

## Assignment #2 — Regression

資工三 409410114 周述君

- Execution description: steps how to execute your codes.

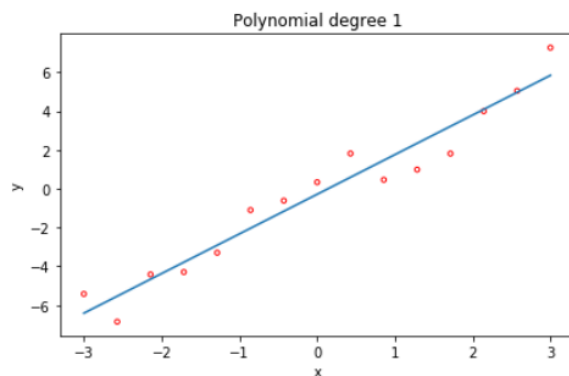
- 使用 Colab，ctrl+F9 執行。
- 將參數存放於 Parameter 可以修改 num\_points，num\_points=15。
- 參數調整好後，Generate Samples 利用 linspace 函數調整 sample 產生的區間(np.linspace(0, 1)在 0 到 1 間均勻產生 num\_points 個 sample，np.linspace(-3, 3)在-3 到 3 間均勻產生 num\_points 個 sample)。
- 產生 sample 後，plot.scatter 將生成的 sample 繪製成二維散布圖。
- 調整 regression 中的 degree 顯示 linear regression(degree=1)和 polynomial regression(degree=5, degree=10, degree=14)的 fitting plot。
- 調整 lmda 來做正則化(lmda=0, lmda=0.001/m, lmda=1/m, lmda=1000/m)。

- Experimental results: As specified in the assignment

- Perform Linear Regression. Show the fitting plot, the training error, and the five-fold cross-validation errors.

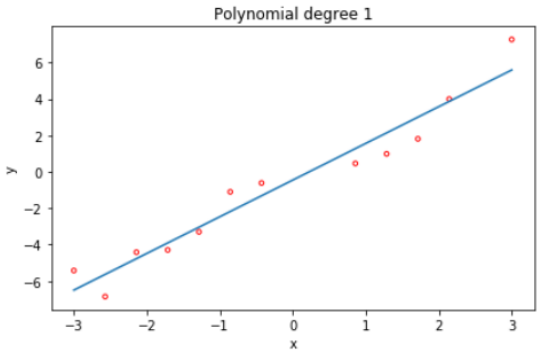
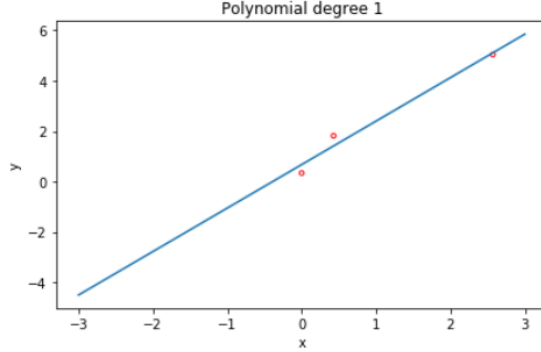
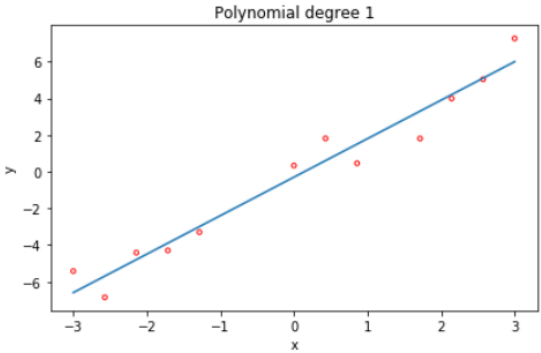
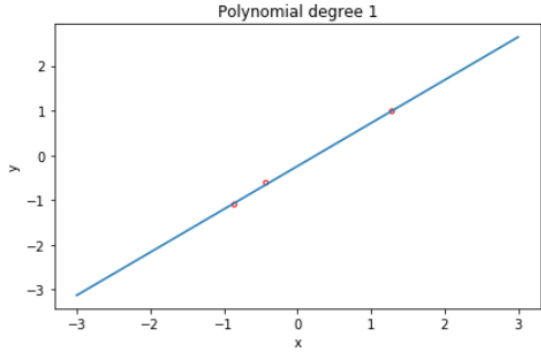
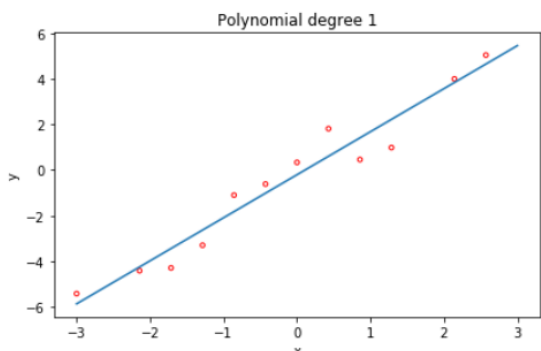
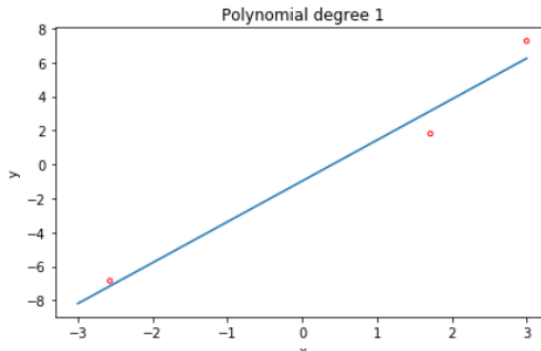
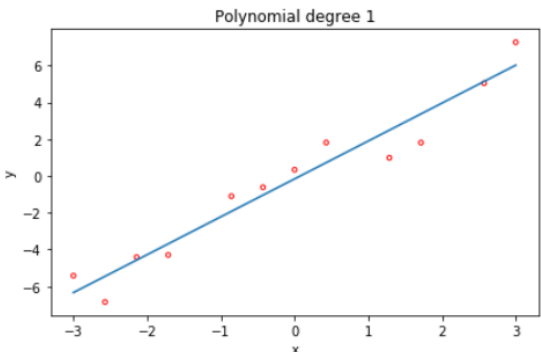
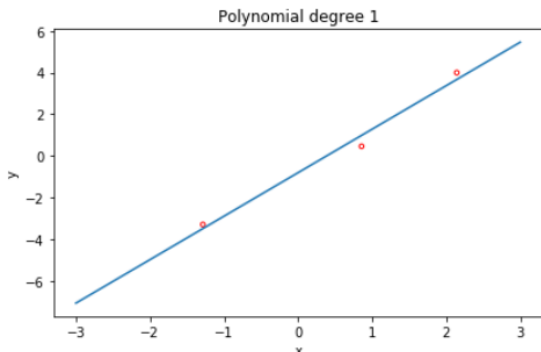
### Training error

Degree=1, RMSE=0.6608466592234012



### Five-fold cross-validation errors

	training	valid
1	Degree=1, RMSE=0.6910131637989462 	Degree=1, RMSE=0.0826318079876958 

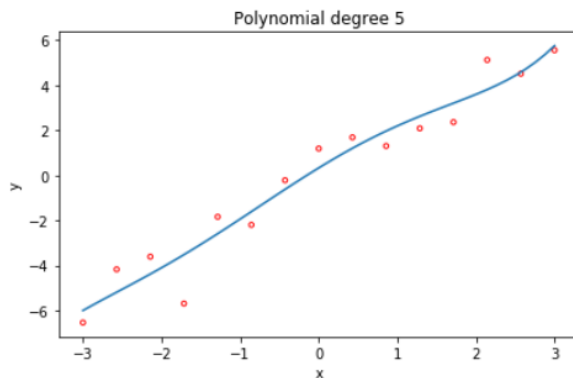
2	Degree=1, RMSE=0.6730182469722918 	Degree=1, RMSE=0.2159158078822646 
3	Degree=1, RMSE=0.6433321731453897 	Degree=1, RMSE=0.02136965063621975 
4	Degree=1, RMSE=0.5184623171237048 	Degree=1, RMSE=0.6996092850236884 
5	Degree=1, RMSE=0.6992285722498134 	Degree=1, RMSE=0.2646937947593224 

- Perform Polynomial Regression with degree 5, 10 and 14, respectively. For each case, show the fitting plot, the training error, and the five-fold cross-validation errors. (Hint: Arrange the polynomial regression equation as follows and solve the model parameter vector  $w$ .)

■ Degree=5

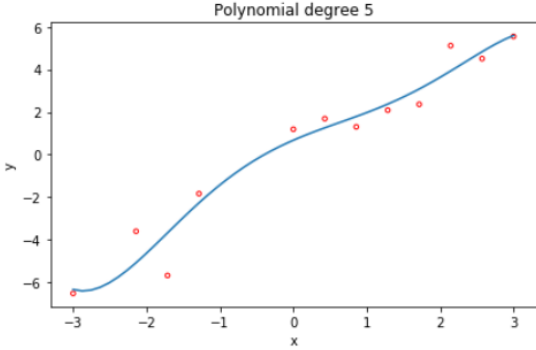
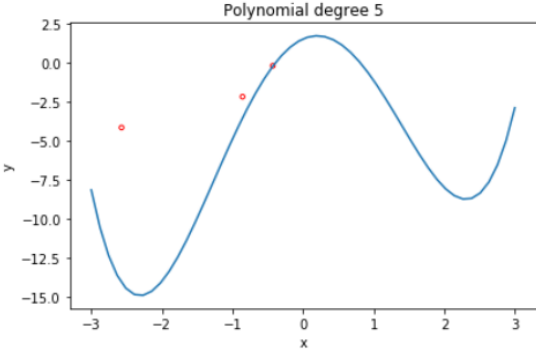
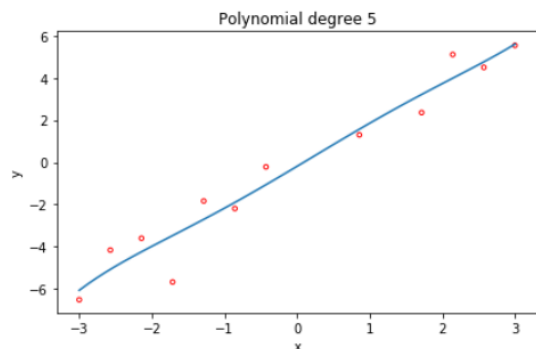
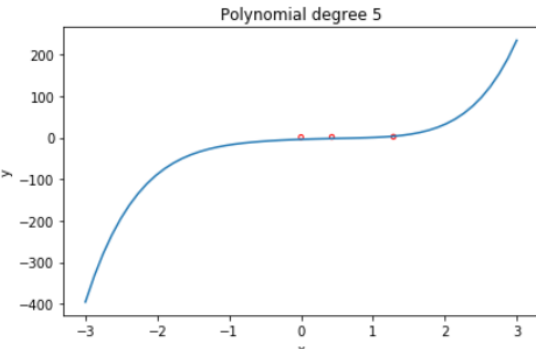
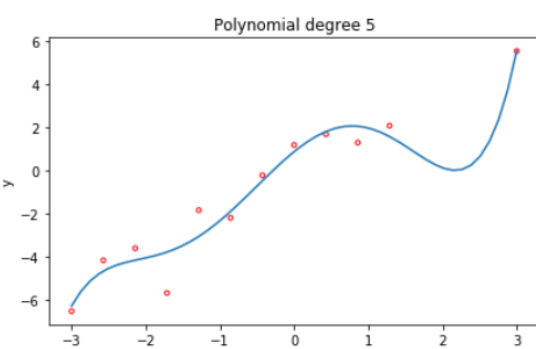
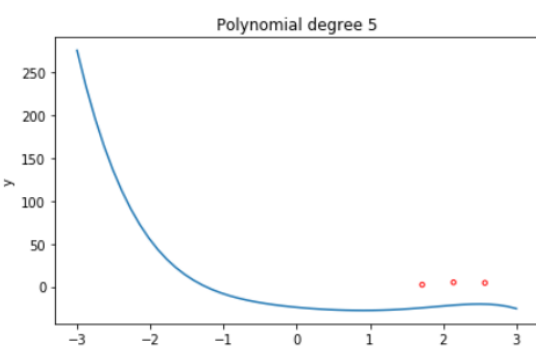
Training error

Degree=5, RMSE=0.6283404445398233



Five-fold cross-validation errors

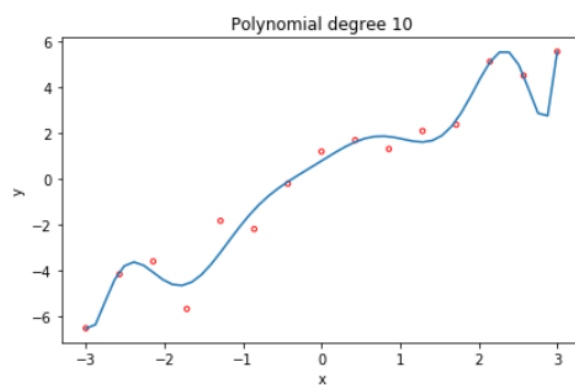
	training	valid
1	Degree=5, RMSE=0.5298677053453411	Degree=5, RMSE=11.660247515332363
2	Degree=5, RMSE=0.4491836894065976	Degree=5, RMSE=3.718096464816144
3	Degree=5, RMSE=0.6199437903881957	Degree=5, RMSE=4.086817144242787

		
4	Degree=5, RMSE=0.6458423938667138	Degree=5, RMSE=2.742977409302637
		
5	Degree=5, RMSE=0.5327445180295234	Degree=5, RMSE=18.476768791115596
		

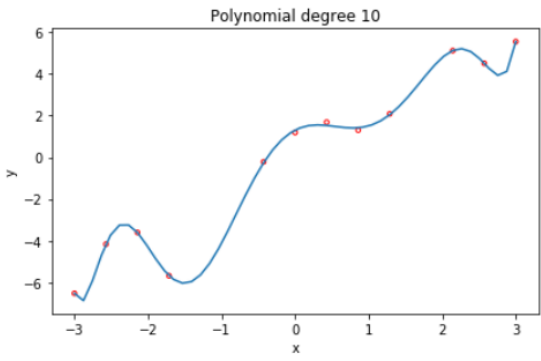
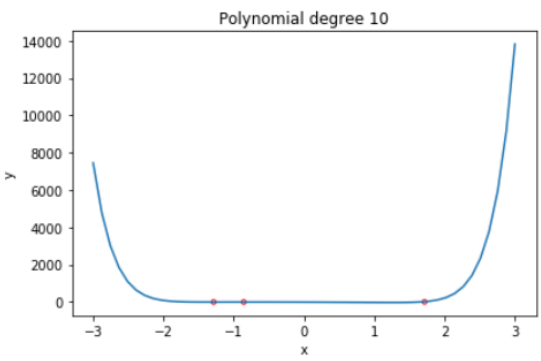
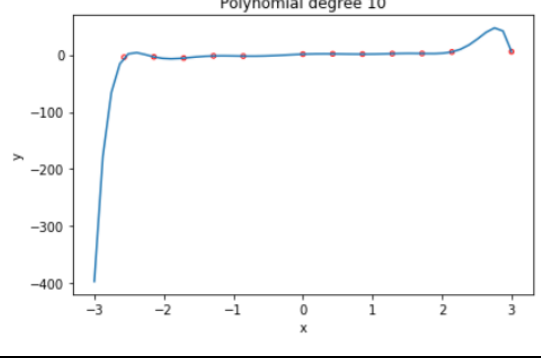
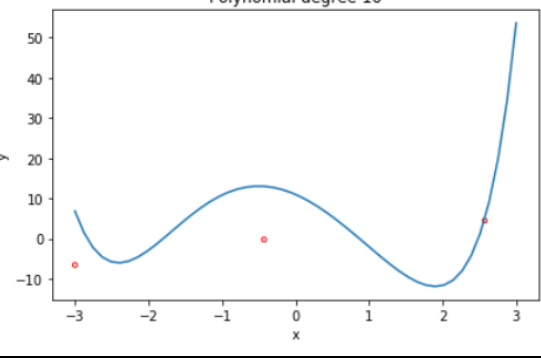
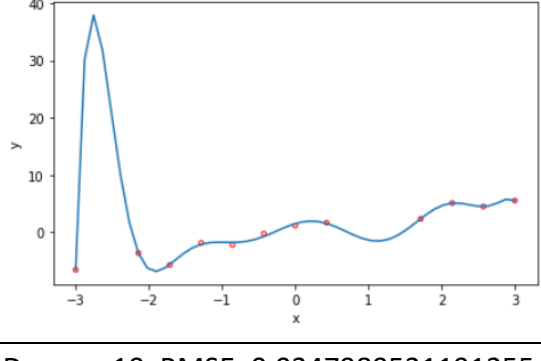
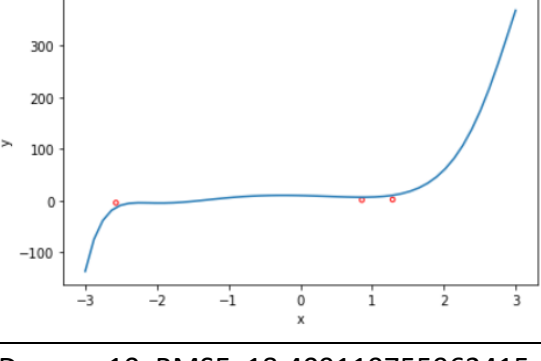
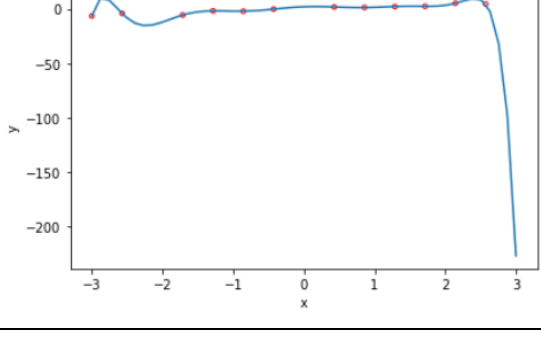
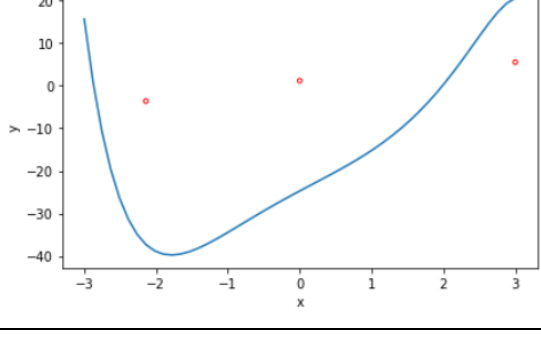
■ Degree=10

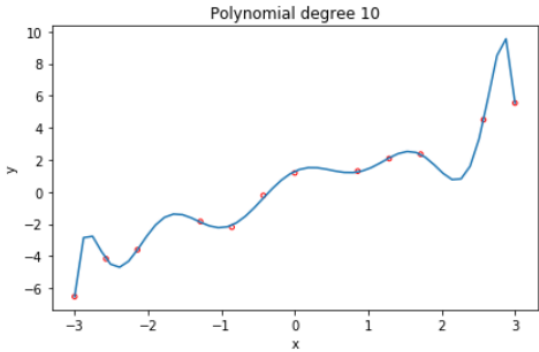
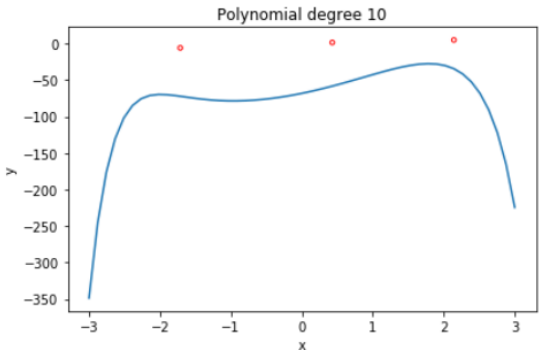
Training error

Degree=10, RMSE=0.3975256867402641



### Five-fold cross-validation errors

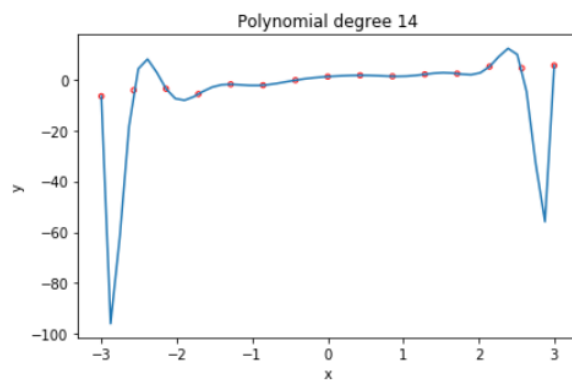
	training	valid
1	Degree=10, RMSE=0.048156581435724	Degree=10, RMSE=5.381509451890489
		
2	Degree=10, RMSE=0.0328826889355411	Degree=10, RMSE=7.689798264979682
		
3	Degree=10, RMSE=0.1636742848310417	Degree=10, RMSE=5.395815093751934
		
4	Degree=10, RMSE=0.0347989531191255	Degree=10, RMSE=18.409119755962415
		

5	Degree=10, RMSE=0.0498893169748178	Degree=10, RMSE=39.83707536208904
	 <p>Polynomial degree 10</p> <p>The plot shows a blue curve fitting 15 red data points. The curve is highly oscillatory, passing through every point, which is characteristic of overfitting. The x-axis ranges from -3 to 3, and the y-axis ranges from -6 to 10.</p>	 <p>Polynomial degree 10</p> <p>The plot shows the same blue curve from the training set applied to new data points (red dots). The curve fails to capture the general trend of the new data, showing a significant discrepancy, which is characteristic of overfitting. The x-axis ranges from -3 to 3, and the y-axis ranges from -350 to 0.</p>

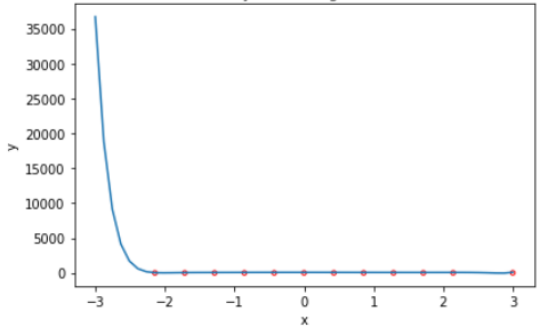
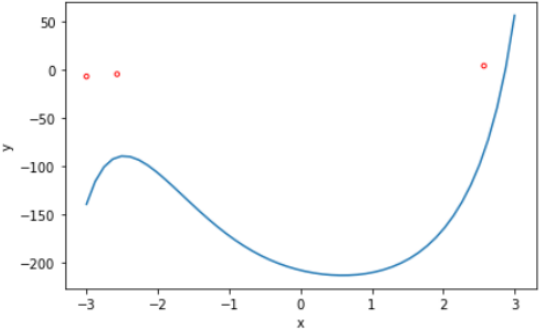
■ Degree=14

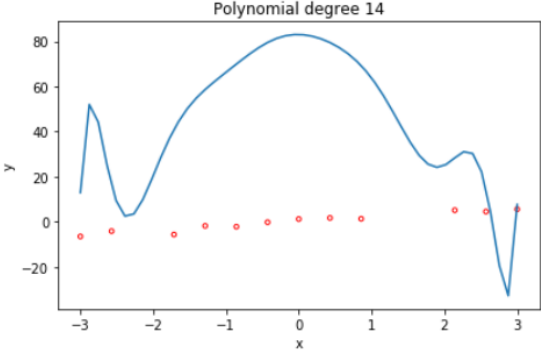
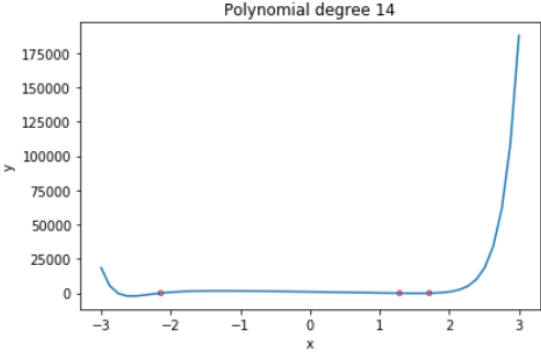
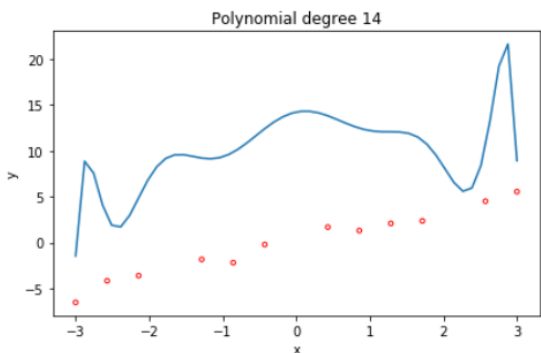
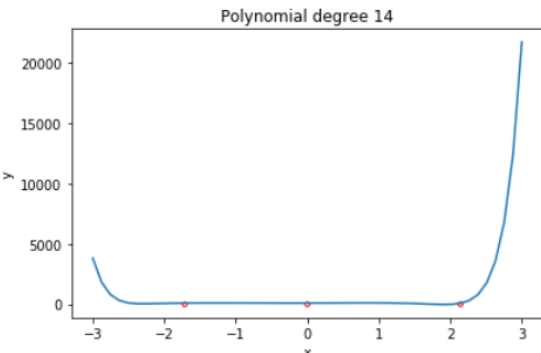
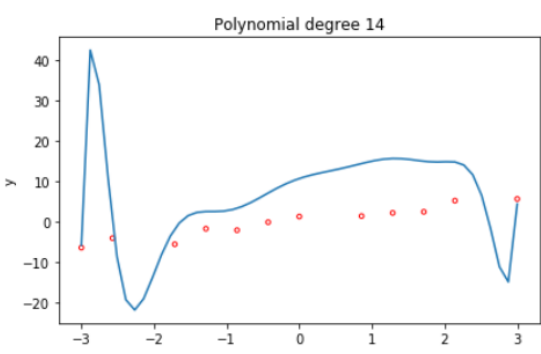
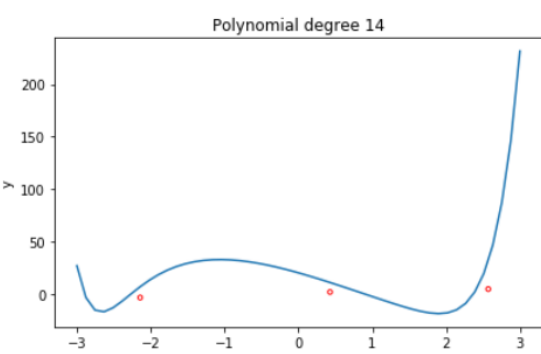
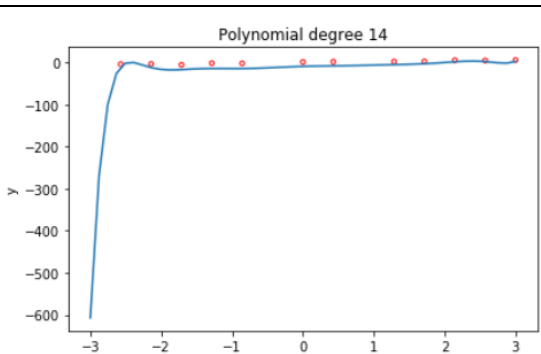
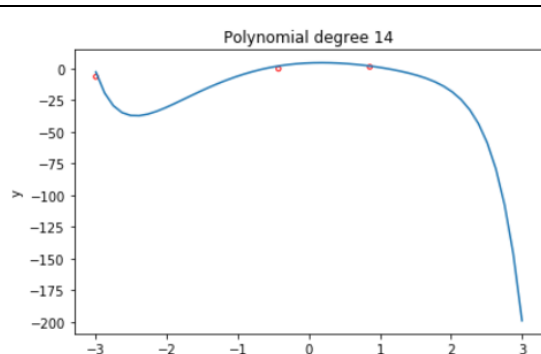
Training error

Degree=14, RMSE=7.136053752255e-06



Five-fold cross-validation errors

	training	valid
1	Degree=14, RMSE=45.23296143162496	Degree=14, RMSE=74.34497108158169
	 <p>Polynomial degree 14</p> <p>The plot shows a blue curve fitting 15 red data points. The curve is highly oscillatory, passing through every point, which is characteristic of overfitting. The x-axis ranges from -3 to 3, and the y-axis ranges from 0 to 35000.</p>	 <p>Polynomial degree 14</p> <p>The plot shows the same blue curve from the training set applied to new data points (red dots). The curve fails to capture the general trend of the new data, showing a significant discrepancy, which is characteristic of overfitting. The x-axis ranges from -3 to 3, and the y-axis ranges from -200 to 50.</p>
2	Degree=14, RMSE=38.82962743875646	Degree=14, RMSE=24.01688485555414

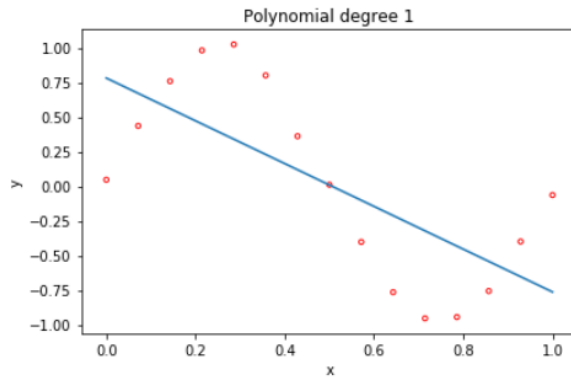
		
3	Degree=14, RMSE=6.673386260196252	Degree=14, RMSE=66.16691285242946
		
4	Degree=14, RMSE=5.8527239630012495	Degree=14, RMSE=12.35147984701117
		
5	Degree=14, RMSE=6.048017752680724	Degree=14, RMSE=1.807152057472771
		

- Change the model to  $y = \sin(2\pi x) + \epsilon$  with the noise  $\epsilon \sim N(0, 0.04)$  and (equal spacing)  $x \in [0, 1]$ . Then repeat those stated in 2) and 3). Compare the results with linear/polynomail regression on different datasets.

■ Degree=1

Training error

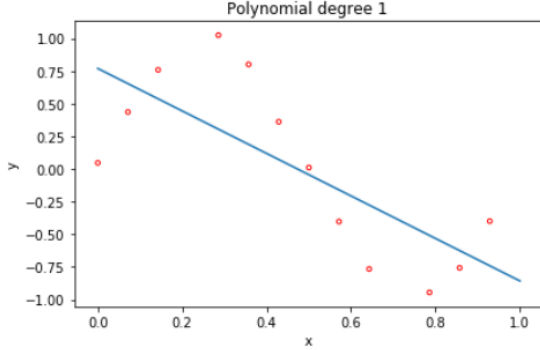
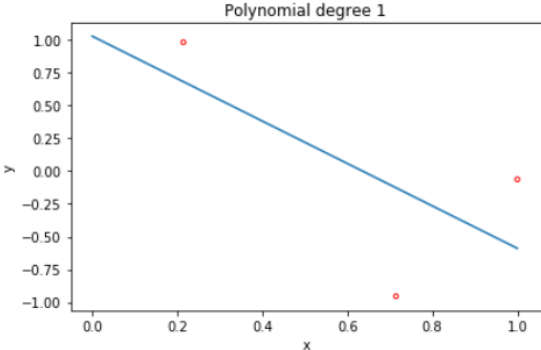
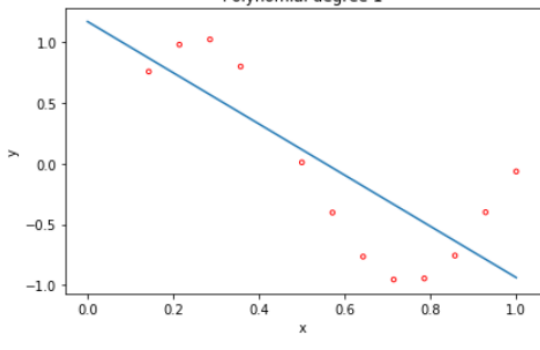
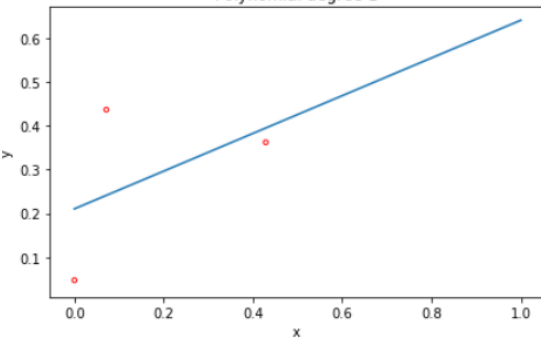
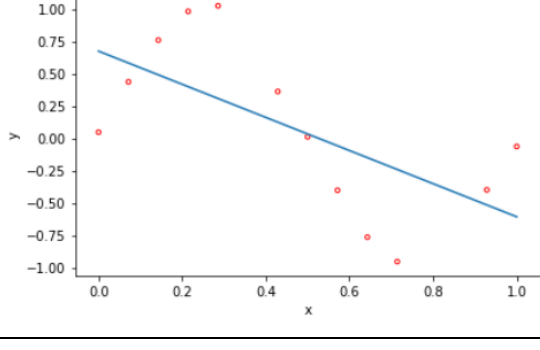
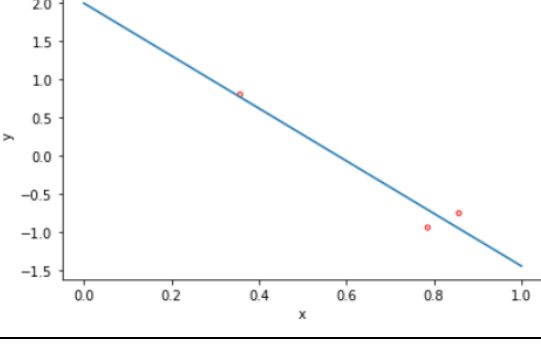
Degree=1, RMSE=0.3381117366078396



Five-fold cross-validation errors

	training	valid
1	Degree=1, RMSE=0.3513241594046753	Degree=1, RMSE=0.00841157733406936
2	Degree=1, RMSE=0.33763457280060355	Degree=1, RMSE=0.1488449843647542
3	Degree=1, RMSE=0.30421926826372075	Degree=1, RMSE=0.4184455303571044

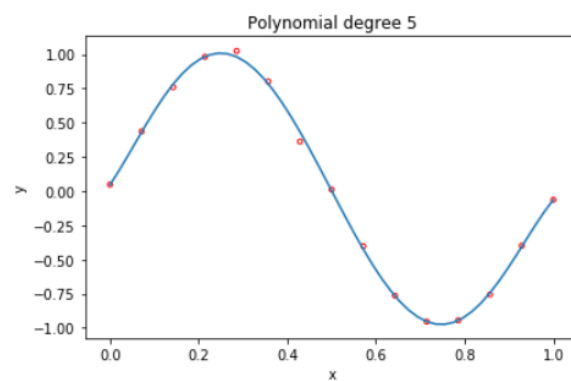


		
4	Degree=1, RMSE=0.319042404406932	Degree=1, RMSE=0.10478813807483199
		
5	Degree=1, RMSE=0.33607727847437785	Degree=1, RMSE=0.12657315863893567
		

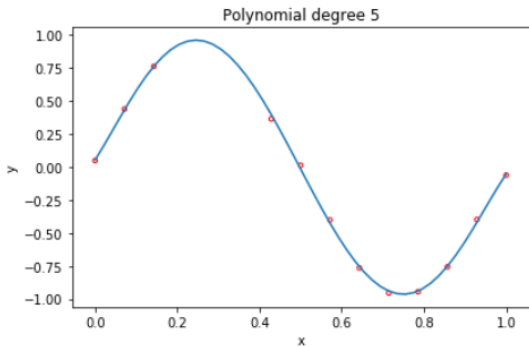
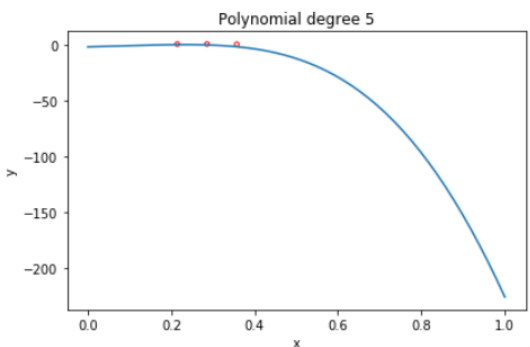
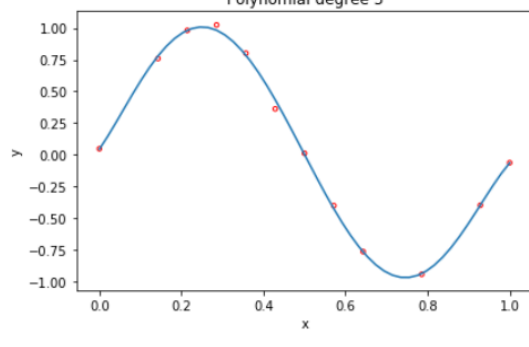
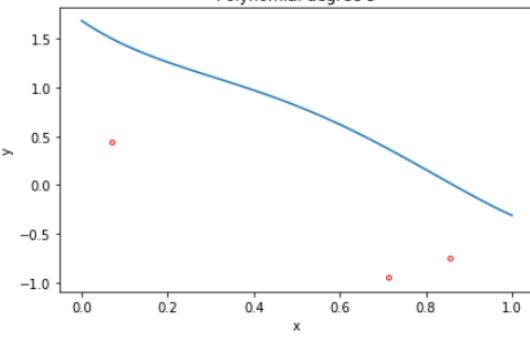
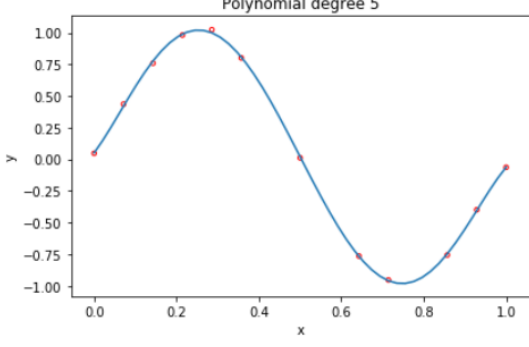
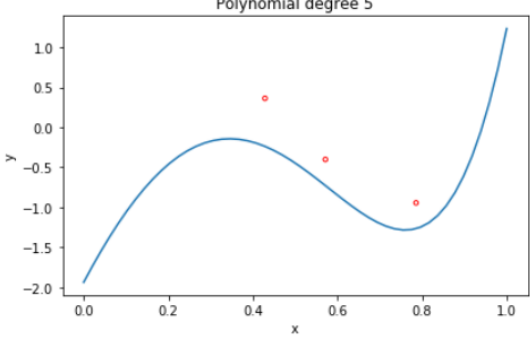
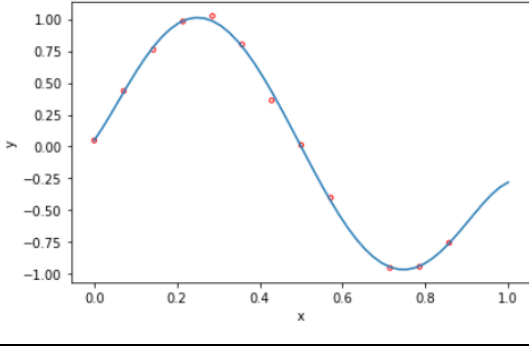
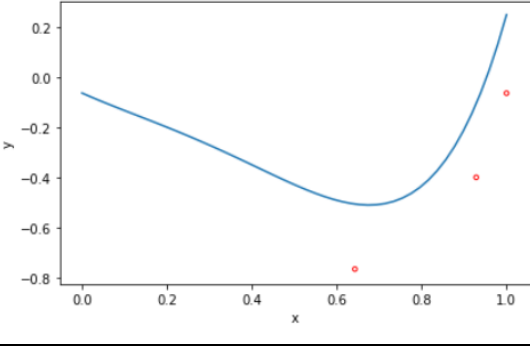
## ■ Degree=5

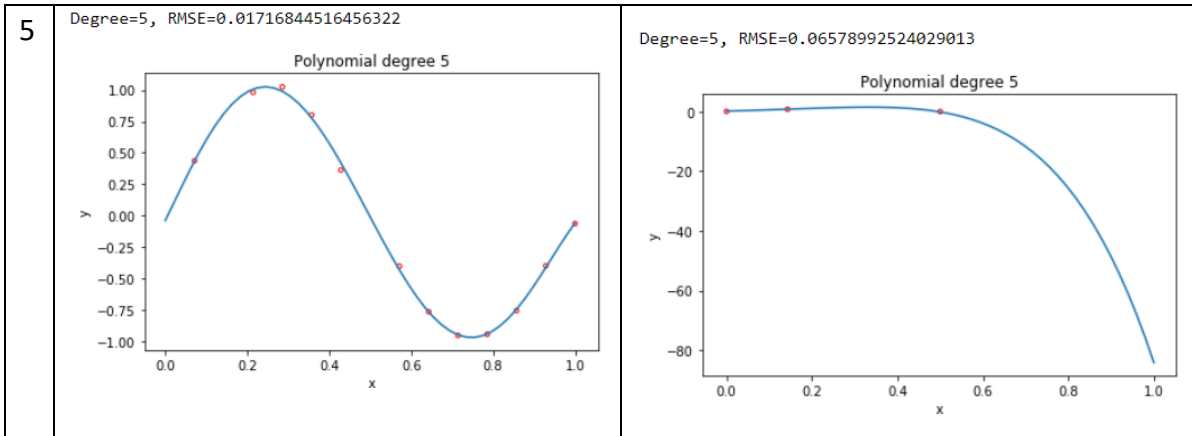
### Training error

Degree=5, RMSE=0.017025404734936914



## Five-fold cross-validation errors

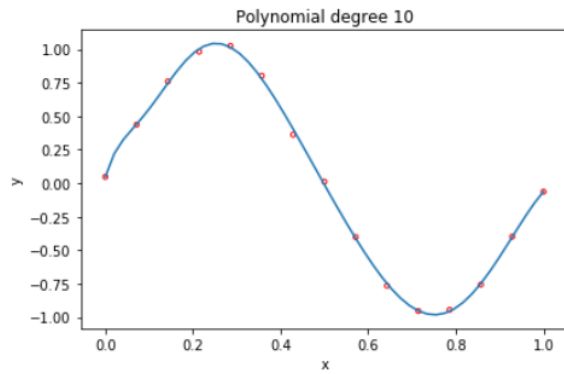
	training	valid
1	<p>Degree=5, RMSE=0.011951427621339135</p> 	<p>Degree=5, RMSE=0.9824565288713851</p> 
2	<p>Degree=5, RMSE=0.018657195409370313</p> 	<p>Degree=5, RMSE=0.7608216066334113</p> 
3	<p>Degree=5, RMSE=0.0079938152690887</p> 	<p>Degree=5, RMSE=0.30861245183275676</p> 
4	<p>Degree=5, RMSE=0.018468231922172892</p> 	<p>Degree=5, RMSE=0.20483656554802698</p> 



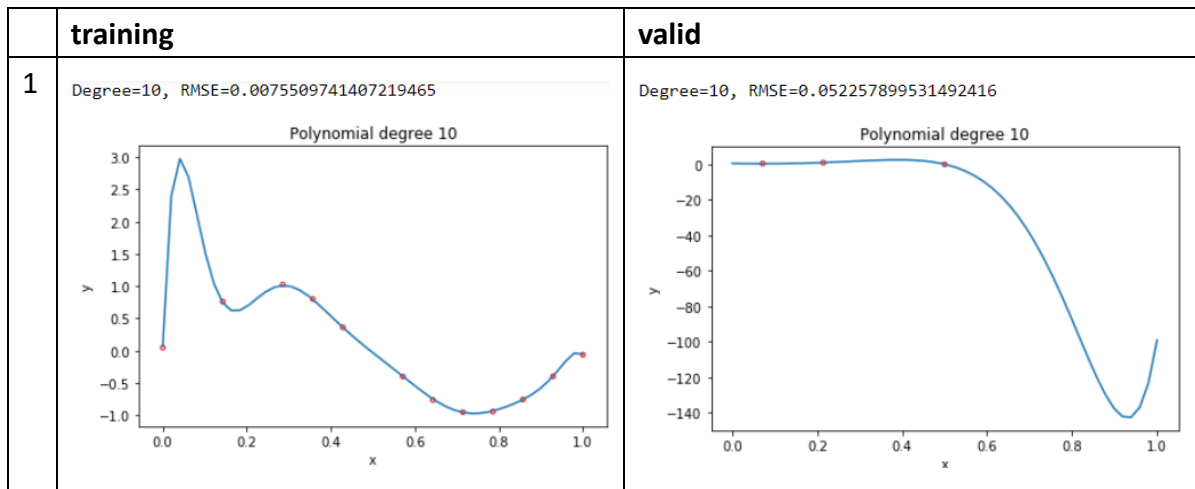
■ Degree=10

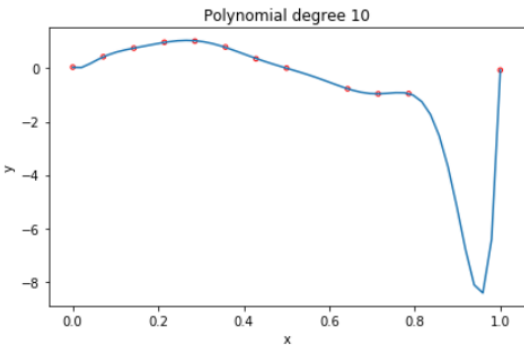
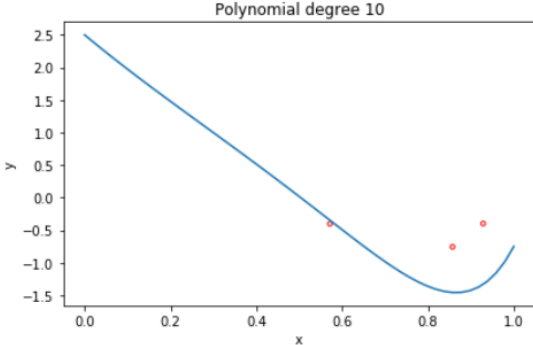
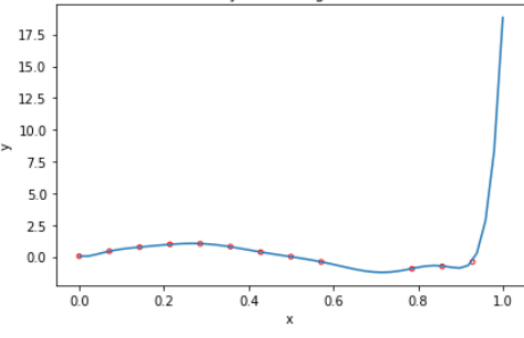
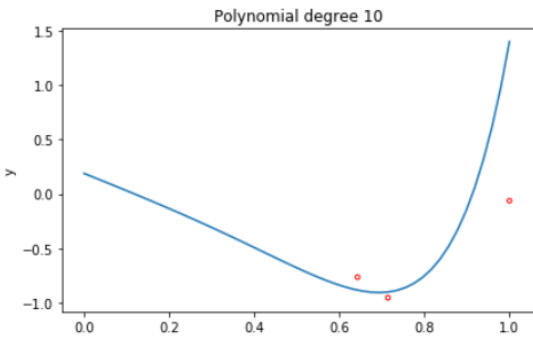
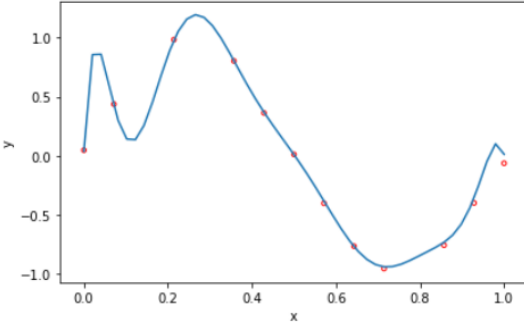
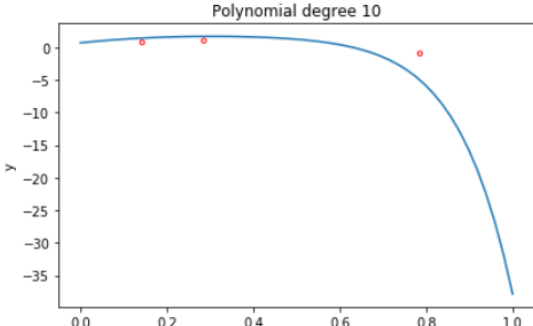
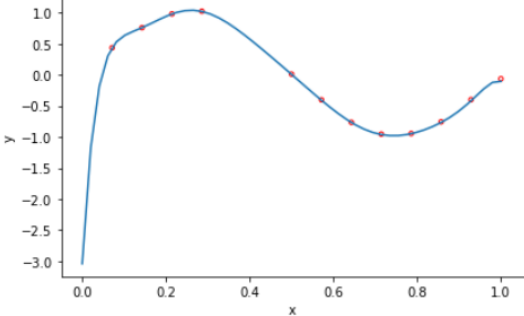
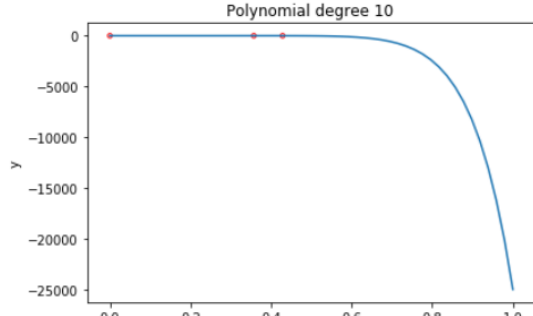
### Training error

Degree=10, RMSE=0.011747387744402643



### Five-fold cross-validation errors

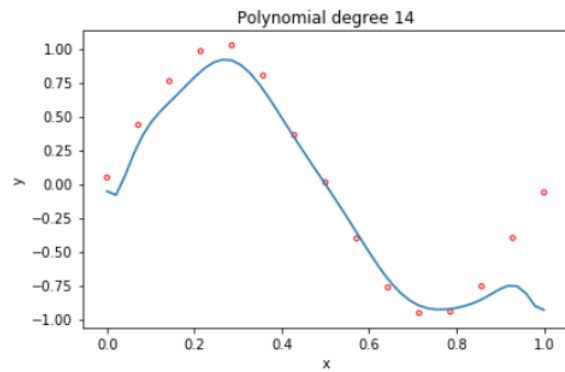


2	<p>Degree=10, RMSE=0.005509561461011673</p> 	<p>Degree=10, RMSE=0.47734765454801276</p> 
3	<p>Degree=10, RMSE=0.016849982581915942</p> 	<p>Degree=10, RMSE=0.5990049713882291</p> 
4	<p>Degree=10, RMSE=0.01966295027334867</p> 	<p>Degree=10, RMSE=1.7085018100298333</p> 
5	<p>Degree=10, RMSE=0.009822547482078545</p> 	<p>Degree=10, RMSE=0.7537526978417897</p> 

■ Degree=14

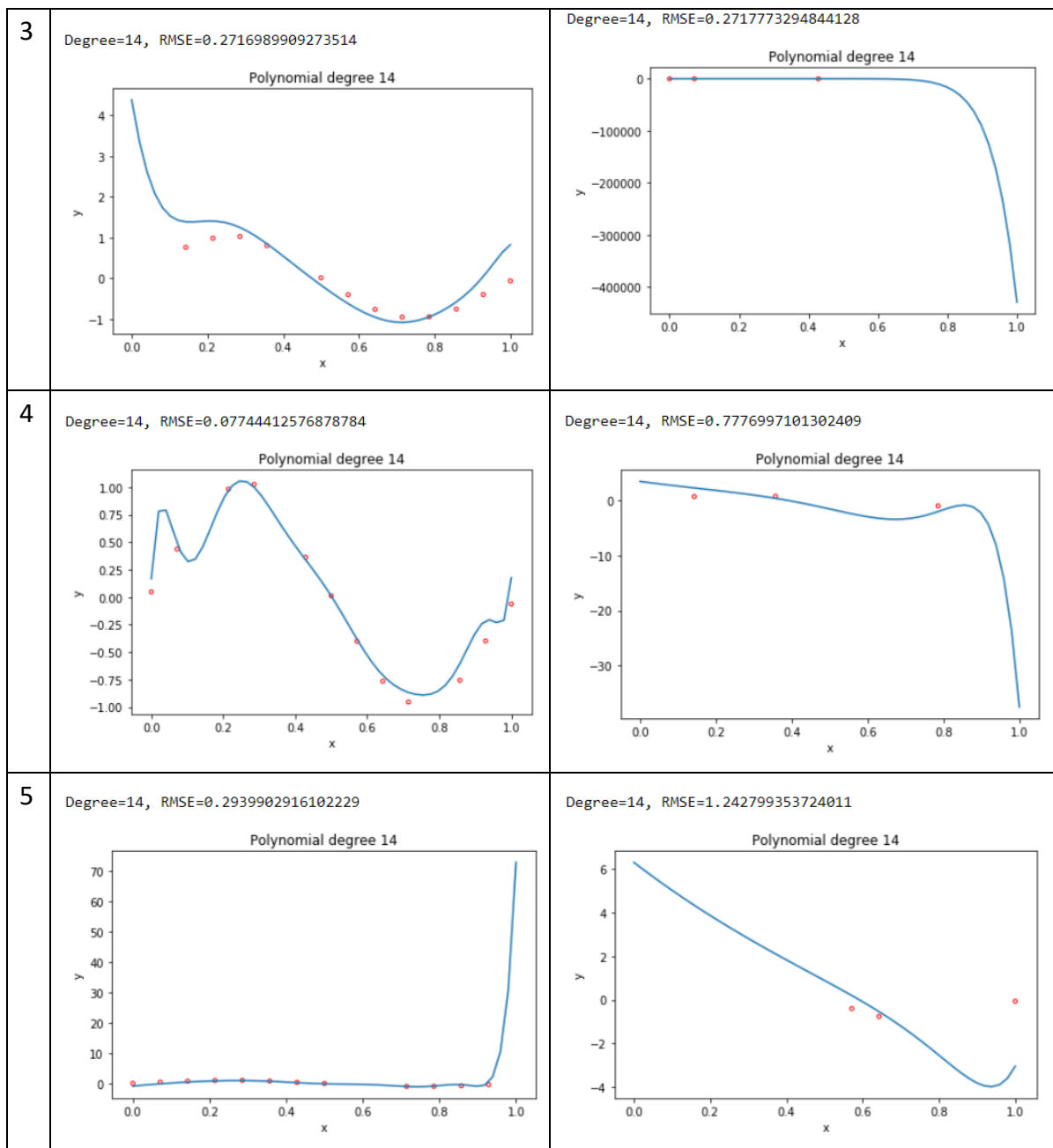
Training error

Degree=14, RMSE=0.18215695742583007



Five-fold cross-validation errors

	training	valid
1	<p>Degree=14, RMSE=0.04286383164989321</p> <p>A plot titled "Polynomial degree 14" showing the training error for fold 1. The x-axis ranges from 0.0 to 1.0, and the y-axis ranges from -30 to 0. Red dots represent the training data points, and a blue line represents the polynomial fit. The fit is very smooth and follows the general trend of the data points, indicating a good fit to the training data.</p>	<p>Degree=14, RMSE=2.612141125753481</p> <p>A plot titled "Polynomial degree 14" showing the validation error for fold 1. The x-axis ranges from 0.0 to 1.0, and the y-axis ranges from -16 to 0. Red dots represent the validation data points, and a blue line represents the polynomial fit. The fit is very smooth and follows the general trend of the data points, indicating a good fit to the validation data.</p>
2	<p>Degree=14, RMSE=2.196955962159729</p> <p>A plot titled "Polynomial degree 14" showing the training error for fold 2. The x-axis ranges from 0.0 to 1.0, and the y-axis ranges from -8 to 2. Red dots represent the training data points, and a blue line represents the polynomial fit. The fit is very smooth and follows the general trend of the data points, indicating a good fit to the training data.</p>	<p>Degree=14, RMSE=4.124748085494175</p> <p>A plot titled "Polynomial degree 14" showing the validation error for fold 2. The x-axis ranges from 0.0 to 1.0, and the y-axis ranges from -50 to 0. Red dots represent the validation data points, and a blue line represents the polynomial fit. The fit is very smooth and follows the general trend of the data points, indicating a good fit to the validation data.</p>

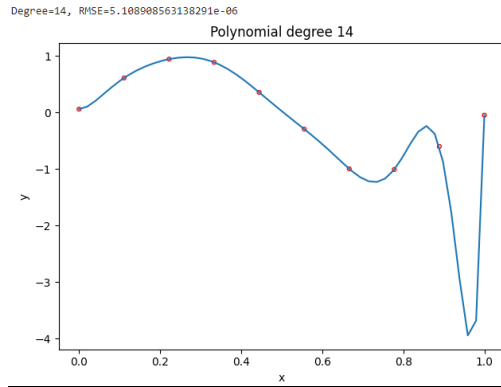


■ 比較: linear regression 對曲線類型的 dataset(Ex: sin)不適用，應用 polynomial regression 解決曲線類型的 dataset。

➤ Following 4), perform polynomial regression with degree 14 by varying the number of training data points  $m = 10, 80, 320$ . Show the five-fold cross-validation errors and the fitting plots. Compare the results to those in 4).

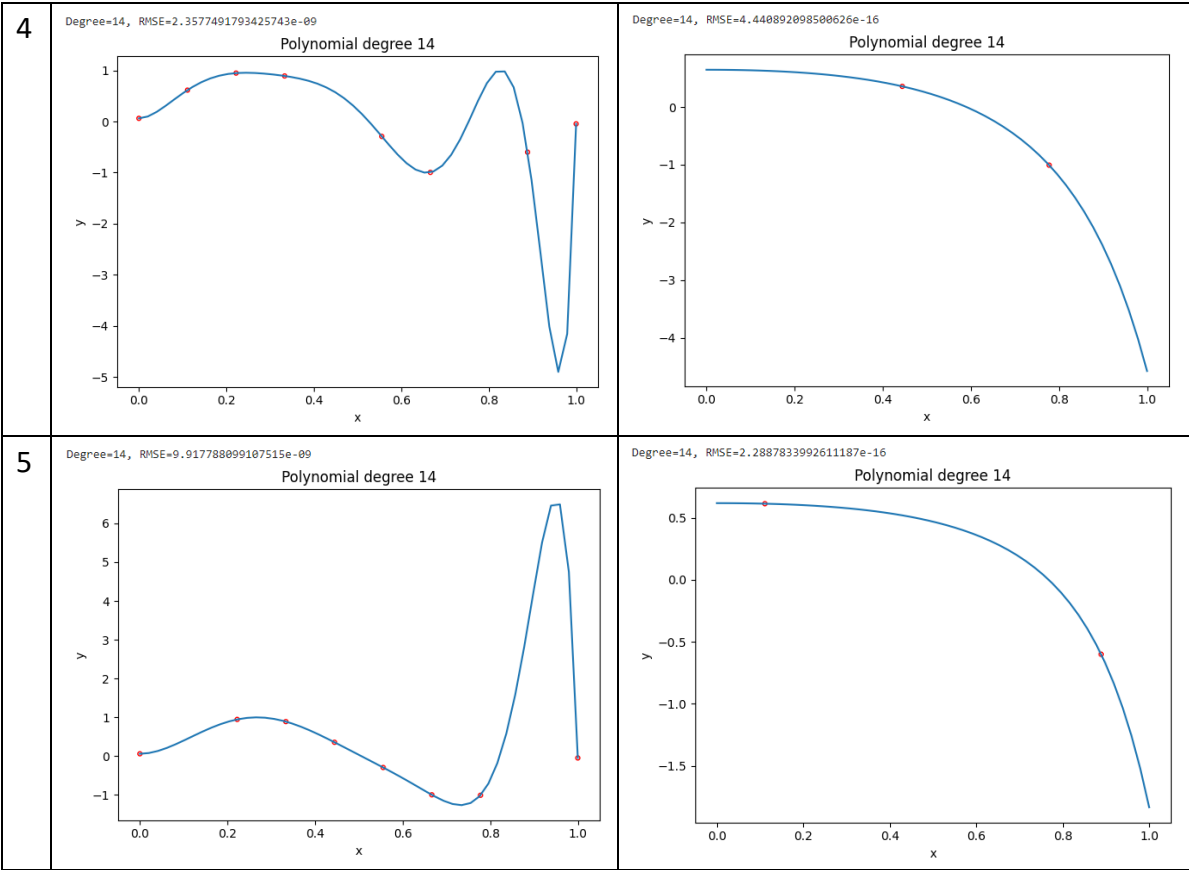
■  $m=10$

Training error



### Five-fold cross-validation errors

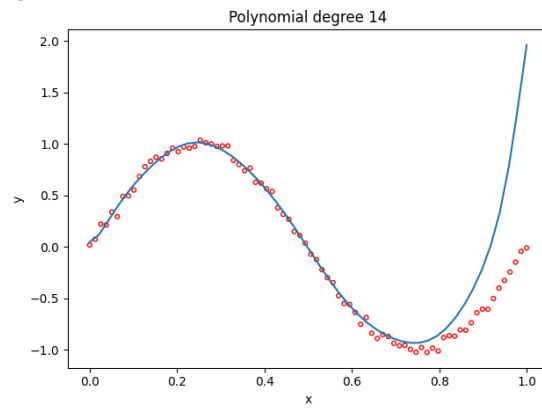
	training	valid
1	<p>Degree=14, RMSE=1.5959613378821164e-09</p> <p>Polynomial degree 14</p>	<p>Degree=14, RMSE=4.652682298944613e-16</p> <p>Polynomial degree 14</p>
2	<p>Degree=14, RMSE=2.8736337424499165e-09</p> <p>Polynomial degree 14</p>	<p>Degree=14, RMSE=5.551115123125783e-16</p> <p>Polynomial degree 14</p>
3	<p>Degree=14, RMSE=8.556286486679673e-10</p> <p>Polynomial degree 14</p>	<p>Degree=14, RMSE=2.0816681711721685e-17</p> <p>Polynomial degree 14</p>



■ m=80

### Training error

Degree=14, RMSE=0.2750086613851848



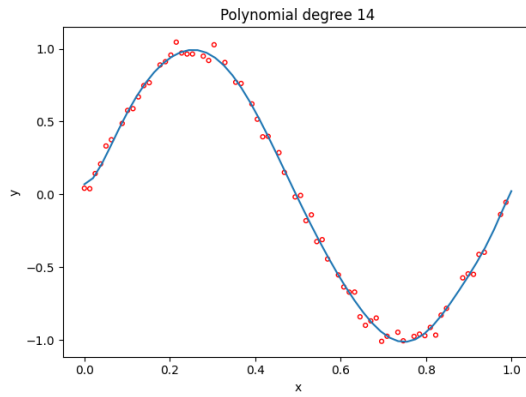
### Five-fold cross-validation errors

training	valid
----------	-------

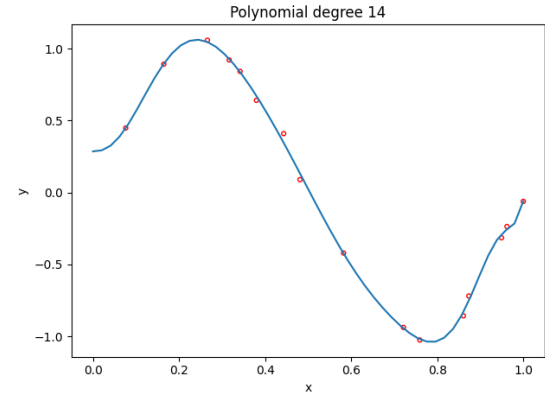


1

Degree=14, RMSE=0.027374141048660316

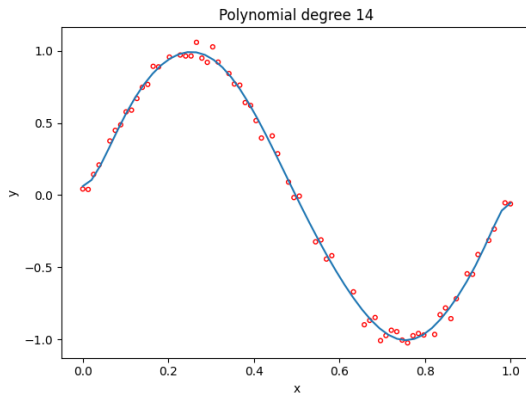


Degree=14, RMSE=0.01771066867462375

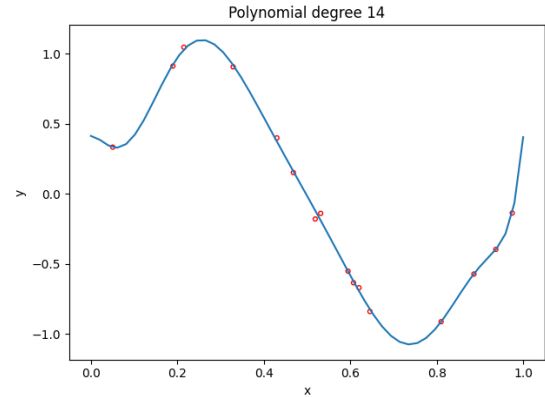


2

Degree=14, RMSE=0.02822626236667774

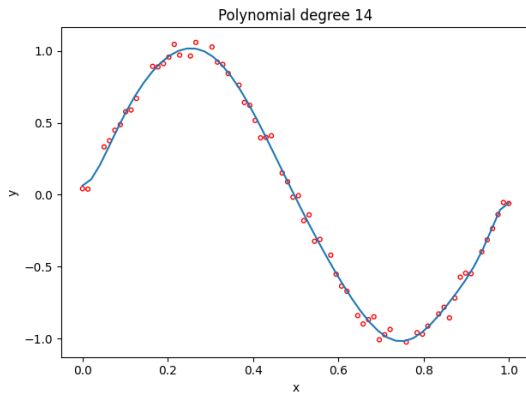


Degree=14, RMSE=0.0162729834785817

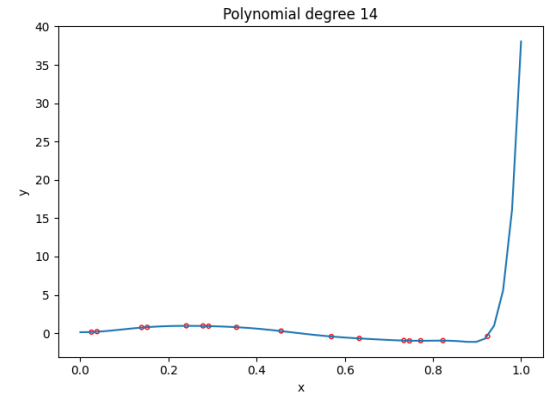


3

Degree=14, RMSE=0.02873690379912878

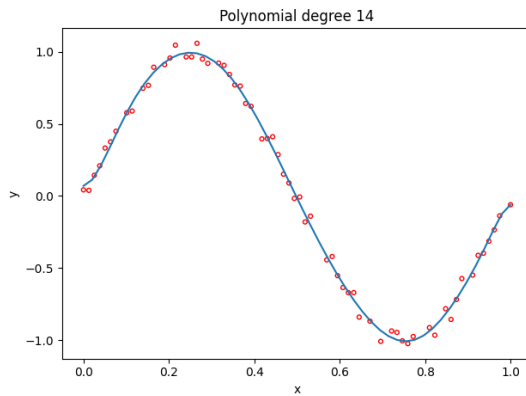


Degree=14, RMSE=0.009838102510673859

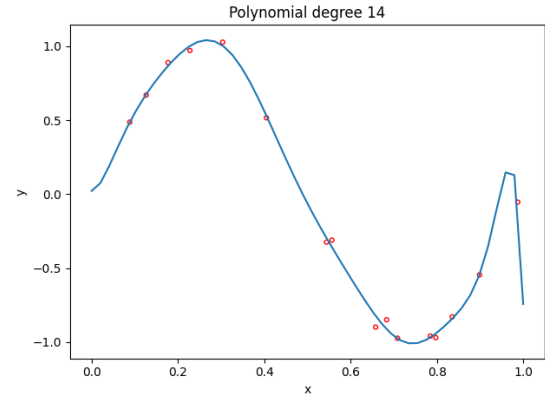


4

Degree=14, RMSE=0.028903071572375698

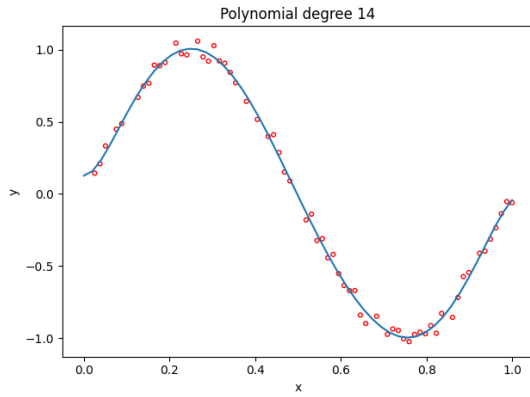


Degree=14, RMSE=0.02254432455853114

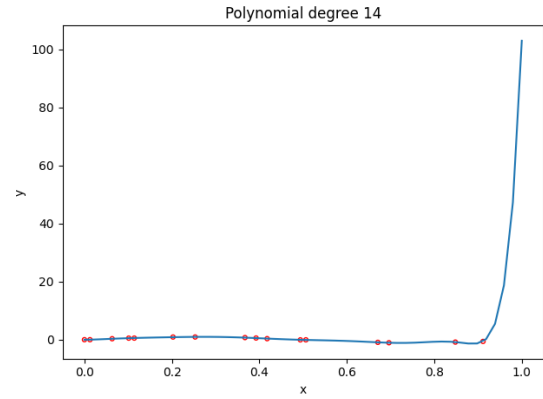


5

Degree=14, RMSE=0.02826525293673486



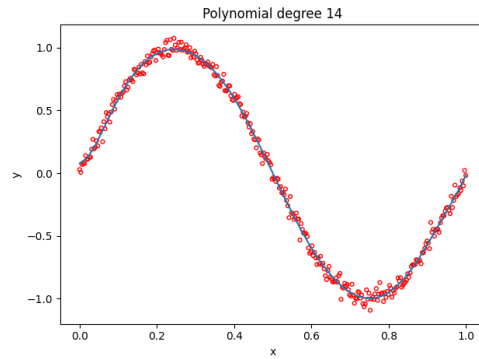
Degree=14, RMSE=0.017961554547374737



■ m=320

Training error

Degree=14, RMSE=0.02822448148805892

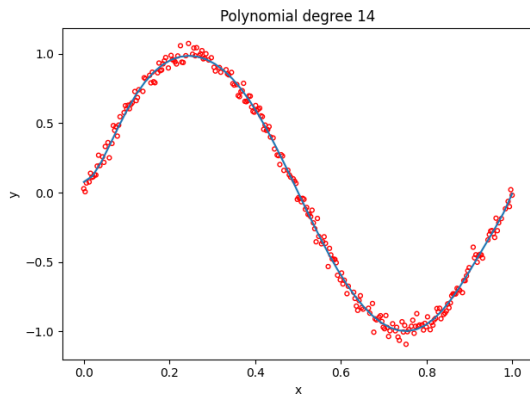


Five-fold cross-validation errors

**training**

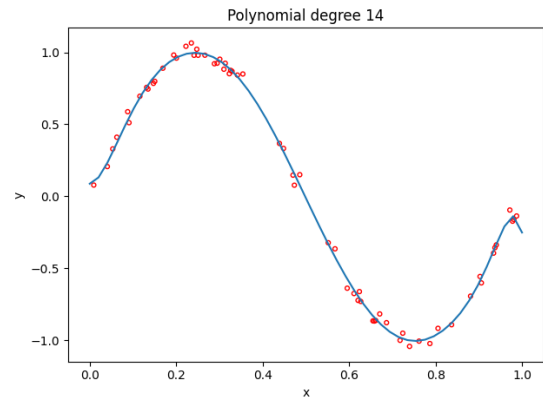
1

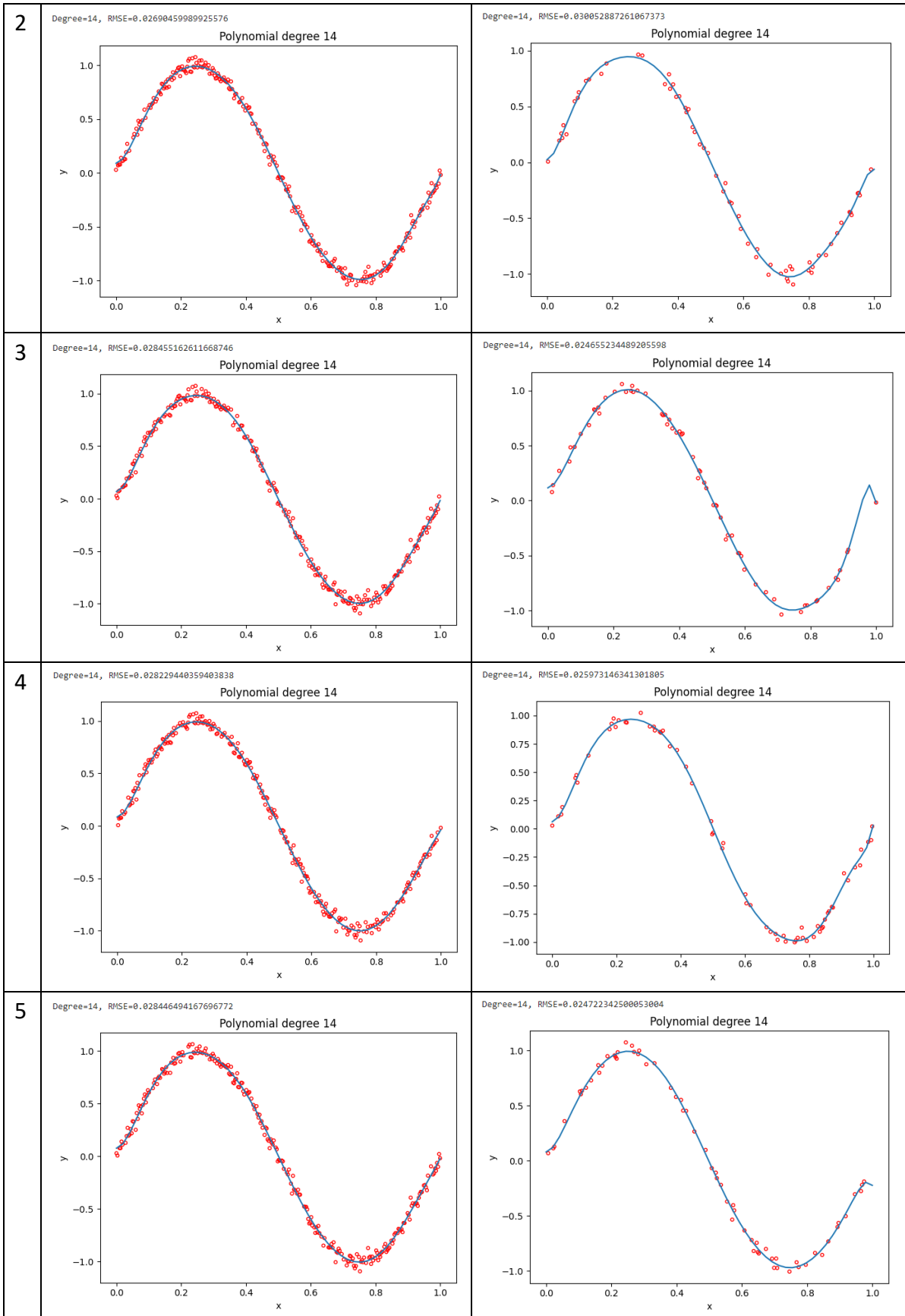
Degree=14, RMSE=0.028539683347882526



**valid**

Degree=14, RMSE=0.024675652812679184





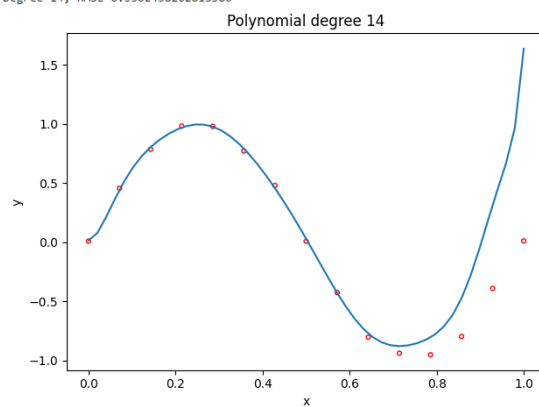
- 比較:  $m=10$  時，因為 data 很少，所以模型遇到沒看過的資料時會隨機猜測，導致預測結果不如預期。對比  $m=320$  的波形明顯可看出在 num\_point 數量多的情況下模型預測準確率較 num\_point 少的高很多。

- Following 4), perform polynomial regression of degree 14 via regularization. Compare the results by setting  $\lambda = 0, 0.001/m, 1/m, 1000/m$ , where  $m = 15$  is the number of data points (with  $x = 0, 1/(m-1), 2/(m-1), \dots, 1$ ). Show the five-fold cross-validation errors and the fitting plots.

- $\lambda = 0$

### Training error

Degree=14, RMSE=0.3302438202813586

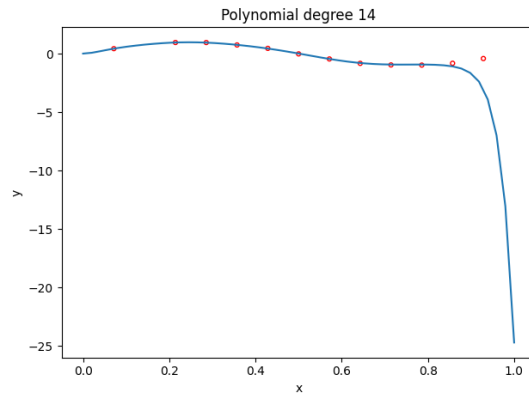


### Five-fold cross-validation errors

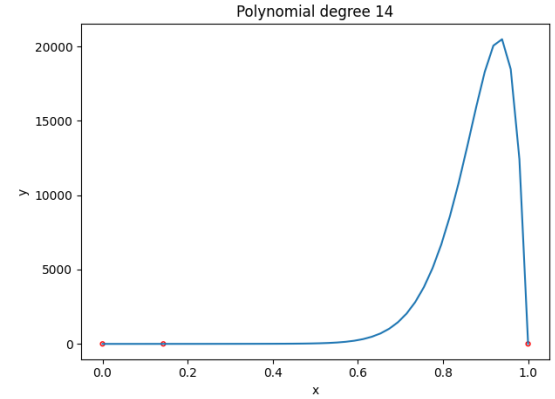
	training	valid
1	<p>Degree=14, RMSE=0.10171318410510395</p> <p>A plot titled "Polynomial degree 14" showing the training error for fold 1. The x-axis is labeled 'x' and ranges from 0.0 to 1.0. The y-axis is labeled 'y' and ranges from -2 to 12. Red dots represent the data points, and a blue line represents the polynomial fit. The fit is very smooth and follows the general trend of the data points.</p>	<p>Degree=14, RMSE=0.2658217545935833</p> <p>A plot titled "Polynomial degree 14" showing the validation error for fold 1. The x-axis is labeled 'x' and ranges from 0.0 to 1.0. The y-axis is labeled 'y' and ranges from 0 to 120. Red dots represent the data points, and a blue line represents the polynomial fit. The fit is very smooth and follows the general trend of the data points.</p>

2

Degree=14, RMSE=0.5340411136129597

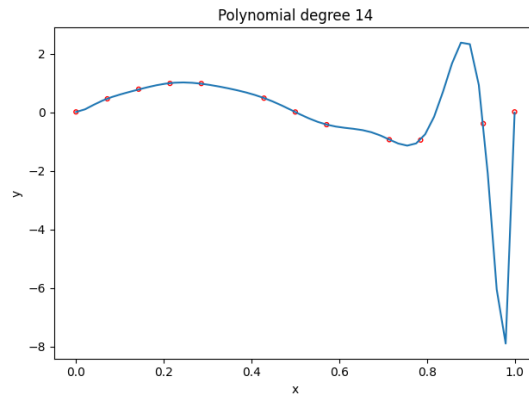


Degree=14, RMSE=0.0833297798658421

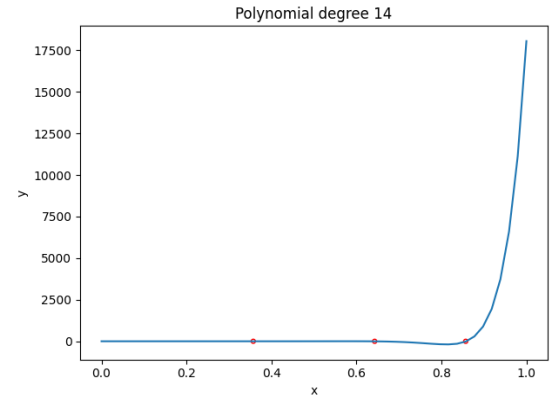


3

Degree=14, RMSE=0.006705815936392849

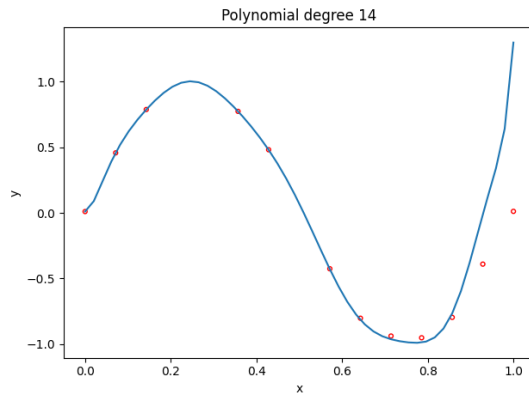


Degree=14, RMSE=5.642128015915632

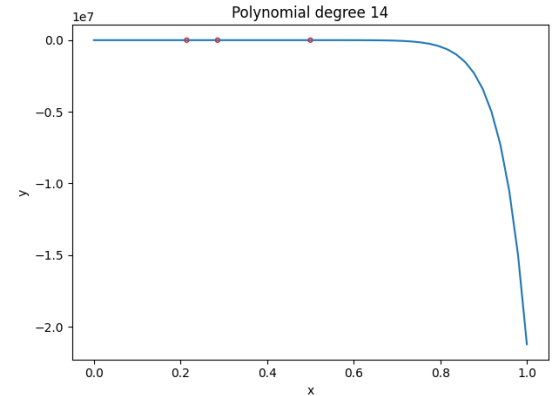


4

Degree=14, RMSE=0.27368353929368183

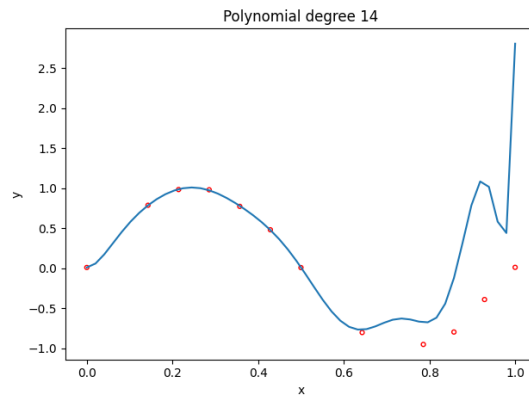


Degree=14, RMSE=0.2366764681273377

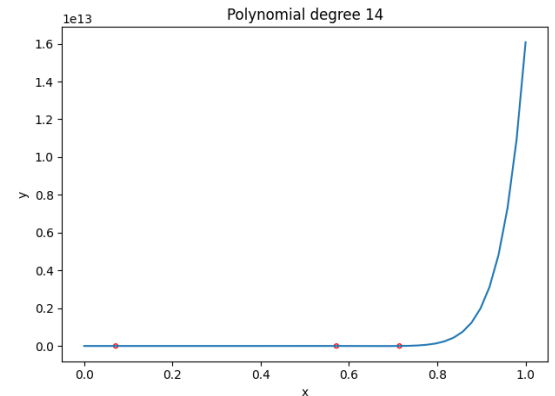


5

Degree=14, RMSE=0.0635628481090699



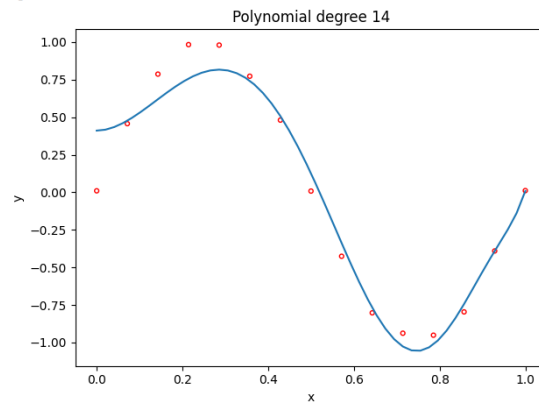
Degree=14, RMSE=214.84346220129635



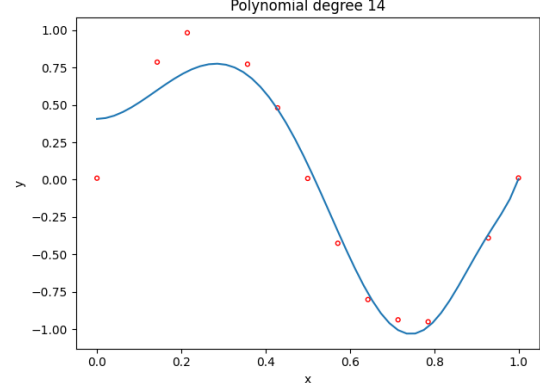
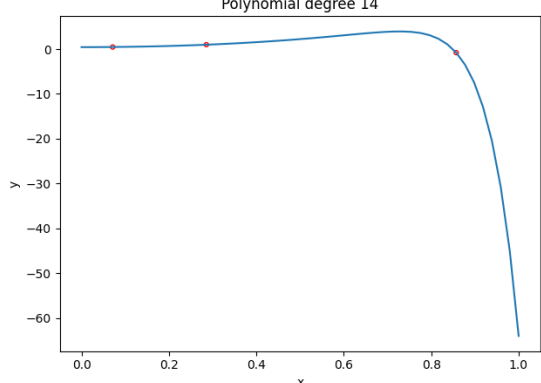
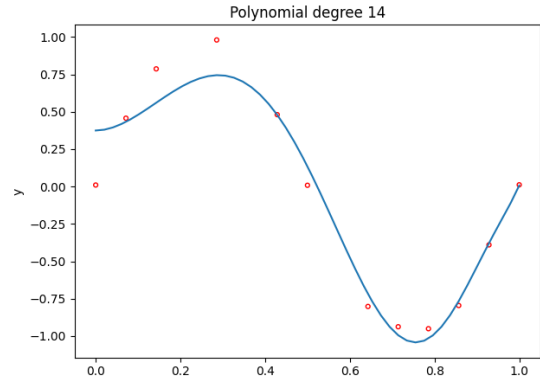
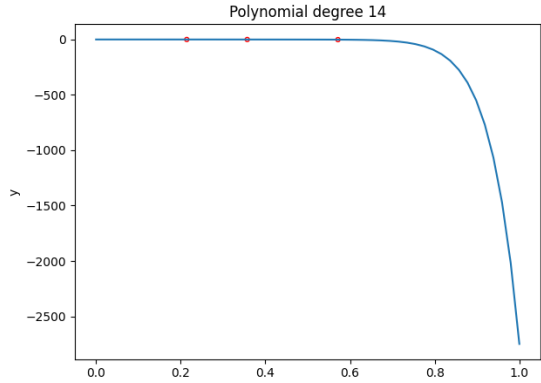
■  $\lambda = 0.001/15$

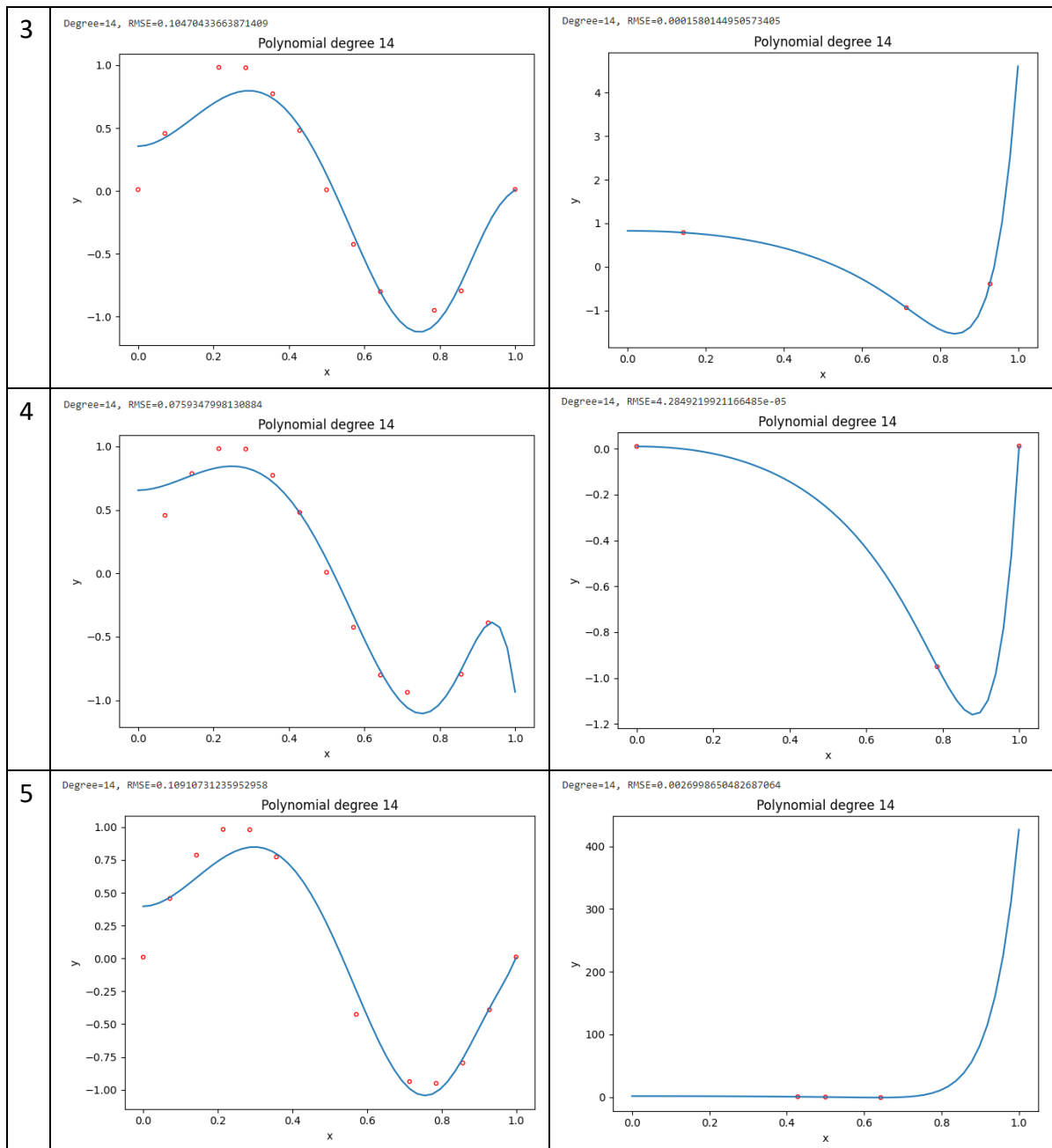
### Training error

Degree=14, RMSE=0.100907961947229



### Five-fold cross-validation errors

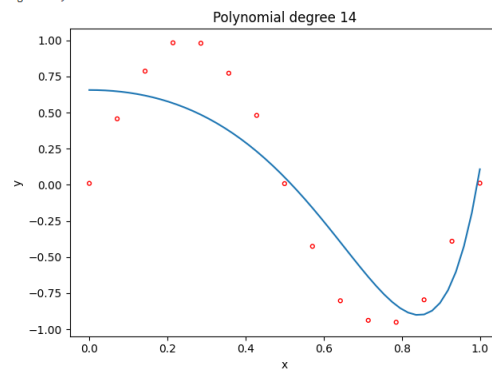
	training	valid
1	<p>Degree=14, RMSE=0.10957180894168682</p> <p>Polynomial degree 14</p> 	<p>Degree=14, RMSE=0.003058046077530753</p> <p>Polynomial degree 14</p> 
2	<p>Degree=14, RMSE=0.10638270483906974</p> <p>Polynomial degree 14</p> 	<p>Degree=14, RMSE=0.001015813918407533</p> <p>Polynomial degree 14</p> 



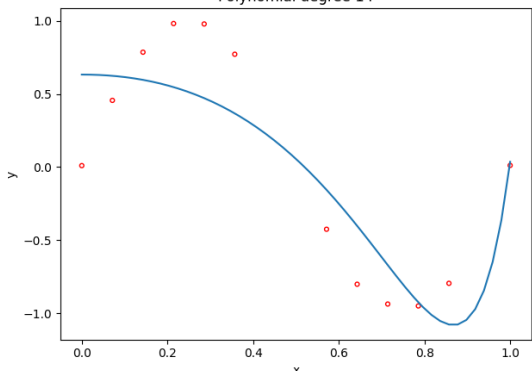
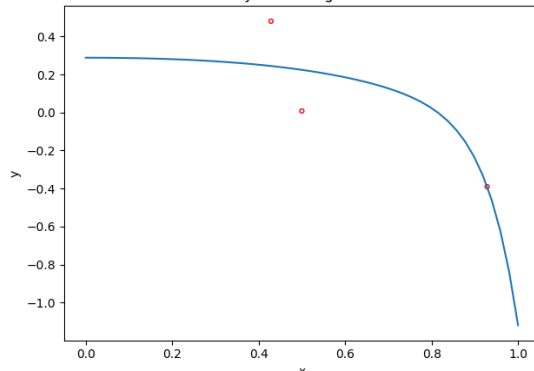
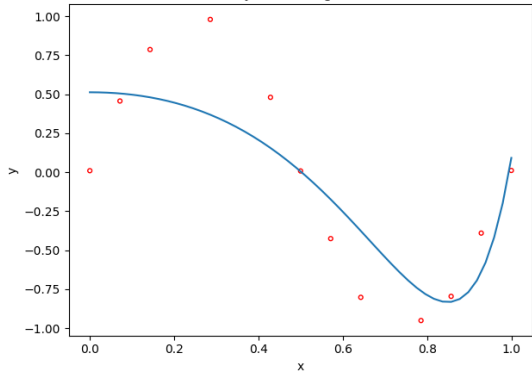
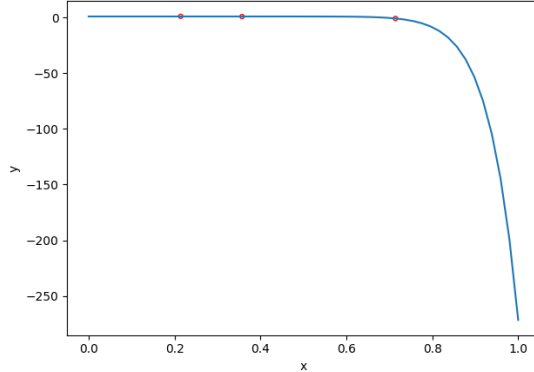
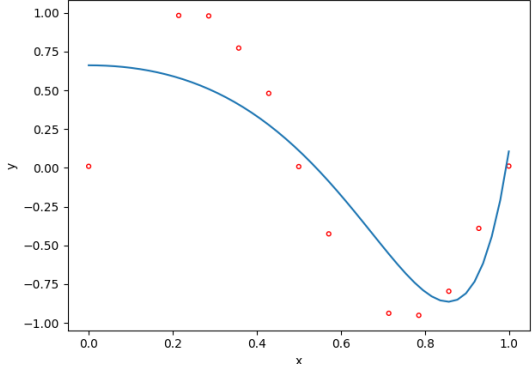
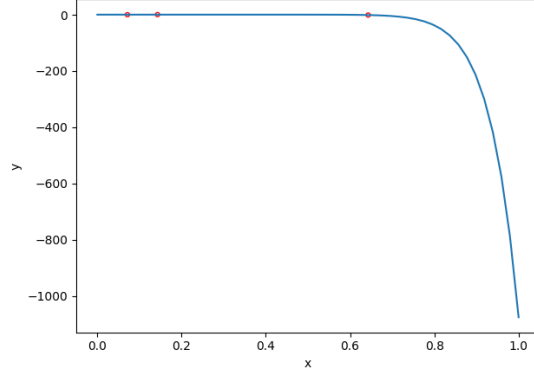
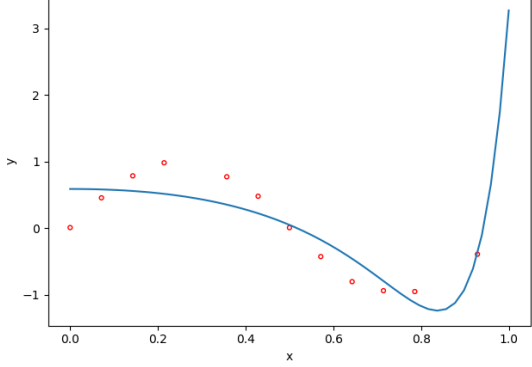
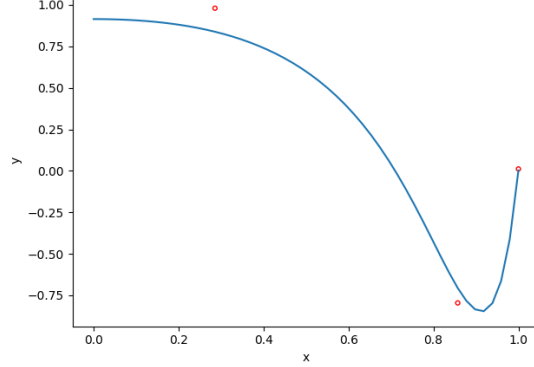
■  $\lambda = 1/15$

### Training error

Degree=14, RMSE=0.22811683626227194



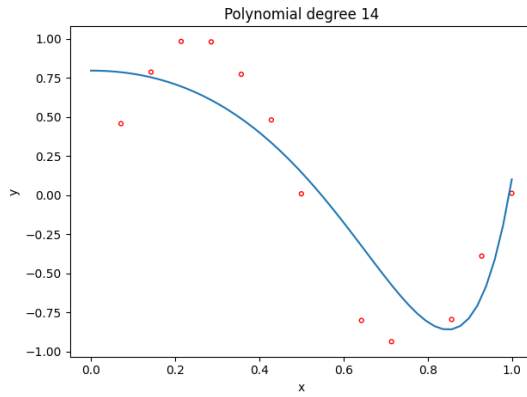
## Five-fold cross-validation errors

	training	valid
1	<p>Degree=14, RMSE=0.245748467882831</p> <p>Polynomial degree 14</p> 	<p>Degree=14, RMSE=0.13060864499556255</p> <p>Polynomial degree 14</p> 
2	<p>Degree=14, RMSE=0.2217018496391975</p> <p>Polynomial degree 14</p> 	<p>Degree=14, RMSE=0.061317219546042204</p> <p>Polynomial degree 14</p> 
3	<p>Degree=14, RMSE=0.24837169564131938</p> <p>Polynomial degree 14</p> 	<p>Degree=14, RMSE=0.09573030141252975</p> <p>Polynomial degree 14</p> 
4	<p>Degree=14, RMSE=0.21361585322383317</p> <p>Polynomial degree 14</p> 	<p>Degree=14, RMSE=0.06863846922599011</p> <p>Polynomial degree 14</p> 

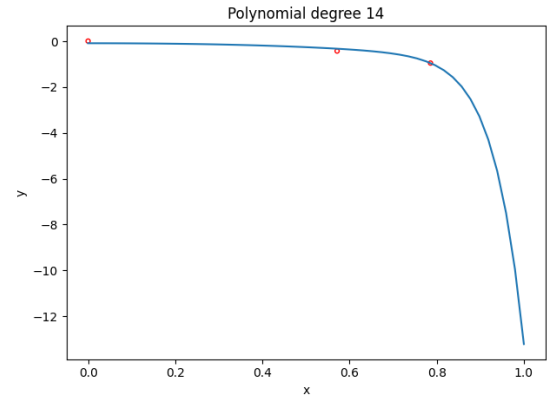


5

Degree=14, RMSE=0.19258320481981725



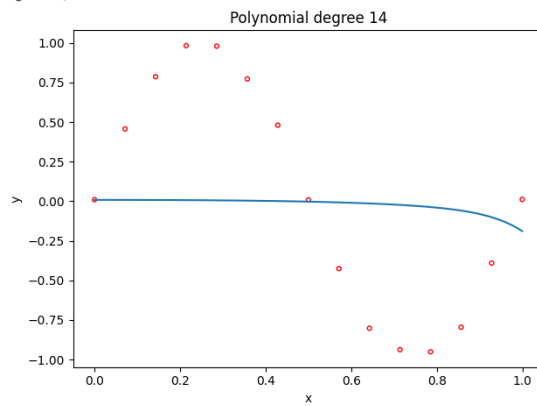
Degree=14, RMSE=0.056757140482775505



■  $\lambda = 1000/15$

### Training error

Degree=14, RMSE=0.4713576802594253

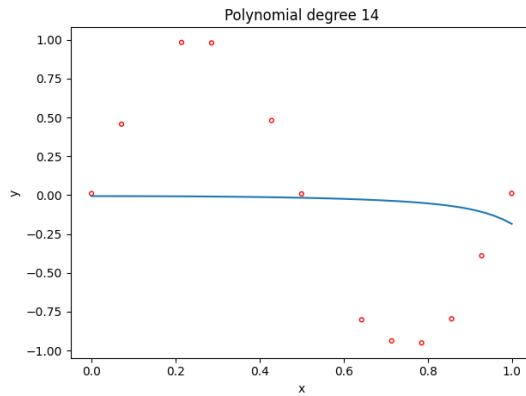


### Five-fold cross-validation errors

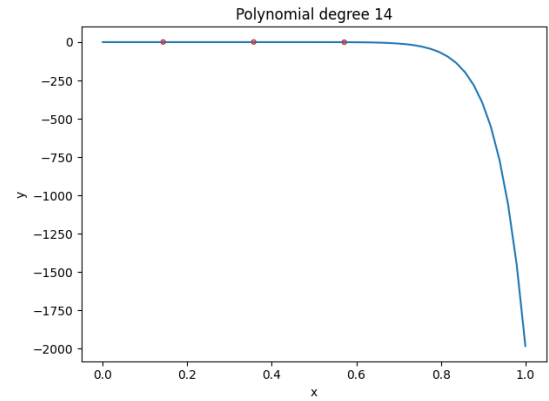
	training	valid
1	<p>Degree=14, RMSE=0.4438510777454821</p>	<p>Degree=14, RMSE=0.4300883302368353</p>

2

Degree=14, RMSE=0.4693961737558944

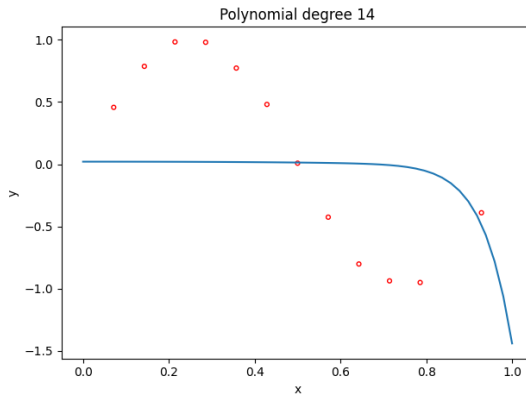


Degree=14, RMSE=0.4371590663187278

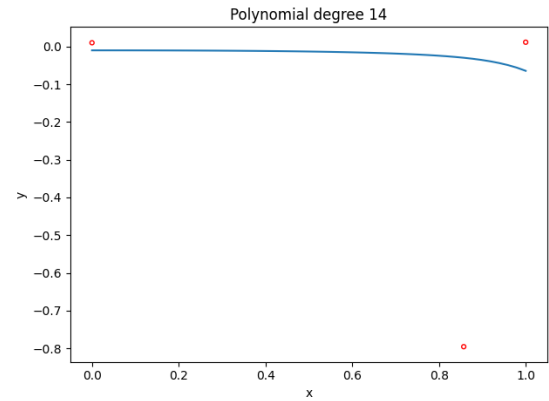


3

Degree=14, RMSE=0.49830922753483975

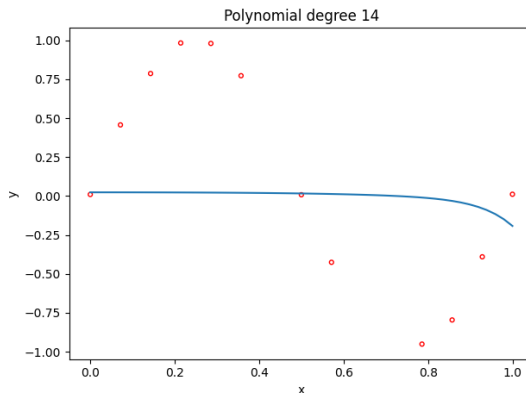


Degree=14, RMSE=0.31438418511465227

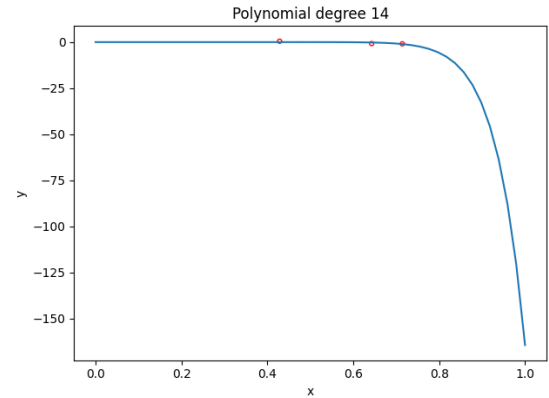


4

Degree=14, RMSE=0.45536663852178516

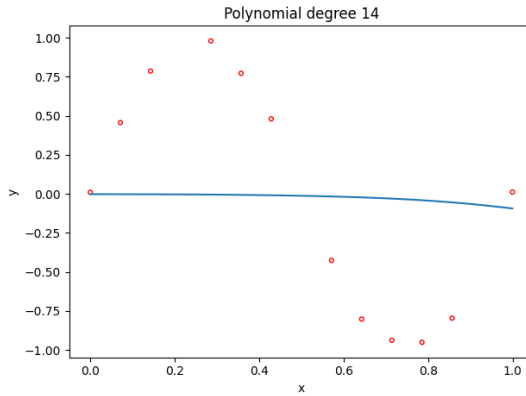


Degree=14, RMSE=0.3127302511377

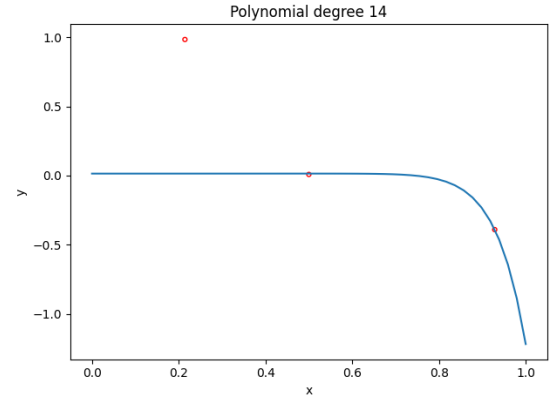


5

Degree=14, RMSE=0.4846834158948426



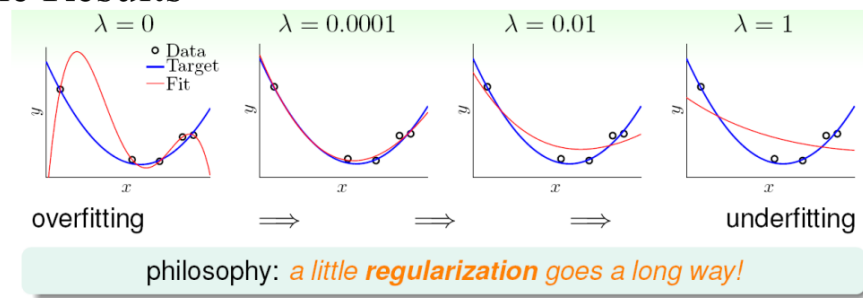
Degree=14, RMSE=0.39540147527277564



- Conclusion: The observation from your results.

- Linear regression 只適用於線性關係的資料，但當我們的資料是非線性的就必須用 polynomial regression 才能畫出最佳擬合線。龐大的資料集在做 five-fold 時，比較不會因資料量不足找不到逆矩陣而出現錯誤，但過少的資料集也可能造成訓練資料不足，因而導致 overfitting。
- 改變  $\lambda$  值， $\lambda = 0$  時模型無法彈性的做出判斷 (模型複雜度太高)，由擬合圖可看出有 overfitting 的問題； $\lambda = 1000/m$  時彈性過大因此模型無法判斷所學習的資料，由擬合圖看出有 underfitting 的問題； $\lambda = 0.001/m$  最為精準，能畫出最佳擬合線。

## The Results



實驗結果完全符合課堂簡報中的結果。

- Discussion: The questions or the difficulties you met during the implementation.

原本以為是因為點太少而一直出現找不到逆矩陣的問題

```
86
87 def _raise_linalgerror_singular(err, flag):
--> 88     raise LinAlgError("Singular matrix")
89
90 def _raise_linalgerror_nonposdef(err, flag):

LinAlgError: Singular matrix
```

將 `np.linalg.inv` 改成 `np.linalg.pinv`。

參考資料: [https://blog.csdn.net/weixin\\_43977640/article/details/109908976](https://blog.csdn.net/weixin_43977640/article/details/109908976)