

## Regression Week 3: Polynomial Regression Quiz

In this notebook you will compare different regression models in order to assess which model fits best. We will be using polynomial regression as a means to examine this topic. In particular you will:

- Write a function to take an array and a degree and return a data frame where each column is the array to a polynomial value up to the total degree.
- Use a plotting tool (e.g. matplotlib) to visualize polynomial regressions
- Use a plotting tool (e.g. matplotlib) to visualize the same polynomial degree on different subsets of the data
- Use a validation set to select a polynomial degree
- Assess the final fit using test data

## If you are doing the assignment with IPython Notebook

An IPython Notebook has been provided below to you for this quiz. This notebook contains the instructions, quiz questions and partially-completed code for you to use as well as some cells to test your code.

## What you need to download

If you are using GraphLab Create:

- Download the King County House Sales data In SFrame format: `kc_house_data.gl.zip` ([https://eventing.coursera.org/api/redirectStrict/-1pk22CEF7l8uzcZPRsJEj1xhPYfAmwqlwbhIWPAKMye35rtKMVN0KLR3xgwX2fn4GFaXLocCM3HnBzS-0kPA.NoIW2rMx3qn030-RP9C7vQ.aMu1RyFPI25JHQObjoqkMTut9sciptgQhpAyj\\_tCjHa0Pn5sM28tjv9YEn1z7uAdNr bfS0DwKjl00BzQV3bOlqkPmYmUfw\\_5FQfaA5PMbpAN9O4ZkpF9sPGRBIV33vExRucRjuXPgJ6avYLoMQNVwg-3sofoyA2Y4DKvL169w4KwCylH5sW2V2nNUec36FkBR0KcDdfqHozZHhX8NmNH\\_om0zwLeziQILaXrkADH0lZeZhpVAIkC9bOoxd9OUFe7XtLHyBi4h0ovs-Vg-4EXnqOf4rmCSbF6oqjaw-87V48ymN3BdO65jqB\\_H9gNm nj7-5dH6VGQj2YiKBmm2g9ANEGMeHK9aH855BISYzxDNezpljAFXVhpEarLyu8xiHbSeZ0U5wYUs2CaUllkqUQHul4mEcYW3lrAnsgu2mKZzcgZmYRk12XobUY8iU8KbBG](https://eventing.coursera.org/api/redirectStrict/-1pk22CEF7l8uzcZPRsJEj1xhPYfAmwqlwbhIWPAKMye35rtKMVN0KLR3xgwX2fn4GFaXLocCM3HnBzS-0kPA.NoIW2rMx3qn030-RP9C7vQ.aMu1RyFPI25JHQObjoqkMTut9sciptgQhpAyj_tCjHa0Pn5sM28tjv9YEn1z7uAdNr bfS0DwKjl00BzQV3bOlqkPmYmUfw_5FQfaA5PMbpAN9O4ZkpF9sPGRBIV33vExRucRjuXPgJ6avYLoMQNVwg-3sofoyA2Y4DKvL169w4KwCylH5sW2V2nNUec36FkBR0KcDdfqHozZHhX8NmNH_om0zwLeziQILaXrkADH0lZeZhpVAIkC9bOoxd9OUFe7XtLHyBi4h0ovs-Vg-4EXnqOf4rmCSbF6oqjaw-87V48ymN3BdO65jqB_H9gNm nj7-5dH6VGQj2YiKBmm2g9ANEGMeHK9aH855BISYzxDNezpljAFXVhpEarLyu8xiHbSeZ0U5wYUs2CaUllkqUQHul4mEcYW3lrAnsgu2mKZzcgZmYRk12XobUY8iU8KbBG))
- Download the companion IPython Notebook: `week-3-polynomial-regression-assignment-blank.ipynb` (<https://eventing.coursera.org/api/redirectStrict/2y4-3imsil7Kz6frMTBg--Jkl4Spn26xjTPoBz2o2BLTeloyKhtLDAak63jWcQacrKIKo->

hh6VxJevtokrYiXA.TtbD5KGp5M2LditlHfsVQg.U7WkDHjexlA0rSTAVLdj7Ro6S24pL9-  
 Qs\_rwwQgd3JA6AsqPIP9CL6ni6PWHfvnKYzpxevOdIV-  
 JFzJbavoa01oYtc\_qjOOPDGGxuUChZQVFD5F42Mo0Cfzrjc3y3yL9YU7CenGahlrPAhdiXTnO  
 LyjIY3IGXYa2PJe2DMQ\_7tahIEl4GQAZC173lWShTPKD944gFGN\_uzWAmDxR-  
 NSXdKsW\_bp2Sb7VyGQPfmKUS9XtsfSBjV2bceZjZ9dnZx\_6iFer0dPyviAsNSNdzzccV4RHDZ  
 NlgiNzbYEEknnOWwVL1uq7zmR\_s-aelg8tzsV-  
 Y7E8ajtVqnForwdMjcnlxtIfwEyl051okO40NI8llgqfgEVIBRROeclLkj-  
 e6wyo9WJJYnc7l\_z4AWUfnzIWWAG9z3wj9imzLUUsReu\_4wW76-  
 X\_paxy5hLH9T5QPqmTqiQx1iQsij6wa-  
 xz5\_dNgDizJcK3UuOrKV29bbRr4ncoOAgBnS3lUkxcmgaOUI4bsoEjkqux4lXblvOLZb4cTaw)

- Save both of these files in the same directory (where you are calling IPython notebook from) and unzip the data file.

If you are not using GraphLab Create:

- Download the King County House Sales data csv file: kc\_house\_data.csv  
 (https://eventing.coursera.org/api/redirectStrict/gMx1iPy2Qc66Dr65dyt61BLNjogOQr9w  
 WJv4p\_nJaQSfodP3U3llyR7YVetBw2P63-tNmQ7kB\_GKfIdal58-  
 iw.n5r6OQ0yEr1x5kby4VakJg.fBbnUSlwc4Pj53g7kvwDcEpZpUHboJhRkLBpbC-  
 pKpWBoq8iiGdKLZ9OMvpvemkrGEL4Nyh0zAVnPk1HQJOcgq7PVoaLcd-  
 F8pFZq3lkm0pDUxhRbdlQZUR-ai3xWhulkHeg5q-Q2YPO-  
 2ACTwMxb0P4WZ9J9fNvSmDfIVPEfpL84A0t99j8\_6x8MijWbAP7kZ\_wYf4ESq-WnMF\_8po-  
 6EBgCaYC\_QMjm6jVwTDbomCACKABE2on4ODSKR0x-  
 hkkaljQtm\_Eh0MWxV5Yw5ORBv8hG1ago-OrTgZ4jLuF0Gc8oUqUf-lyn1W3L-  
 uF94Tupw0fya7h1cg\_pHAQooujn8CgrscttYqHzdgo8SvPaT\_r3tPqkNzBYRID328SfZMQufbl  
 ZF4RDy7d3J1QhBZe4YswrmbBunRJD8ULJQKDMCd6fwkcS97fvSKFq5rkN-t)

NOTE: The following files are different from weeks 1 and 2

- Download the King County House Sales training data csv file:  
 wk3\_kc\_house\_train\_data.csv  
 (https://eventing.coursera.org/api/redirectStrict/URqNOBmcokRr8Anp\_Eqvg3nWITAF9\_X  
 MdP3w63q8gA948Q-  
 Mgl9nnRWAPsTCB\_YiQ5oysAdpGtVmOBIXNlidTg.Zf2QjRRhGnV4VBL01EtHUw.TPraNXSGy  
 GguJFIS\_Fb8eK6CnSFYUc4FAwgcRI2452EEcfbzKRISStKC\_ljNVtxQ66U0VN0L4PPU25hgmq0E  
 s1H1pPWSELk8Lm77EpKv6eFISmuOmYTJl0PZYWx9ox2iiskkoqRoLA73-  
 rsO07vlzMmtNOSU8TWel8cb\_sEOzouM67oe4CF-  
 r3BxB5nqQNC7rmcVcgaP46vc\_gQX0RZ0JDSH8udFgU8QVOOjBuTe8jkKpO1OF6nZWchrxA  
 HeVz0tYwofgbrmXtJlRnbk7ITKLsXMOqp-  
 BBFh11hA0\_vWqfG1Cd2BA7GSlvOb5pHoMNX8LN7x5TFFg8Af\_9UZyuyI-  
 u0zfYbaoAXuxT0MX-KCHNDTq06iVBLeCeUp8cVhllAF8RcXI2HNowZ\_9cdM7Kt-  
 0iueq\_5CL864qswValPDjSAodWn\_lplT97aW67XxNBYJ91LXfWrV6XCKz0I6NkeKdTg)
- Download the King County House Sales validation data csv file:  
 wk3\_kc\_house\_valid\_data.csv  
 (https://eventing.coursera.org/api/redirectStrict/tmgNZD\_k822F3SBWx4xKfb1zS\_sJgXEhq

gM3XLZ2EuDzAOg5SI8Im3KsMLL44tb-yon8Wp8Sj-  
 JLzbWz1gc2Jw.HbpShjZGwqZl3Ycr4bCoCg.95l1CV5lqbF7dbYgbedNF4znQtLmsijBBaW\_nq  
 4Njh6L8S21d-Ey69Qv9qDG1-x3Lxtsb7-eTx6lxv8sOAvKxDIFaTgYWM-  
 nHJewgP81o5NZ1NLSXqOM9\_OV7TgaFHzP0Jf-Yu-  
 H5gogcpaeCw5hKU\_6Fwfl5iWFovLshmnzoqp5WhONsgFWhKXCG-  
 TfHjjmwRW1PWzjVZvpmJeXgyhPT\_0cx4Nkt5R2XHAcMxMqs679ClpK1cU6liM4wQLMb2E1  
 53xNC4p5OhhMdgfctmwCoY6Su554zZqMvlj5KHCE4gX4Qij1VI0JgJXvEOA\_Vt0weXwL\_5Ln\_  
 KBFwf4aAA6HRFHvjrJSatYkl4wskf-g\_x7WpI03nOPCN5UVIsUzuIYsXHm15V\_PLUkZyggmc-  
 SBtKGdnCBG5dghCLrMOCB0xWeE-Rvm7unBx2vFgYWGT74Ry3tNmh1pwoV5i6UpLNQTW)

- Download the King County House Sales testing data csv file: wk3\_kc\_house\_test\_data.csv  
 (https://eventing.coursera.org/api/redirectStrict/2JgmCnN-TXM43xYPnmzhRLxu6XR0gaFXhrAZs7g\_LnXFOuqsD7hPPyg\_7QMmZM4pVs0jyZ-NtsmFnUxhI13Gmw.cQxWSrIhWi3EHvWbe\_K8Ag.jKhcS6C78DSLGrE3zQ9rhOHNanzgPenN7V4T3\_203SkjEbuhc3LsMcC6YQpZSigrNS6b2I9vMIGiHmVEmKr8MEv63sD6o1jUETHxdmIgw6Te8FQx8icy2vBnjpvE2jUU\_kvV8aVHTbPpDA\_b\_DkFOM2KZAOhQyfVqetwxKWzb4SjICjCi2Zxb1Ycl9dcM\_8wiCQ7dHREFnaNxxhxAQm3ID7uMmgX-AbH4UOd2vHh5XMyxQGOqgwA2B6F6vR099BLRMI4tYWfMqJ\_HAU57vDaPJ6GhNYHHR748YyuijQidPFegniEBzjy67zC5qwcSyi90DWhcYvEASHa6quE4W8m9C6aF62YvS0\_voMnAKo0rSF89H9HeQbp-BMMihUpUy4JiTU9fcZlLkIYxikO\_\_cpo-e7Qr\_xmGA\_YSruub4zjFV-xPSKcqsptjL19dAI8hOPRH6xdlya5pE5dYowqjYnbA)
- Download the King County House Sales subset 1 data csv file:  
 wk3\_kc\_house\_set\_1\_data.csv  
 (https://eventing.coursera.org/api/redirectStrict/ow13NwOesXMCwkllDeh0UaElf9X\_D7c1HnOs2NgzlUN8xPid5XwjRTsI7\_Jzboxcc4lxuZGQF0F9NZ\_LNsJtg.Kf-n3fu5eE1c7zD\_vD7daw.Tumk6Moh\_x8bn9dHJRnSiA9EkjQFIMpzdDNjvyLwai0UolcnJiKsuKTK\_bW\_HHVuDJa6syoY9PITloi\_0wKShd11ozSpkLVGEbtEQfs-g4S3DKqgk7UoiVj74ct7thGjymTvgAxtsnY6Y23b-9dtDBEKRmXKyxFeVeaEzQU2SjvCOTyGj-uAHKQCdtAJWsm8W7Lc-IVqrUDq-qzOceEwAzZdV3Cg5zxFHt7SbxeJ7UaHN9t6cAy\_trAAgWxPUC36awLLAulCiH1-ze2niY4v0uqFUfOgZTjXYPlilgDVwxq92JGuFZAw1opsvdrIR4qsl7SKyxsqEqRrbaejgX50-f1RSDLrvH1BX-tfQR7pLdb\_nQy1Tc5kljxtkLwwWXMcuNGzmur2jBYxWTW6bPFc69E10BaAnWD6fuvng8soH1qlafr6cxln0H0vj5503PMzyhSfEE-mBB9jal3BYDwg)
- Download the King County House Sales subset 2 data csv file:  
 wk3\_kc\_house\_set\_2\_data.csv  
 (https://eventing.coursera.org/api/redirectStrict/7P4F4sA1NS3WjqmK75evoM14AU9I6vrj3XQXNhYhGXf2ObWg8B\_jpEFoJDjrWa2c1B4qlMjzD-2SBP-n7gD3Gg.67aOl4wIGbjWamwWU1dDoQ.Qocwy530bUaH1d8k2wxx\_RqKgFvA2Bcu4YqTRLVQp6zEhWKOnPStIW\_hYF-ru4i7aVgho7FR0H-Ygk869PptueiECb45Z4gnD83zz18AKz9aolE65OWUzoHdkxBwaoQF50ieS8M8EE9ZbMpGj2zJJQf7xDpdnMceriGsnf4u07HAPcRNw8TT9pbgSYM7j8G-2Dq0Ja5XnXYLOKbggUMhAzJyV6Vh-mhD3SnyOfS-wQ9ns75kqTZOfCWzjkMp3-FSGn1axhXDJHQO1lyq4su4tt5CftCud4iKc9f-ewD-i\_e39QLltNzqYY3Ro1dMQ\_1af4zExk7UvT9a36e3SxEqpQtdUt0t8uyG1i\_93MNHMqZbW8ZKE5ptig8wnnysje3BgZKgfZ7COuoQZLilcn56axVyQ9bljTYZ88keV74wxQQ3gE4hfQ-

EAj5m2zJKDSGBVblaTaP6BifBKvyZELA)

- Download the King County House Sales subset 3 data csv file:  
wk3\_kc\_house\_set\_3\_data.csv  
([https://eventing.coursera.org/api/redirectStrict/65Z2NTElg7Phb9WnEtoOIUqtOqoUfjLHB DGSxu0MYZ8hh-YJ\\_CiSjyaiy8LXTbKjvIGnOmMYTIAV1PuMaos5LA.vx53WUK0wUBkl4vgqiRbdA.vcno4aQePT H2etP1ewl40LLsLITYPI6C2hl5VhnxK-SzRNyn0So-uSm7C4kJJSqEF-SEicrvWpAfwyBHzyCII8PEyxjalf7wn-e1T66\\_MncbhlLXlQdQ50cRuCv6773\\_JaikhhZoV\\_r8R4y0GwnfxOTFGw825jzYZlnSla5Ylv3rd0 OxNX5wwH1yUfryxa8e9GdiiRtP5493XQR3D4z\\_WMAD6zxWIGJOez7VR7Qs0ZX5elDOxvKrfi sHXzfg6KW2AYKHUB\\_Pi7JCgn2cChlddX5u699ApOsV2sAAh-KDAV095bcfEyFDr4v5JfVW8Ro90u66\\_wct03qoFVOMN7ukh1G8EHuy4jK1bnfDKwrep6rx5Z Hlt4o95geXzZ5eln88f2gDmSitqv5C0YofDX6A3b\\_17lei45BxAJdd8yzMW9A4m9loFWm-9UdkU-4L2f0km2dP3J1mLk6tLQCK749gg](https://eventing.coursera.org/api/redirectStrict/65Z2NTElg7Phb9WnEtoOIUqtOqoUfjLHB DGSxu0MYZ8hh-YJ_CiSjyaiy8LXTbKjvIGnOmMYTIAV1PuMaos5LA.vx53WUK0wUBkl4vgqiRbdA.vcno4aQePT H2etP1ewl40LLsLITYPI6C2hl5VhnxK-SzRNyn0So-uSm7C4kJJSqEF-SEicrvWpAfwyBHzyCII8PEyxjalf7wn-e1T66_MncbhlLXlQdQ50cRuCv6773_JaikhhZoV_r8R4y0GwnfxOTFGw825jzYZlnSla5Ylv3rd0 OxNX5wwH1yUfryxa8e9GdiiRtP5493XQR3D4z_WMAD6zxWIGJOez7VR7Qs0ZX5elDOxvKrfi sHXzfg6KW2AYKHUB_Pi7JCgn2cChlddX5u699ApOsV2sAAh-KDAV095bcfEyFDr4v5JfVW8Ro90u66_wct03qoFVOMN7ukh1G8EHuy4jK1bnfDKwrep6rx5Z Hlt4o95geXzZ5eln88f2gDmSitqv5C0YofDX6A3b_17lei45BxAJdd8yzMW9A4m9loFWm-9UdkU-4L2f0km2dP3J1mLk6tLQCK749gg))
- Download the King County House Sales subset 4 data csv file:  
wk3\_kc\_house\_set\_4\_data.csv  
([https://eventing.coursera.org/api/redirectStrict/mYxJbJoatn2pcs-vJWnkvOyPDWMK-lhko1A-xZ-HJjphsvvQqaM7W1SneQXuTwhG0IEeFh-R5urVRGJho0ZVmg.jE3f4PWLxCGI-fggnoiOgA.\\_e5ulPz5JwtGMO6O4uma4TQ4\\_\\_9Z6mSPyRTPD3kSDH\\_ssBhMvVWljwpjdiOug6 HcNLWbklObD\\_M994unmo1OVZV6JsU3bXBaKyY5YwA4oq3IEyaQF\\_jmfym-x0om6LSpl6D3Z2muOZ4t003pnRK8rhjks5yDGmg0UvTYWSHZlbnS7FdI2ANXk1aLOGFDEQ gxcNvAwipOzhvV1qPhOYp\\_Nb0hpDpS2p54r6V6Xu3DAEFenGefxoKpPeo1l6aDCI\\_1G1viC p\\_AK6LG49JWpL6ft4NbRGvymBN9PRm3jvrkYiTw8zGNZ9OU6tovON-D3C9IHqLjMTCqOK7X1tmo9ZJ4xBzKQHgclaNn8Lpw7NV6Cj6Ah0WblcNCitrqxfZRR6mdvRz WZ1Jp7idybkj9gfACUntunL7\\_KIHLnh3JzjVQqIE\\_9ZhklSzbucf-loH6xVuEH5WUal6HbwKPXJnylf\\_0A](https://eventing.coursera.org/api/redirectStrict/mYxJbJoatn2pcs-vJWnkvOyPDWMK-lhko1A-xZ-HJjphsvvQqaM7W1SneQXuTwhG0IEeFh-R5urVRGJho0ZVmg.jE3f4PWLxCGI-fggnoiOgA._e5ulPz5JwtGMO6O4uma4TQ4__9Z6mSPyRTPD3kSDH_ssBhMvVWljwpjdiOug6 HcNLWbklObD_M994unmo1OVZV6JsU3bXBaKyY5YwA4oq3IEyaQF_jmfym-x0om6LSpl6D3Z2muOZ4t003pnRK8rhjks5yDGmg0UvTYWSHZlbnS7FdI2ANXk1aLOGFDEQ gxcNvAwipOzhvV1qPhOYp_Nb0hpDpS2p54r6V6Xu3DAEFenGefxoKpPeo1l6aDCI_1G1viC p_AK6LG49JWpL6ft4NbRGvymBN9PRm3jvrkYiTw8zGNZ9OU6tovON-D3C9IHqLjMTCqOK7X1tmo9ZJ4xBzKQHgclaNn8Lpw7NV6Cj6Ah0WblcNCitrqxfZRR6mdvRz WZ1Jp7idybkj9gfACUntunL7_KIHLnh3JzjVQqIE_9ZhklSzbucf-loH6xVuEH5WUal6HbwKPXJnylf_0A))
- **IMPORTANT: use the following types for columns when importing the csv files.**  
**Otherwise, they may not be imported correctly: [str, str, float, float, float, float, int, str, int, int, int, int, int, int, int, str, float, float, float, float].** If your tool of choice requires a dictionary of types for importing csv files (e.g. Pandas), use:

```
dtype_dict = {'bathrooms':float, 'waterfront':int, 'sqft_above':int, 'sqft_living15':float, 'grade':int, 'yr_renovated':int, 'price':float, 'bedrooms':float, 'zipcode':str, 'long':float, 'sqft_lot15':float, 'sqft_living':float, 'floors':str, 'condition':int, 'lat':float, 'date':str, 'sqft_basement':int, 'yr_built':int, 'id':str, 'sqft_lot':int, 'view':int}
```

## Useful resources

You may need to install the software tools or use the free Amazon EC2 machine. Instructions for both options are provided in the reading for Module 1.

If you are following the IPython Notebook and/or are new to numpy then you might find the following tutorial helpful: [numpy-tutorial.ipynb](https://eventing.coursera.org/api/redirectStrict/Yn4ZKLBK5fVEf_HU1k_DG4E31Mailhgq7KpwIM6-s75EL3hHBGmoaFMqGQa3uCCTFOjR9ExHQ2QBBMofQ_fXuA.8ievP8byLdkELOuz8wS8BQ.pHBjxc1YEWxsnAvSAksLNRAMsa-9gmcBlm6ot_fAjeWtUKl3lYChasDXyjoSzKgSzeGoyrKbRsyLvbr6TAcSJ9vKmjHxRmmAIWMoTOE5AYOSKoSaRZzM-9Z1no-8aW5dTEmv87N3k9lmYD451etQZQmoKsxczGwmmx6GZ93fjLul1PSR-3es9XdO1yC56SnCrAkfd35eCREBjhiyKOgK74E77B7vsY97XZHBVi39hoW6l_x9Qfrl3VK_7mLzEZHV1f65rw5zilQOf14NGijHoKVCbb9DahpMiw8JosaUY9OwB7LsuzhpcQoUrDlkP273CQljcKaaIRywOi40E2Ujvn-SzzGjAsX-BYJgjYeQdO3JMtj8F6D7eqvLrRdUZRhggNLWa6ZBiQFybzX4MvdD0bb6ZLVjIVws0tIo5694uSOmP_fWCvbeHtT7tUiY4hX)

([https://eventing.coursera.org/api/redirectStrict/Yn4ZKLBK5fVEf\\_HU1k\\_DG4E31Mailhgq7KpwIM6-](https://eventing.coursera.org/api/redirectStrict/Yn4ZKLBK5fVEf_HU1k_DG4E31Mailhgq7KpwIM6-s75EL3hHBGmoaFMqGQa3uCCTFOjR9ExHQ2QBBMofQ_fXuA.8ievP8byLdkELOuz8wS8BQ.pHBjxc1YEWxsnAvSAksLNRAMsa-9gmcBlm6ot_fAjeWtUKl3lYChasDXyjoSzKgSzeGoyrKbRsyLvbr6TAcSJ9vKmjHxRmmAIWMoTOE5AYOSKoSaRZzM-9Z1no-8aW5dTEmv87N3k9lmYD451etQZQmoKsxczGwmmx6GZ93fjLul1PSR-3es9XdO1yC56SnCrAkfd35eCREBjhiyKOgK74E77B7vsY97XZHBVi39hoW6l_x9Qfrl3VK_7mLzEZHV1f65rw5zilQOf14NGijHoKVCbb9DahpMiw8JosaUY9OwB7LsuzhpcQoUrDlkP273CQljcKaaIRywOi40E2Ujvn-SzzGjAsX-BYJgjYeQdO3JMtj8F6D7eqvLrRdUZRhggNLWa6ZBiQFybzX4MvdD0bb6ZLVjIVws0tIo5694uSOmP_fWCvbeHtT7tUiY4hX)

[s75EL3hHBGmoaFMqGQa3uCCTFOjR9ExHQ2QBBMofQ\\_fXuA.8ievP8byLdkELOuz8wS8BQ.pHBjxc1YEWxsnAvSAksLNRAMsa-](https://eventing.coursera.org/api/redirectStrict/Yn4ZKLBK5fVEf_HU1k_DG4E31Mailhgq7KpwIM6-s75EL3hHBGmoaFMqGQa3uCCTFOjR9ExHQ2QBBMofQ_fXuA.8ievP8byLdkELOuz8wS8BQ.pHBjxc1YEWxsnAvSAksLNRAMsa-9gmcBlm6ot_fAjeWtUKl3lYChasDXyjoSzKgSzeGoyrKbRsyLvbr6TAcSJ9vKmjHxRmmAIWMoTOE5AYOSKoSaRZzM-9Z1no-8aW5dTEmv87N3k9lmYD451etQZQmoKsxczGwmmx6GZ93fjLul1PSR-3es9XdO1yC56SnCrAkfd35eCREBjhiyKOgK74E77B7vsY97XZHBVi39hoW6l_x9Qfrl3VK_7mLzEZHV1f65rw5zilQOf14NGijHoKVCbb9DahpMiw8JosaUY9OwB7LsuzhpcQoUrDlkP273CQljcKaaIRywOi40E2Ujvn-SzzGjAsX-BYJgjYeQdO3JMtj8F6D7eqvLrRdUZRhggNLWa6ZBiQFybzX4MvdD0bb6ZLVjIVws0tIo5694uSOmP_fWCvbeHtT7tUiY4hX)

[9gmcBlm6ot\\_fAjeWtUKl3lYChasDXyjoSzKgSzeGoyrKbRsyLvbr6TAcSJ9vKmjHxRmmAIWMoTOE5AYOSKoSaRZzM-9Z1no-](https://eventing.coursera.org/api/redirectStrict/Yn4ZKLBK5fVEf_HU1k_DG4E31Mailhgq7KpwIM6-s75EL3hHBGmoaFMqGQa3uCCTFOjR9ExHQ2QBBMofQ_fXuA.8ievP8byLdkELOuz8wS8BQ.pHBjxc1YEWxsnAvSAksLNRAMsa-9gmcBlm6ot_fAjeWtUKl3lYChasDXyjoSzKgSzeGoyrKbRsyLvbr6TAcSJ9vKmjHxRmmAIWMoTOE5AYOSKoSaRZzM-9Z1no-8aW5dTEmv87N3k9lmYD451etQZQmoKsxczGwmmx6GZ93fjLul1PSR-3es9XdO1yC56SnCrAkfd35eCREBjhiyKOgK74E77B7vsY97XZHBVi39hoW6l_x9Qfrl3VK_7mLzEZHV1f65rw5zilQOf14NGijHoKVCbb9DahpMiw8JosaUY9OwB7LsuzhpcQoUrDlkP273CQljcKaaIRywOi40E2Ujvn-SzzGjAsX-BYJgjYeQdO3JMtj8F6D7eqvLrRdUZRhggNLWa6ZBiQFybzX4MvdD0bb6ZLVjIVws0tIo5694uSOmP_fWCvbeHtT7tUiY4hX)

[8aW5dTEmv87N3k9lmYD451etQZQmoKsxczGwmmx6GZ93fjLul1PSR-](https://eventing.coursera.org/api/redirectStrict/Yn4ZKLBK5fVEf_HU1k_DG4E31Mailhgq7KpwIM6-s75EL3hHBGmoaFMqGQa3uCCTFOjR9ExHQ2QBBMofQ_fXuA.8ievP8byLdkELOuz8wS8BQ.pHBjxc1YEWxsnAvSAksLNRAMsa-9gmcBlm6ot_fAjeWtUKl3lYChasDXyjoSzKgSzeGoyrKbRsyLvbr6TAcSJ9vKmjHxRmmAIWMoTOE5AYOSKoSaRZzM-9Z1no-8aW5dTEmv87N3k9lmYD451etQZQmoKsxczGwmmx6GZ93fjLul1PSR-3es9XdO1yC56SnCrAkfd35eCREBjhiyKOgK74E77B7vsY97XZHBVi39hoW6l_x9Qfrl3VK_7mLzEZHV1f65rw5zilQOf14NGijHoKVCbb9DahpMiw8JosaUY9OwB7LsuzhpcQoUrDlkP273CQljcKaaIRywOi40E2Ujvn-SzzGjAsX-BYJgjYeQdO3JMtj8F6D7eqvLrRdUZRhggNLWa6ZBiQFybzX4MvdD0bb6ZLVjIVws0tIo5694uSOmP_fWCvbeHtT7tUiY4hX)

[3es9XdO1yC56SnCrAkfd35eCREBjhiyKOgK74E77B7vsY97XZHBVi39hoW6l\\_x9Qfrl3VK\\_7mLzEZHV1f65rw5zilQOf14NGijHoKVCbb9DahpMiw8JosaUY9OwB7LsuzhpcQoUrDlkP273CQljcKaaIRywOi40E2Ujvn-SzzGjAsX-](https://eventing.coursera.org/api/redirectStrict/Yn4ZKLBK5fVEf_HU1k_DG4E31Mailhgq7KpwIM6-s75EL3hHBGmoaFMqGQa3uCCTFOjR9ExHQ2QBBMofQ_fXuA.8ievP8byLdkELOuz8wS8BQ.pHBjxc1YEWxsnAvSAksLNRAMsa-9gmcBlm6ot_fAjeWtUKl3lYChasDXyjoSzKgSzeGoyrKbRsyLvbr6TAcSJ9vKmjHxRmmAIWMoTOE5AYOSKoSaRZzM-9Z1no-8aW5dTEmv87N3k9lmYD451etQZQmoKsxczGwmmx6GZ93fjLul1PSR-3es9XdO1yC56SnCrAkfd35eCREBjhiyKOgK74E77B7vsY97XZHBVi39hoW6l_x9Qfrl3VK_7mLzEZHV1f65rw5zilQOf14NGijHoKVCbb9DahpMiw8JosaUY9OwB7LsuzhpcQoUrDlkP273CQljcKaaIRywOi40E2Ujvn-SzzGjAsX-BYJgjYeQdO3JMtj8F6D7eqvLrRdUZRhggNLWa6ZBiQFybzX4MvdD0bb6ZLVjIVws0tIo5694uSOmP_fWCvbeHtT7tUiY4hX)

[BYJgjYeQdO3JMtj8F6D7eqvLrRdUZRhggNLWa6ZBiQFybzX4MvdD0bb6ZLVjIVws0tIo5694uSOmP\\_fWCvbeHtT7tUiY4hX](https://eventing.coursera.org/api/redirectStrict/Yn4ZKLBK5fVEf_HU1k_DG4E31Mailhgq7KpwIM6-s75EL3hHBGmoaFMqGQa3uCCTFOjR9ExHQ2QBBMofQ_fXuA.8ievP8byLdkELOuz8wS8BQ.pHBjxc1YEWxsnAvSAksLNRAMsa-9gmcBlm6ot_fAjeWtUKl3lYChasDXyjoSzKgSzeGoyrKbRsyLvbr6TAcSJ9vKmjHxRmmAIWMoTOE5AYOSKoSaRZzM-9Z1no-8aW5dTEmv87N3k9lmYD451etQZQmoKsxczGwmmx6GZ93fjLul1PSR-3es9XdO1yC56SnCrAkfd35eCREBjhiyKOgK74E77B7vsY97XZHBVi39hoW6l_x9Qfrl3VK_7mLzEZHV1f65rw5zilQOf14NGijHoKVCbb9DahpMiw8JosaUY9OwB7LsuzhpcQoUrDlkP273CQljcKaaIRywOi40E2Ujvn-SzzGjAsX-BYJgjYeQdO3JMtj8F6D7eqvLrRdUZRhggNLWa6ZBiQFybzX4MvdD0bb6ZLVjIVws0tIo5694uSOmP_fWCvbeHtT7tUiY4hX)

## If instead you are using other tools to do your homework

You are welcome, however, to write your own code and use any other libraries, like Pandas or R, to help you in the process. If you would like to take this path, follow the instructions below.

1. You're going to write a function that adds powers of a feature to columns of a data frame. For those using SFrames:

Recall that if we have an SArray 'tmp' we can get a new SArray with all the values to the third power with:

```
tmp_cubed = tmp.apply(lambda x: x**3)
```

We can create an empty SFrame with:

```
my_SFrame = graphlab.SFrame()
```

And append the tmp to it with:

```
my_SFrame['power_1'] = tmp
```

Where here 'power\_1' will refer to the power our feature was raised to.

2. Write your own function called 'polynomial\_sframe' (or otherwise) which accepts an array 'feature' and a maximal 'degree' and returns an data frame (e.g. SFrame) with the first column equal to 'feature' and the remaining columns equal to 'feature' to increasing integer powers up to 'degree'.

e.g. if you're using SFrames, you can complete the following function:

```
def polynomial_sframe(feature, degree):
    # assume that degree >= 1
    # initialize the SFrame:
    poly_sframe = graphlab.SFrame()
    # and set poly_sframe['power_1'] equal to the passed feature
    ...
    # first check if degree > 1
    if degree > 1:
        # then loop over the remaining degrees:
        for power in range(2, degree+1):
            # first we'll give the column a name:
            name = 'power_' + str(power)
            # assign poly_sframe[name] to be feature^power
            ...
    return poly_sframe
```

e.g. if you're using Pandas, you can complete the following function:

```
def polynomial_dataframe(feature, degree): # feature is pandas.Series type
    # assume that degree >= 1
    # initialize the dataframe:
    poly_dataframe = pandas.DataFrame()
    # and set poly_dataframe['power_1'] equal to the passed feature
    ...
    # first check if degree > 1
    if degree > 1:
        # then loop over the remaining degrees:
        for power in range(2, degree+1):
            # first we'll give the column a name:
            name = 'power_' + str(power)
            # assign poly_dataframe[name] to be feature^power; use apply(*)
            ...
    return poly_dataframe
```

**3.** For the remainder of the assignment we will be working with the house Sales data as in the previous notebooks. Load in the data and also sort the sales SFrame by 'sqft\_living'. When we plot the fitted values we want to join them up in a line and this works best if the variable on the X-axis (which will be 'sqft\_living') is sorted. For houses with identical square footage, we break the tie by their prices.

e.g. if you're using SFrames

```
sales = graphlab.SFrame('kc_house_data.gl/')
sales = sales.sort(['sqft_living', 'price'])
```

e.g. if you're using Pandas

```
sales = pandas.read_csv('kc_house_data.csv', dtype=dtype_dict)
sales = sales.sort(['sqft_living', 'price'])
```

4. Make a 1 degree polynomial SFrame with sales['sqft\_living'] as the the feature. Call it 'poly1\_data'.

5. Add sales['price'] to poly1\_data as this will be our output variable. e.g. if you're using SFrames

```
poly1_data = polynomial_sframe(sales['sqft_living'], 1)
poly1_data['price'] = sales['price']
```

6. Use graphlab.linear\_regression.create (or another linear regression library) to compute the regression weights for predicting sales['price'] based on the 1 degree polynomial feature 'sqft\_living'. The result should be an intercept and slope. e.g if you're using graphlab create:

```
model1 = graphlab.linear_regression.create(poly1_data, target = 'price', features = ['power_1'], validation_set = None)
```

*If you use graphlab.linear\_regression.create() to estimate these models please ensure that you set validation\_set = None. This way you will get the same answer every time you run the code.*

7. Next use the produce a scatter plot of the training data (just square feet vs price) and add the fitted model. e.g. with matplotlib and SFrames:

```
import matplotlib.pyplot as plt
%matplotlib inline
plt.plot(poly1_data['power_1'], poly1_data['price'], '.',
poly1_data['power_1'], model1.predict(poly1_data), '-')
```

The resulting plot should look like a cloud of points with a straight line passing through.

8. Now that you have plotted the results using a 1st degree polynomial, try it again using a 2nd degree and 3rd degree polynomial. Look at the fitted lines, do they appear as you would expect?

9. Now try a 15th degree polynomial. Print out the coefficients and look at the resulted fitted line. Do you think this degree is appropriate for these data? If we were to use a different subset of the data do you think we would get pretty much the same curve?

10. If you're using SFrames then create four subsets as follows:

- first split sales into 2 subsets with .random\_split(.5) use seed = 0!
- next split these into 2 more subsets (4 total) using random\_split(0.5) again set seed = 0!

- you should have 4 subsets of (approximately) equal size, call them `set_1`, `set_2`, `set_3`, and `set_4`

If you're not using SFrames then please download the provided csv files for each subset.

**11.** Estimate a 15th degree polynomial on all 4 sets, plot the results and view the coefficients for all four models.

**12. Quiz Question:** Is the sign (positive or negative) for `power_15` the same in all four models?

**13. Quiz Question:** True/False the plotted fitted lines look the same in all four plots

**14.** Since the "best" polynomial degree is unknown to us we will use cross validation to select the best degree. If you're using SFrames then create a training, validation and testing subsets as follows:

- First split sales into training\_and\_validation and testing with `sales.random_split(0.9)` use `seed = 1!`
- Next split training\_and\_validation into training and validation using `.random_split(0.5)` use `seed = 1!`

If you're not using SFrames then please download the provided csv files for training, validation and test data.

**15.** Now for each degree from 1 to 15:

- Build an polynomial data set using `training_data['sqft_living']` as the feature and the current degree
- Add `training_data['price']` as a column to your polynomial data set
- Learn a model on TRAINING data to predict 'price' based on your polynomial data set at the current degree
- Compute the RSS on VALIDATION for the current model (print or save the RSS)

*Hint: in `graphlab.linear_regression.create()` you can set `verbose = False` if you want to suppress the interim output of `linear_regression.create()`.*

**16. Quiz Question:** Which degree (1, 2, ..., 15) had the lowest RSS on Validation data?

**17.** Now that you have selected a degree compute the RSS on TEST data for the model with the best degree from the Validation data.

**18. Quiz Question:** what is the RSS on TEST data for the model with the degree selected from Validation data? (Make sure you got the correct degree from the previous question)



