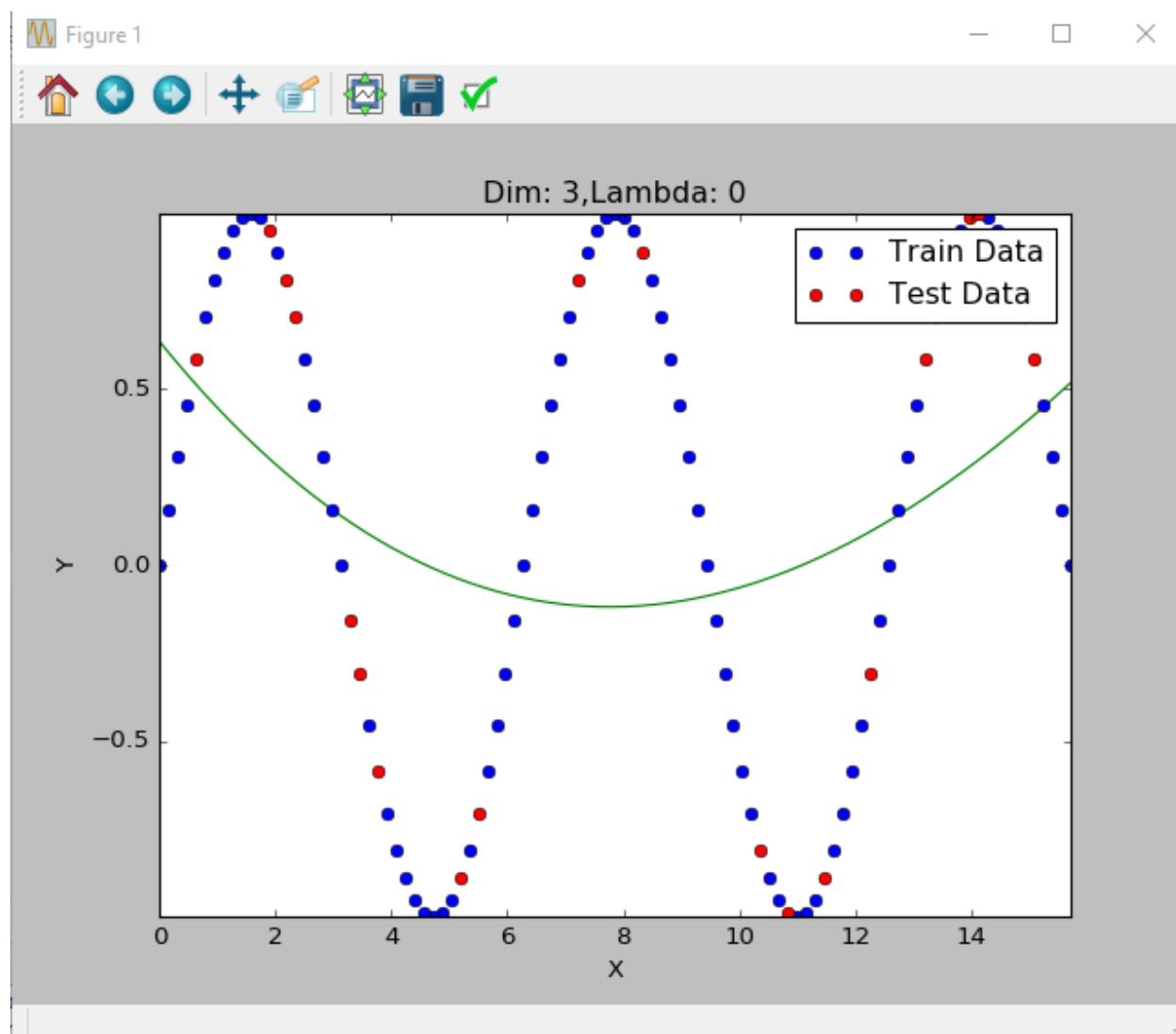


میزان خطای مجموعه‌ی آموزش و تست برای درجه‌های مختلف ۳، ۵ و ۷ با استفاده از گرادینان نزولی با ۱۰۰۰۰ گام:

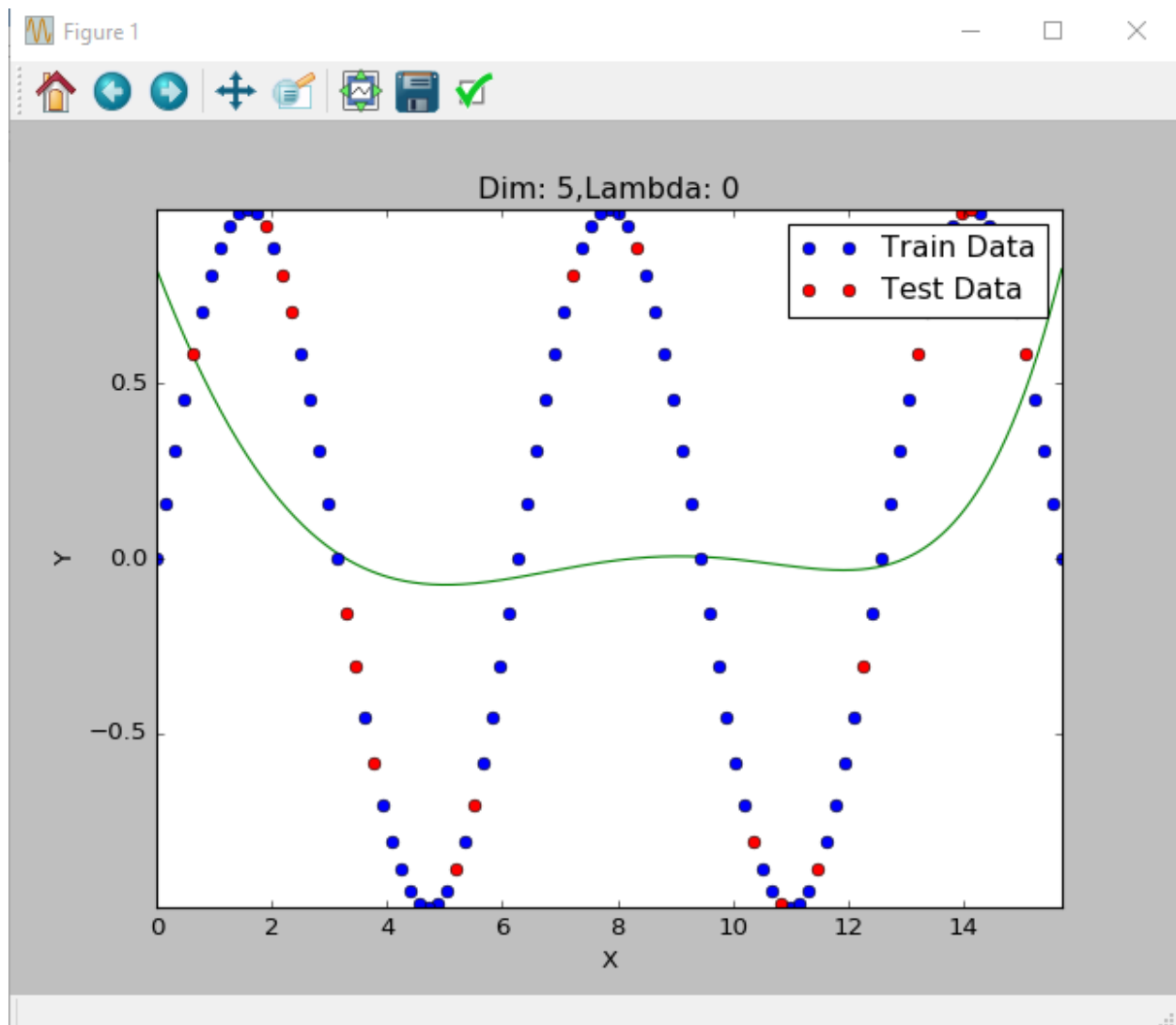


ب) کدهای قسمت ب در فایل `b-normal equation.py` موجود است که به خوبی کامنت گذاری شده است. نتایج حاصل را در تصویرهای زیر مشاهده می کنید:

تصویر زیر منحنی پیدا شده با استفاده از معادله نرمال با لامبدا صفر و `nonlinear transform` به درجه ۳ را نشان می دهد.

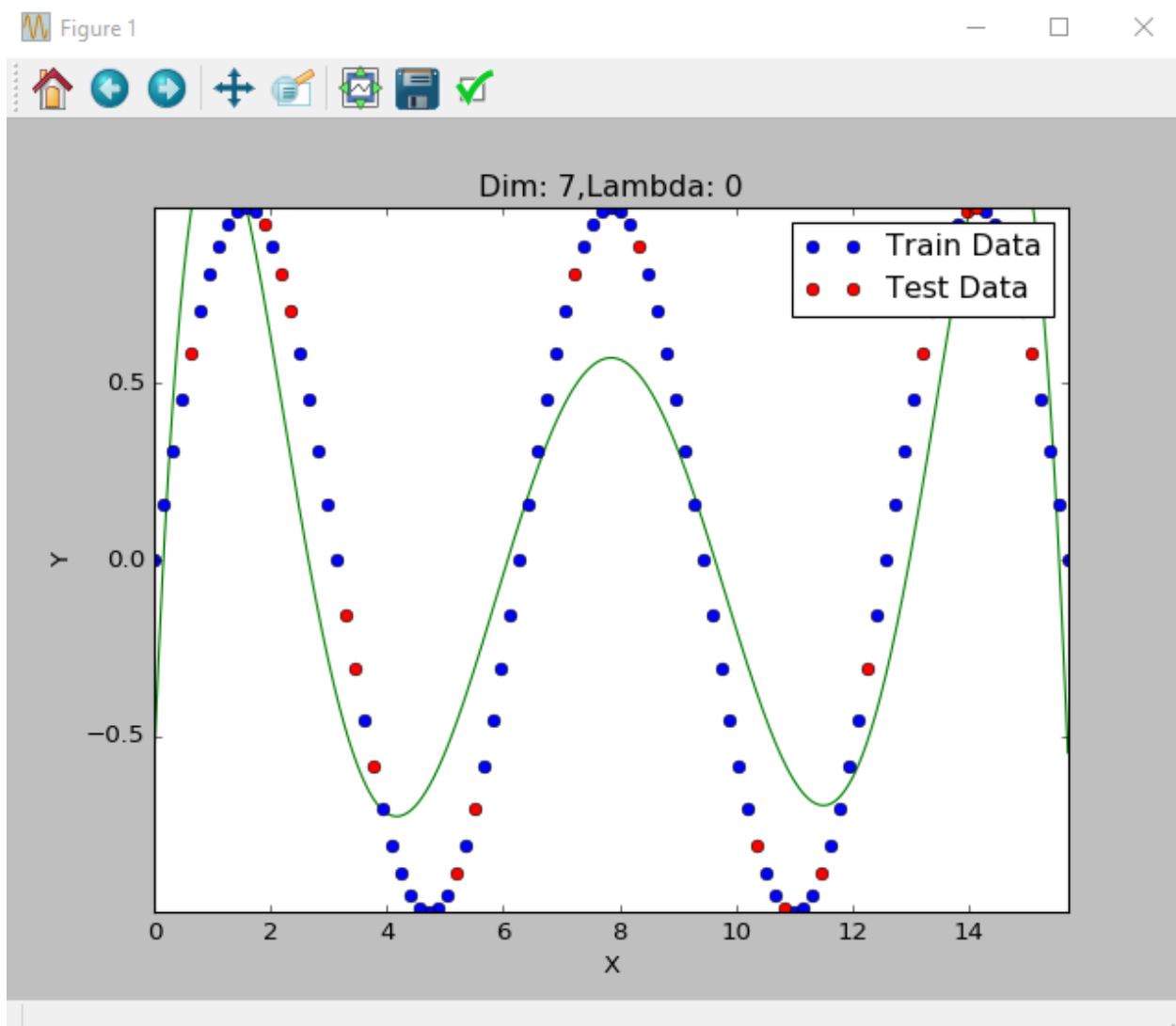


تصویر زیر منحنی پیدا شده با استفاده از معادله نرمال با لامبدا صفر و `nonlinear transform` به درجه ۵ را نشان می دهد.

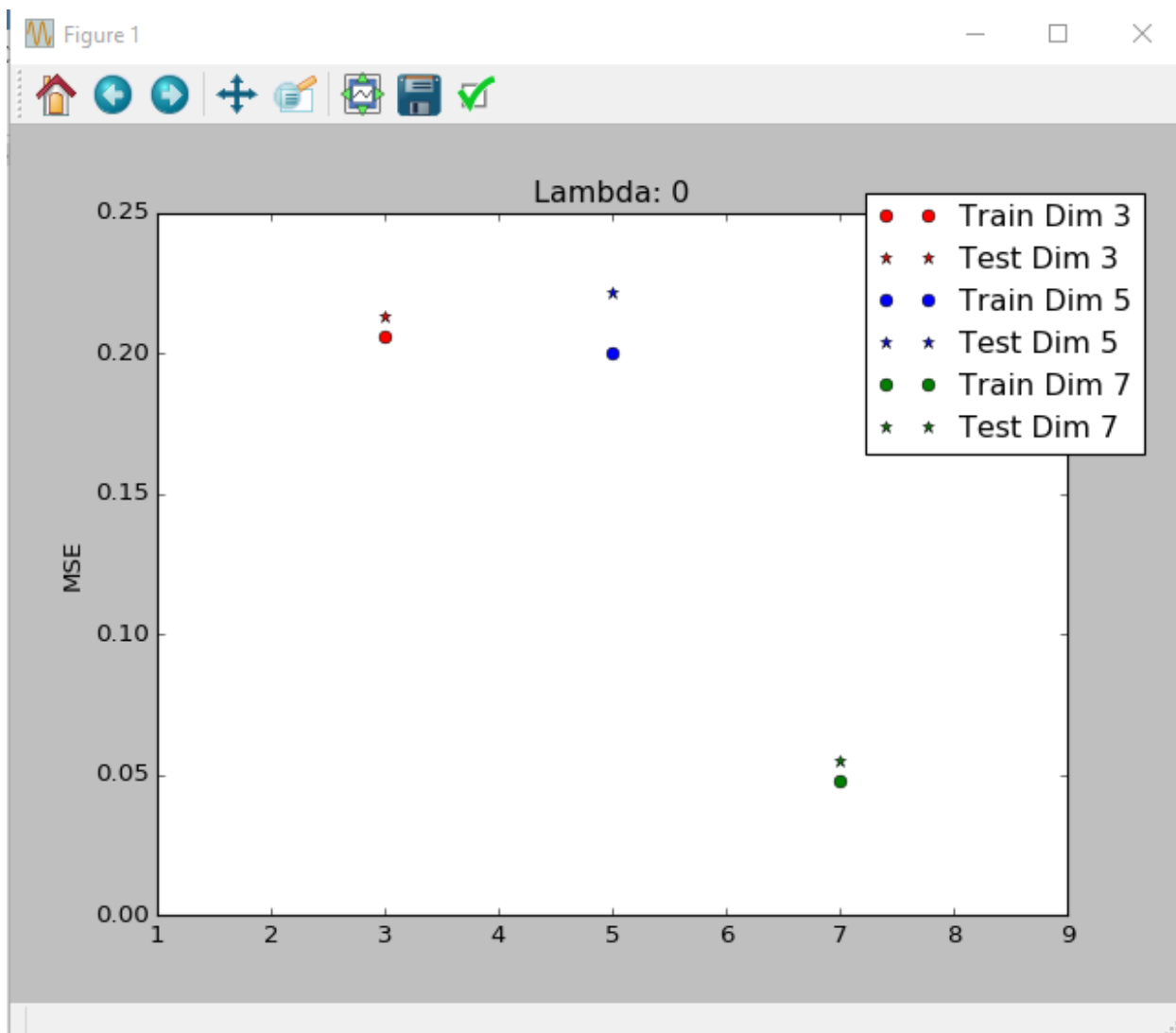


تصویر زیر منحنی پیدا شده با استفاده از معادله نرمال با لامبدا صفر و nonlinear transform به درجه ی ۷ را نشان می دهد.





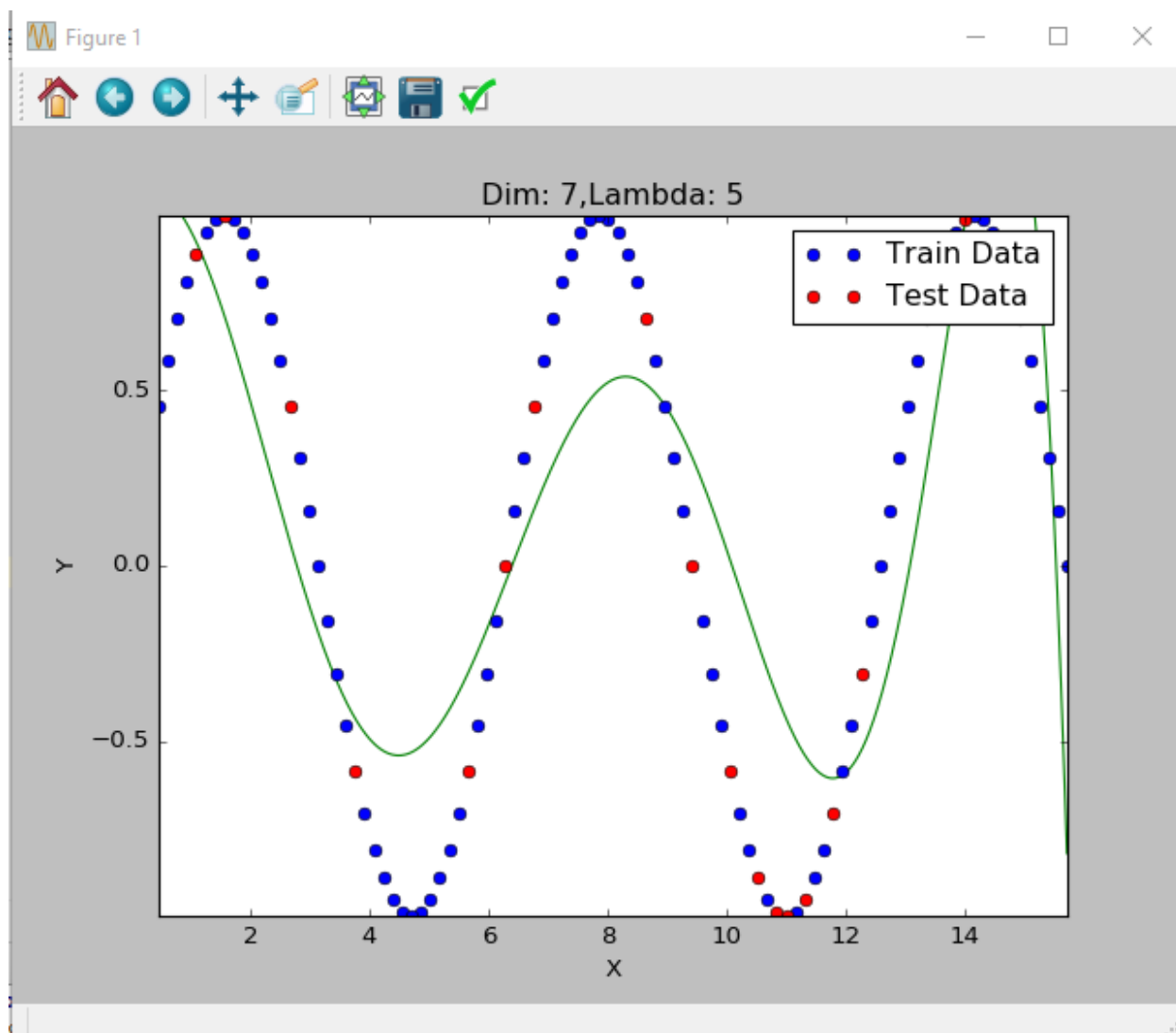
در تصویر زیر خطای مجموعه ی تست و آموزش را برای درجه های ۳ و ۵ و ۷ را مشاهده می کنید:



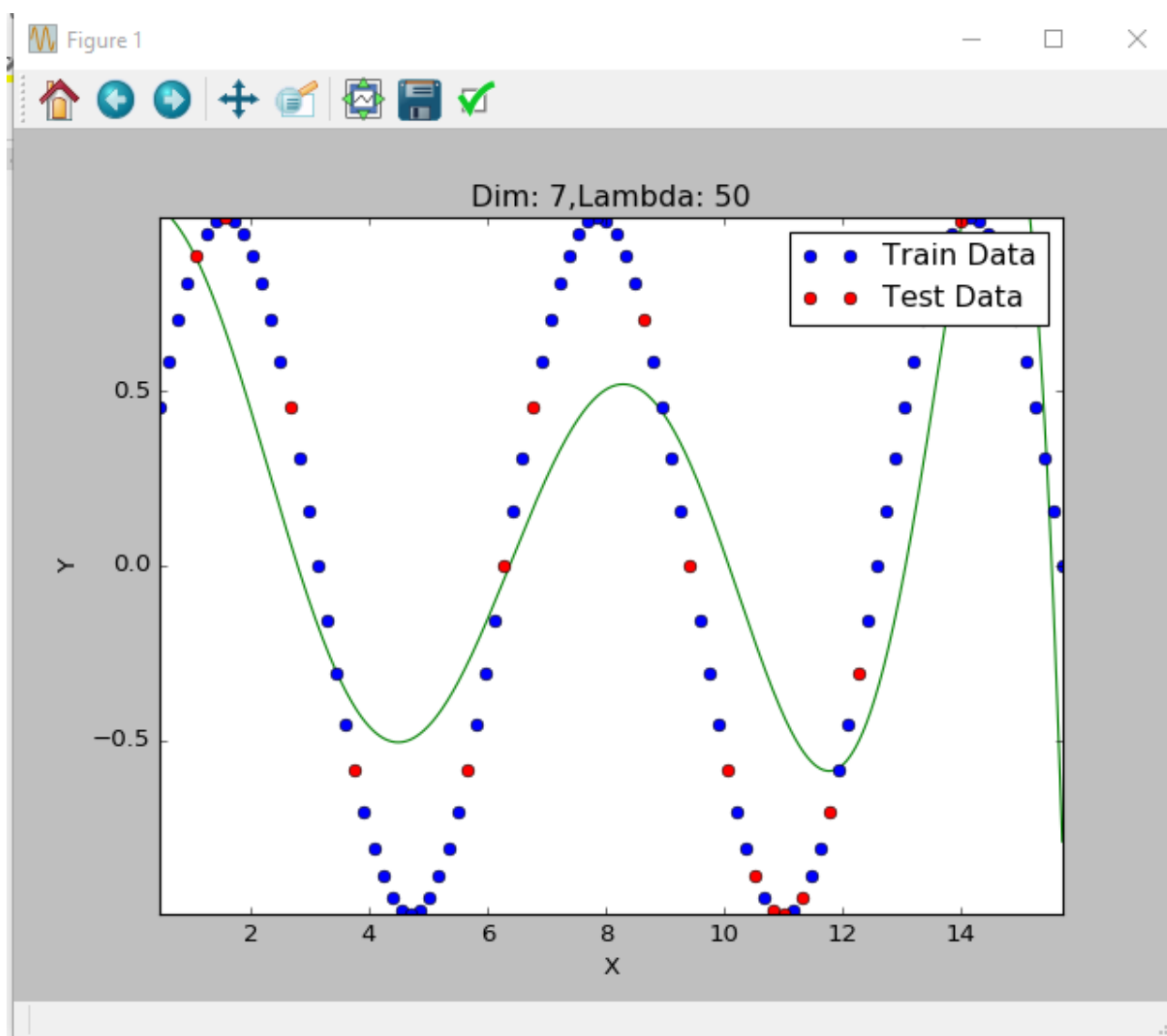
(ج)

کدهای کامنت گذاری شده‌ی این سوال در فایل `c-normal equation with lambda.py` موجود است در زیر نتایج حاصل را مشاهده می کنید:

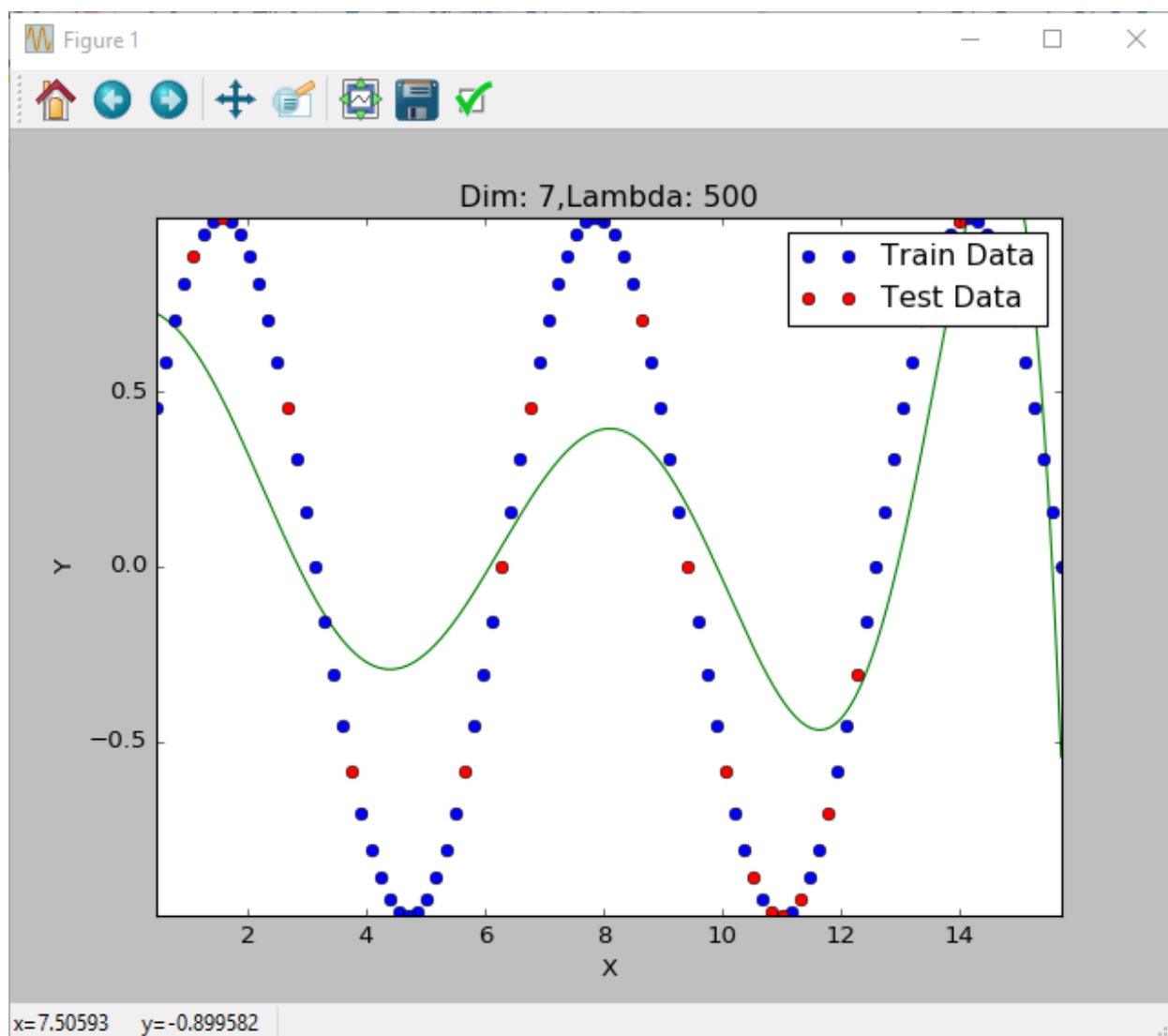
تصویر زیر منحنی پیدا شده با استفاده از معادله نرمال با لامبدا ۵ و nonlinear transform به درجه‌ی ۷ را نشان می‌دهد.



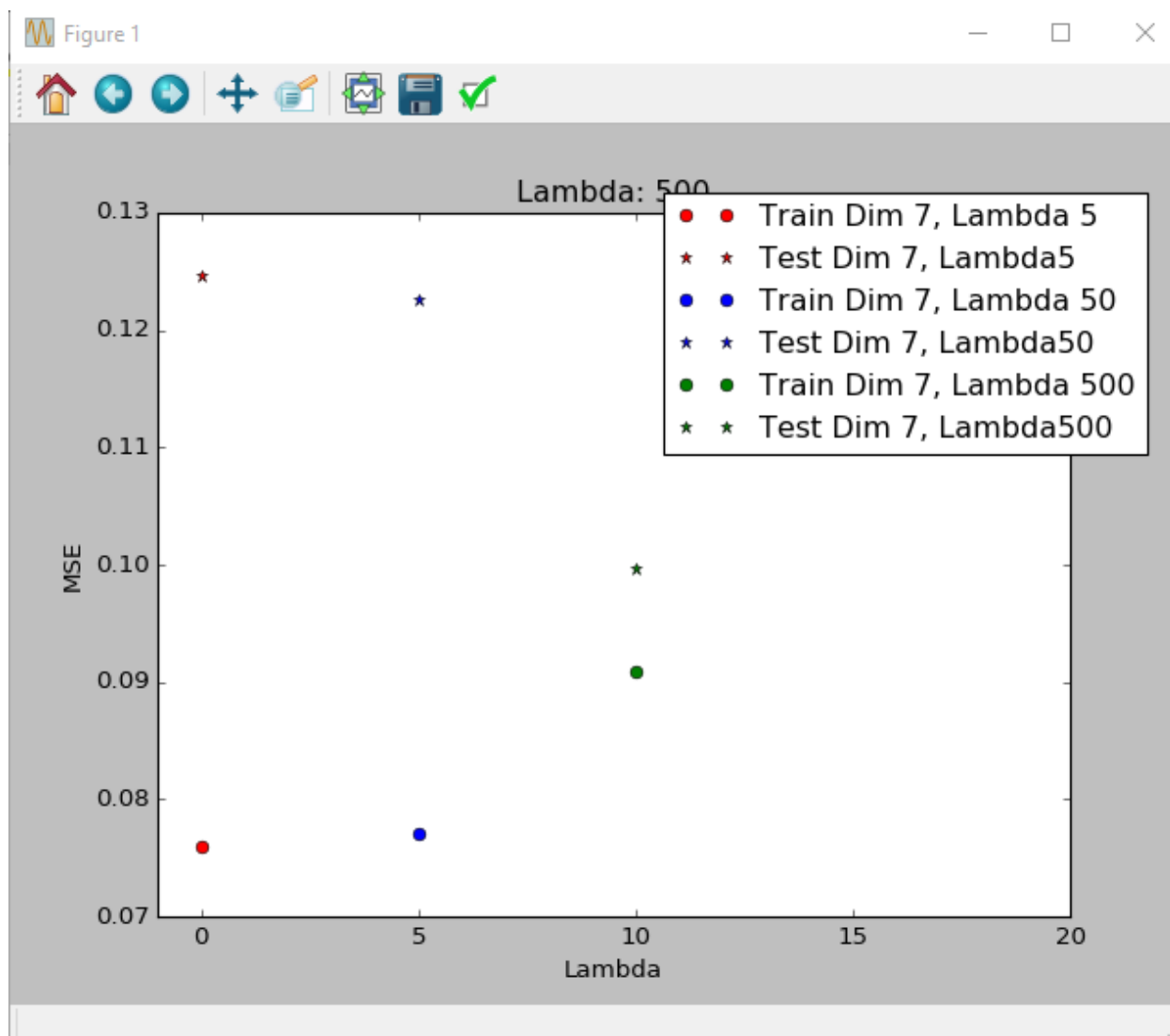
تصویر زیر منحنی پیدا شده با استفاده از معادله نرمال با لامبدا ۵۰ و nonlinear transform به درجه ۷ را نشان می‌دهد.



تصویر زیر منحنی پیدا شده با استفاده از معادله نرمال با لامبدا ۵۰۰ و nonlinear transform به درجه‌ی ۷ را نشان می‌دهد.



در تصویر زیر خطای مجموعه ی تست و آموزش را برای درجه‌های ۷ و لامبدا ۵ و ۵۰۰ را مشاهده می‌کنید:



بردار ضرایب را در تصویر زیر مشاهده می‌کنید:

[0.899918345184096, 0.1236303863411357, -0.17686276738759243, -0.07085961625524564, 0.04011189776971223, -0.006116879169635015, 0.00038291474714322526, -8.60326686415398e-06]  
 [0.8987631698431089, -0.0030818710787711033, -0.07158218883752263, -0.10084796378836713, 0.044083630696047676, -0.006379930207877106, 0.0003911418724021306, -8.69462459019549e-06]  
 [0.6265504234736387, -0.010743062492449038, -0.043277158854587286, -0.07454483475453565, 0.03263561180331699, -0.004783144393980185, 0.0002968005329916662, -6.665154442997611e-06]