

CPSC 340 Assignment 3 (Due 2021-10-15 at 11:59pm)

Important: Submission Format [5 points]

Please make sure to follow the submission instructions posted on the course website. We will deduct marks if the submission format is incorrect, or if you're not using \LaTeX and your handwriting is *at all* difficult to read – at least these 5 points, more for egregious issues. Compared to assignment 1, your name and student number are no longer necessary (though it's not a bad idea to include them just in case, especially if you're doing the assignment with a partner).

1 Matrix Notation and Minimizing Quadratics [12 points]

1.1 Converting to Matrix/Vector/Norm Notation [6 points]

Using our standard supervised learning notation (X, y, w) express the following functions in terms of vectors, matrices, and norms (there should be no summations or maximums).

1. $\max_{i \in \{1, 2, \dots, n\}} |w^T x_i - y_i|$. This is “brittle regression”.

Answer: $\|Xw - y\|_\infty$

2. $\sum_{i=1}^n v_i (w^T x_i - y_i)^2 + \frac{\lambda}{2} \sum_{j=1}^d w_j^2$. This is regularized least squares with a *weight* v_i for each training example: Hint: You can use V to denote a diagonal matrix that has the values v_i along the diagonal. What does $a^T V b$ look like in summation form (for some arbitrary vectors a, b)?

Answer: $\|V^{\frac{1}{2}}(Xw - y)\|_2^2 + \frac{\lambda}{2} \|w\|_2^2$

3. $(\sum_{i=1}^n |w^T x_i - y_i|)^2 + \frac{1}{2} \sum_{j=1}^d \lambda_j |w_j|$. This is L1-regularized least squares with a different regularization strength for each dimension: Hint: You can use Λ to denote a diagonal matrix that has the λ_j values along the diagonal.

Answer: $\|Xw - y\|_1^2 + \frac{1}{2} \|\Lambda w\|_1$

Note: you can assume that all the v_i and λ_i values are non-negative.

1.2 Minimizing Quadratic Functions as Linear Systems [6 points]

Write finding a minimizer w of the functions below as a system of linear equations (using vector/matrix notation and simplifying as much as possible). Note that all the functions below are convex, so finding a w with $\nabla f(w) = 0$ is sufficient to minimize the functions – but show your work in getting to this point.

Answer: References:

http://www.gatsby.ucl.ac.uk/teaching/courses/sntn/sntn-2017/resources/Matrix_derivatives_cribsheet.pdf

1. $f(w) = \frac{1}{2} \|w - v\|^2$ (projection of v onto real space).

Answer: $f(w) = \frac{1}{2} \|w\|^2 - 2 * \frac{1}{2} * w^T v + \frac{1}{2} \|v\|^2$

Thus, we can find that $\nabla f(w) = w - v$. Therefore, $w = v$

2. $f(w) = \frac{1}{2}\|Xw - y\|^2 + \frac{1}{2}w^T\Lambda w$ (least squares with weighted regularization).

Answer: $f(w) = \frac{1}{2}\|Xw\|^2 - 1/2 * 2 * w^T X^T y + \frac{1}{2}\|y\|^2 + \frac{1}{2}w^T\Lambda w$

Thus, we can find that $\nabla f(w) = X^T Xw - X^T y + \Lambda w$

Now set $\nabla f(w) = 0$, which gives us:

$$(X^T X + \Lambda I)w = X^T y$$

Now solve for w .

3. $f(w) = \frac{1}{2} \sum_{i=1}^n v_i (w^T x_i - y_i)^2 + \frac{\lambda}{2} \|w - w^{(0)}\|^2$ (weighted least squares shrunk towards non-zero $w^{(0)}$).

Above we assume that v and $w^{(0)}$ are $d \times 1$ vectors, and Λ is a $d \times d$ diagonal matrix (with positive entries along the diagonal). You can use V as a diagonal matrix containing the v_i values along the diagonal.

Hint: Once you convert to vector/matrix notation, you can use the results from class to quickly compute these quantities term-wise. As a spot check, make sure that the dimensions match for all quantities/operations: to do this, you may need to introduce an identity matrix. For example, $X^T Xw + \lambda w$ can be re-written as $(X^T X + \lambda I)w$.

Answer: First, convert to norm notation:

$$f(w) = \frac{1}{2} \|V^{\frac{1}{2}}(Xw - y)\|^2 + \frac{\lambda}{2} \|w - w^{(0)}\|^2$$

Then expand it out (converted first part to matrix notation).

$$f(w) = \frac{1}{2} w^T X^T V X w - \frac{2*1}{2} w^T X^T V y + \frac{1}{2} y^T V y + \frac{\lambda}{2} \|w\|^2 - 2 * \lambda/2 * w^T w^{(0)} + \frac{\lambda}{2} \|w^{(0)}\|^2$$

Now take the derivative:

$$\nabla f(w) = X^T V X w - X^T V y + \lambda w - \lambda w^{(0)}.$$

Set $\nabla f(w) = 0$:

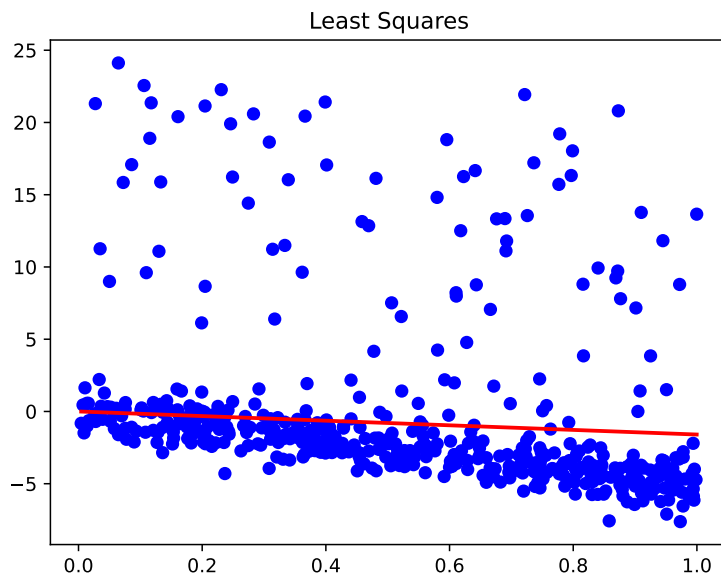
$$X^T V X w + \lambda w = X^T V y + \lambda w^{(0)}$$

$$(X^T V X + \lambda I)w = X^T V y + \lambda w^{(0)}$$

Now solve for w .

2 Robust Regression and Gradient Descent [41 points]

If you run `python main.py -q 2`, it will load a one-dimensional regression dataset that has a non-trivial number of ‘outlier’ data points. These points do not fit the general trend of the rest of the data, and pull the least squares model away from the main downward trend that most data points exhibit:



Note: we are fitting the regression without an intercept here, just for simplicity of the homework question. In reality one would rarely do this. But here it’s OK because the “true” line passes through the origin (by design). In Q3.1 we’ll address this explicitly.

A coding note: when we’re doing math, we always treat y and w as column vectors, i.e. if we’re thinking of them as matrices, then shape $n \times 1$ or $d \times 1$, respectively. This is also what you’d usually do when coding things in, say, Matlab. It is *not* what’s usually done in Python machine learning code, though: we usually have `y.shape == (n,)`, i.e. a one-dimensional array. Mathematically, these are the same thing, but if you mix between the two, you can really easily get confusing answers: if you add something of shape `(n, 1)` to something of shape `(n,)`, then the NumPy broadcasting rules give you something of shape `(n, n)`. This is a very unfortunate consequence of the way the broadcasting rules work. If you stick to either one, you generally don’t have to worry about it; **we’re assuming shape `(n,)` here**. Note that you can go from something of shape `(n, 1)` to `(n,)` by, for example, using `two_d_array.squeeze(1)` (which checks that the axis at index 1, the second one, is length 1 and then removes it). You can go in the other direction with, for instance, `one_d_array[:, np.newaxis]` (which says “give me the whole first axis, then add another axis of length 1 in the second position”).

2.1 Weighted Least Squares in One Dimension [8 points]

One of the most common variations on least squares is *weighted* least squares. In this formulation, we have a weight v_i for every training example. To fit the model, we minimize the weighted squared error,

$$f(w) = \frac{1}{2} \sum_{i=1}^n v_i (w^T x_i - y_i)^2.$$

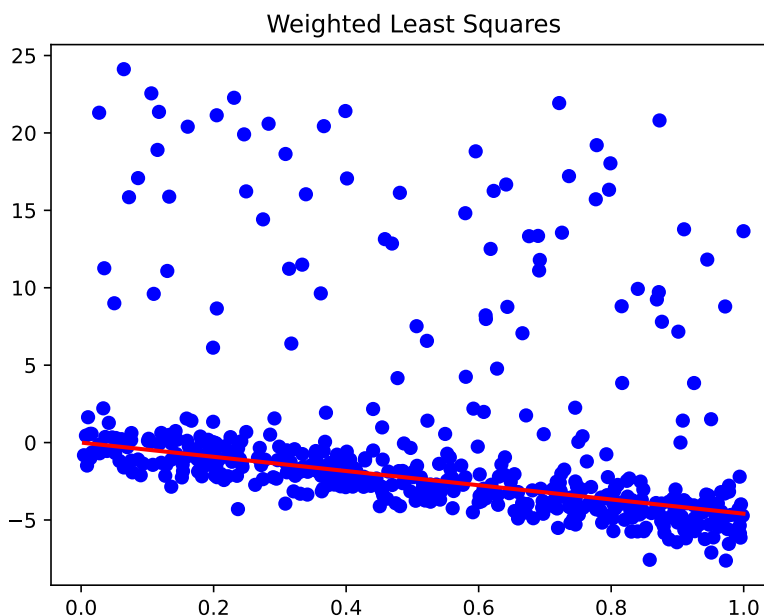
In this formulation, the model focuses on making the error small for examples i where v_i is high. Similarly, if v_i is low then the model allows a larger error. Note: these weights v_i (one per training example) are completely different from the model parameters w_j (one per feature), which, confusingly, we sometimes also call “weights.” The v_i are sometimes called *sample weights* or *instance weights* to help distinguish them.

Complete the model class, `WeightedLeastSquares` (inside `linear_models.py`), to implement this model. (Note that Q1.2.3 asks you to show how a similar formulation can be solved as a linear system.) Apply this model to the data containing outliers, setting $v = 1$ for the first 400 data points and $v = 0.1$ for the last 100 data points (which are the outliers).

Answer:

```
1 def fit(self, X, y, v):
2     """YOUR CODE HERE FOR Q2.1"""
3     weighted_X = v*X
4     weighted_y = v*y
5     self.w = solve(X.T @ weighted_X, X.T @ weighted_y)
```

Hand in your code and the updated plot.



2.2 Smooth Approximation to the L1-Norm [8 points]

Unfortunately, we typically do not know the identities of the outliers. In situations where we suspect that there are outliers, but we do not know which examples are outliers, it makes sense to use a loss function that is more robust to outliers. In class, we discussed using the sum of absolute values objective,

$$f(w) = \sum_{i=1}^n |w^T x_i - y_i|.$$

This is less sensitive to outliers than least squares, but it is non-differentiable and harder to optimize. Nevertheless, there are various smooth approximations to the absolute value function that are easy to

optimize. One possible approximation is to use the log-sum-exp approximation of the max function¹:

$$|r| = \max\{r, -r\} \approx \log(\exp(r) + \exp(-r)).$$

Using this approximation, we obtain an objective of the form

$$f(w) = \sum_{i=1}^n \log(\exp(w^T x_i - y_i) + \exp(y_i - w^T x_i)).$$

which is smooth but less sensitive to outliers than the squared error. Derive the gradient ∇f of this function with respect to w . You should show your work but you do not have to express the final result in matrix notation.

Answer: We take a derivative of r :

$$\frac{d}{dr} [\log(\exp(r) + \exp(-r))] = \frac{\exp(r) - \exp(-r)}{\exp(r) + \exp(-r)}$$

Then, substitute $w^T x_i - y_i$ into r :

$$\nabla f = \frac{df(w)}{dw_j} = \sum_{i=1}^n x_{ij} \frac{(\exp(w^T x_i - y_i) - \exp(y_i - w^T x_i))}{(\exp(w^T x_i - y_i) + \exp(y_i - w^T x_i))}$$

2.3 Gradient Descent: Understanding the Code [5 points]

Recall gradient descent, a derivative-based optimization algorithm that uses gradients to navigate the parameter space until a locally optimal parameter is found. In `optimizers.py`, you will see our implementation of gradient descent, taking the form of a class named `OptimizerGradientDescent`. This class has a similar design pattern as PyTorch, a popular differentiable programming and optimization library. One step of gradient descent is defined as

$$w^{t+1} = w^t - \alpha^t \nabla_w f(w^t).$$

Look at the methods named `get_learning_rate_and_step()` and `break_yes()`, and answer each of these questions, one sentence per answer:

1. Which variable is equivalent to α^t , the step size at iteration t ?

Answer: α^t is equivalent to `alpha`

2. Which variable is equivalent to $\nabla_w f(w^t)$ the current value of the gradient vector?

Answer: $\nabla_w f(w^t)$ is equivalent to `g_old`

3. Which variable is equivalent to w^t , the current value of the parameters?

Answer: w^t is equivalent to `w_old`

4. What is the method `break_yes()` doing?

Answer: This method returns true when the `gradient_norm` becomes smaller than optimization tolerance or number of function evaluation reaches the designated maximum number.

2.4 Robust Regression [20 points]

The class `LinearModelGradientDescent` is the same as `LeastSquares`, except that it fits the least squares model using a gradient descent method. If you run `python main.py -q 2.4` you'll see it produces the same fit as we obtained using the normal equations.

The typical input to a gradient method is a function that, given w , returns $f(w)$ and $\nabla f(w)$. See `funObj` in `LinearModelGradientDescent` for an example. Note that the `fit` function of `LinearModelGradientDescent`

¹Other possibilities are the Huber loss, or $|r| \approx \sqrt{r^2 + \epsilon}$ for some small ϵ .

also has a numerical check that the gradient code is approximately correct, since implementing gradients is often error-prone.²

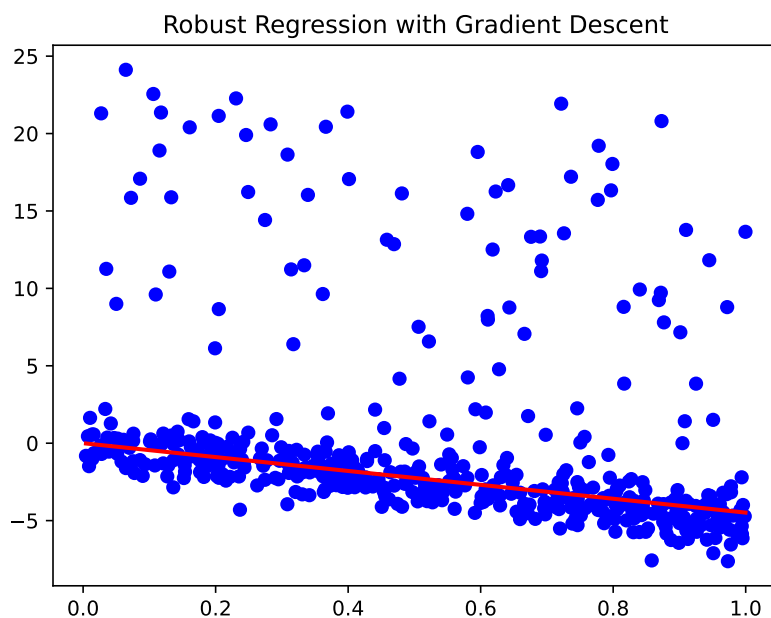
An advantage of gradient-based strategies is that they are able to solve problems that do not have closed-form solutions, such as the formulation from the previous section. The class `LinearModelGradientDescent` has most of the implementation of a gradient-based strategy for fitting the robust regression model under the log-sum-exp approximation.

2.4.1 Implementing the Objective Function [15 points]

Optimizing robust regression parameters is the matter of implementing a function object and using an optimizer to minimize the function object. The only part missing is the function and gradient calculation inside `fun_obj.py`. Inside `fun_obj.py`, complete `FunObjRobustRegression` to implement the objective function and gradient based on the smooth approximation to the absolute value function (from the previous section). Hand in your code, as well as the plot obtained using this robust regression approach.

Answer:

```
1 """YOUR CODE HERE FOR Q2.3"""
2     f = np.sum(np.log(np.exp(X @ w - y) + np.exp(y - X @ w)))
3     g_numerator = np.exp(X @ w - y) - np.exp(y - X @ w)
4     g_denominator = np.exp(X @ w - y) + np.exp(y - X @ w)
5     g = X.T @ (g_numerator / g_denominator)
6     return f, g
```

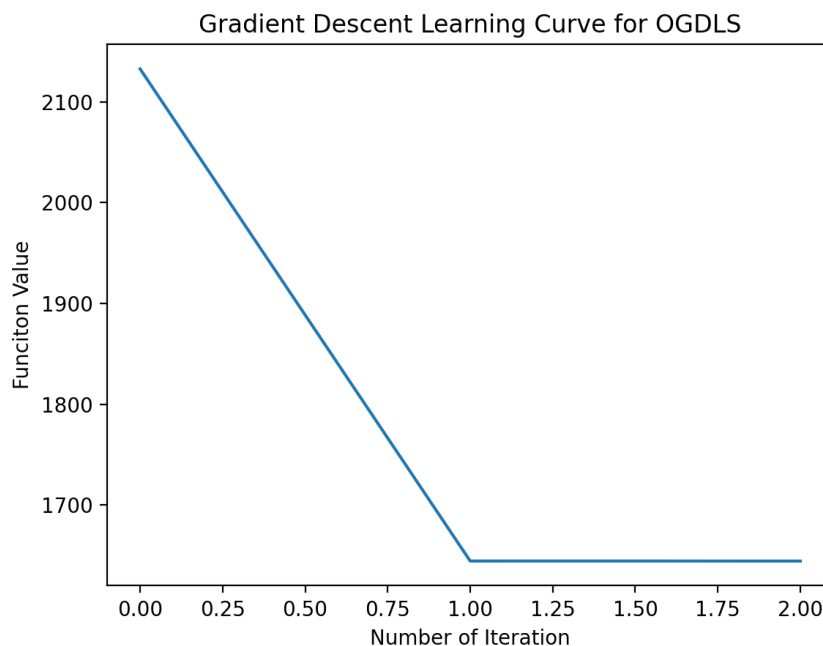
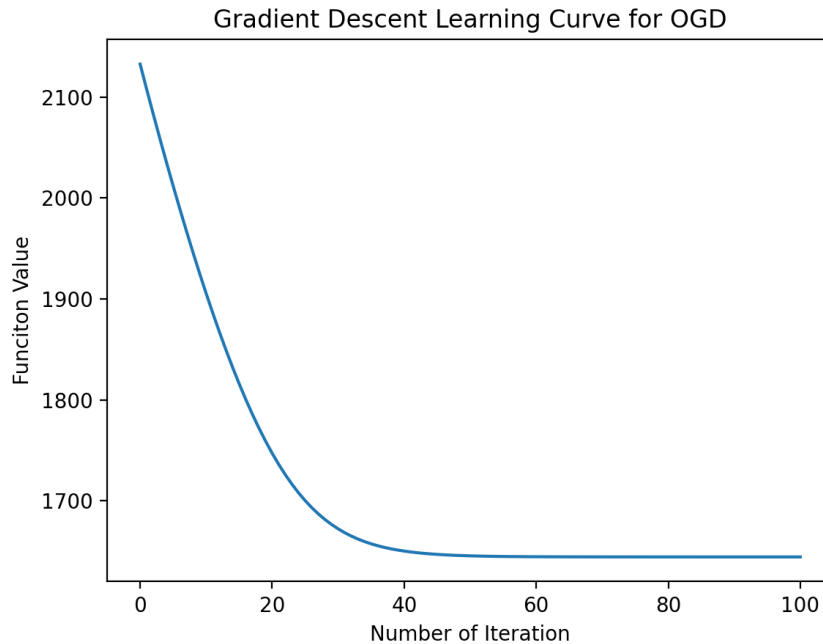


2.4.2 The Learning Curves [5 points]

Using the same dataset as the previous sections, produce the plot of “gradient descent learning curves” to compare the performances of `OptimizerGradientDescent` and `OptimizerGradientDescentLineSearch` for

²Sometimes the numerical gradient checker itself can be wrong. See CPSC 303 for a lot more on numerical differentiation.

robust regression, where **one hundred (100) iterations** of gradient descent are on the x-axis and the **objective function value** corresponding to each iteration is visualized on the y-axis (see gradient descent lecture). Use the default `learning_rate` for `OptimizerGradientDescent`. [Submit this plot.](#) According to this plot, which optimizer is more “iteration-efficient”?



Answer: `OptimizerGradientDescentLineSearch` is more efficient since they have less iterations to reach the

minimum function value.

3 Linear Regression and Nonlinear Bases

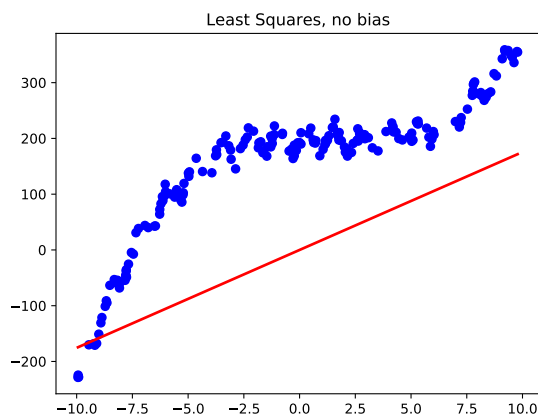
In class we discussed fitting a linear regression model by minimizing the squared error. In this question, you will start with a data set where least squares performs poorly. You will then explore how adding a bias variable and using nonlinear (polynomial) bases can drastically improve the performance. You will also explore how the complexity of a basis affects both the training error and the validation error.

3.1 Adding a Bias Variable [8 points]

If you run `python main.py -q 3`, it will:

1. Load a one-dimensional regression dataset.
2. Fit a least-squares linear regression model.
3. Report the training error.
4. Report the validation error.
5. Draw a figure showing the training data and what the linear model looks like.

Unfortunately, this is an awful model of the data. The average squared training error on the data set is over 28000 (as is the validation error), and the figure produced by the demo confirms that the predictions are usually nowhere near the training data:



The y -intercept of this data is clearly not zero (it looks like it's closer to 200), so we should expect to improve performance by adding a *bias* (a.k.a. intercept) variable, so that our model is

$$y_i = w^T x_i + w_0.$$

instead of

$$y_i = w^T x_i.$$

In file `linear_models.py`, complete the class `LeastSquaresBias`, that has the same input/model/predict format as the `LeastSquares` class, but that adds a *bias* variable (also called an intercept) w_0 (also called β in lecture). Hand in your new class, the updated plot, and the updated training/validation error.

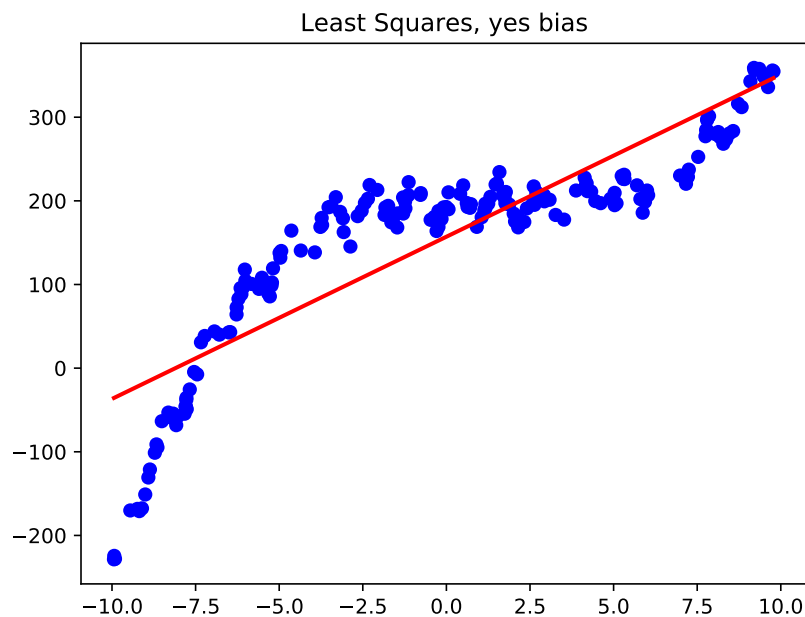
Hint: recall that adding a bias w_0 is equivalent to adding a column of ones to the matrix X . Don't forget that you need to do the same transformation in the `predict` function.

Answer:

```
1 class LeastSquaresBias:
2     "Least Squares with a bias added"
3     def fit(self, X, y):
4         """YOUR CODE HERE FOR Q3.1"""
5         n, d = X.shape
6         new_column = np.ones(n)
7         Z = np.insert(X, 0, new_column, axis=1)
8         self.v = solve(Z.T @ Z, Z.T @ y)
9
10    def predict(self, X_pred):
11
12        """YOUR CODE HERE FOR Q3.1"""
13        n, d = X_pred.shape
14        new_column = np.ones(n)
15        Z_pred = np.insert(X_pred, 0, new_column, axis=1)
16        return Z_pred @ self.v
```

Training error = 3551.3

Validation error = 3393.9



3.2 Polynomial Basis [10 points]

Adding a bias variable improves the prediction substantially, but the model is still problematic because the target seems to be a *non-linear* function of the input. Complete `LeastSquarePoly` class, that takes a data vector x (i.e., assuming we only have one feature) and the polynomial order p . The function should perform a least squares fit based on a matrix Z where each of its rows contains the values $(x_i)^j$ for $j = 0$ up to p .

E.g., `LeastSquaresPoly.fit(x,y)` with $p = 3$ should form the matrix

$$Z = \begin{bmatrix} 1 & x_1 & (x_1)^2 & (x_1)^3 \\ 1 & x_2 & (x_2)^2 & (x_2)^3 \\ \vdots & & & \\ 1 & x_n & (x_n)^2 & (x_n)^3 \end{bmatrix},$$

and fit a least squares model based on it. Submit your code, and a plot showing training and validation error curves for the following values of p : 0, 1, 2, 3, 4, 5, 10, 20, 30, 50, 75, 100. Clearly label your axes, and use a logarithmic scale for y by `plt.yscale("log")` or similar, so that we can still see what's going on if there are a few extremely large errors. Explain the effect of p on the training error and on the validation error.

NOTE: large values of p may cause numerical instability. Your solution may look different from others' even with the same code depending on the OS and other factors. As long as your training and validation error curves behave as expected, you will not be penalized.

Note: you should write the code yourself; don't use a library like sklearn's `PolynomialFeatures`.

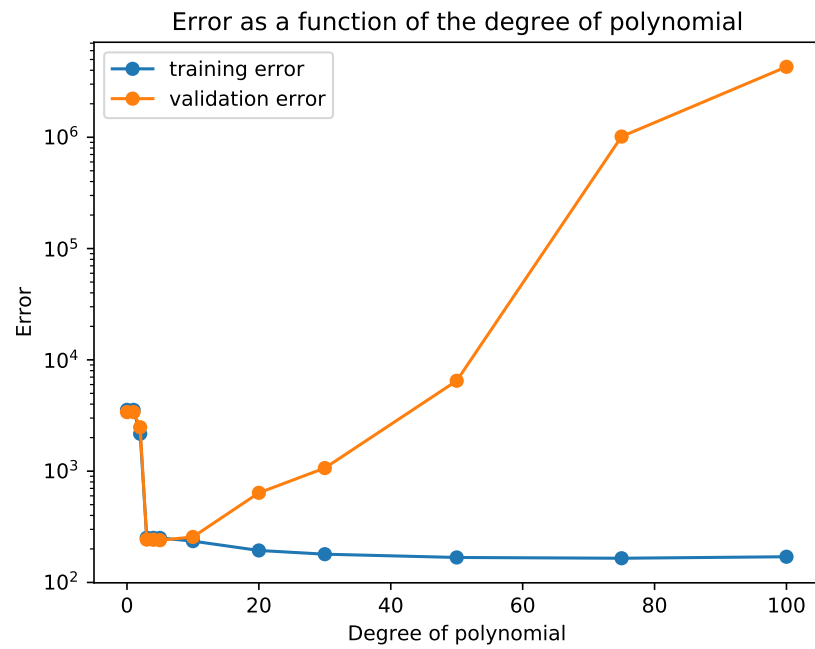
Note: in addition to the error curves, the code also produces a plot of the fits themselves. This is for your information; you don't have to submit it.

Answer:

```

1 class LeastSquaresPoly:
2     "Least Squares with polynomial basis"
3     def __init__(self, p):
4         self.leastSquares = LeastSquares()
5         self.p = p
6
7     def fit(self, X, y):
8         """YOUR CODE HERE FOR Q3.2"""
9         Z = self._poly_basis(X)
10        self.v = solve(Z.T @ Z, Z.T @ y)
11
12
13    def predict(self, X_pred):
14        """YOUR CODE HERE FOR Q3.2"""
15        Z_pred = self._poly_basis(X_pred)
16        return Z_pred @ self.v
17
18    # A private helper function to transform any matrix X into
19    # the polynomial basis defined by this class at initialization
20    # Returns the matrix Z that is the polynomial basis of X.
21    def _poly_basis(self, X):
22        n, d = X.shape
23        for i in range(self.p - 1): # doing minus one because we already have the
24            ↪ original column of X
25            poly_column = np.ones(n)
26            for j in range(n):
27                poly_column[j] = X[j][0] ** (i + 2)
28            X = np.insert(X, i + 1, poly_column, axis=1)
29        new_column = np.ones(n)
30        Z = np.insert(X, 0, new_column, axis=1)
31        return Z

```



The first few error seem to be the same, but as p got larger, the training error seemed to plateau whereas the validation error started increasing.

4 Very-Short Answer Questions [24 points]

Answer the following questions (in a sentence or two).

1. Suppose that a training example is global outlier, meaning it is really far from all other data points. How is the cluster assignment of this example by k -means? And how is it set by density-based clustering?

Answer: For k -means, the outlier is assigned to a cluster, but a human looking at it can tell that it's far away from the rest of the group. Outliers can also appear as a smaller cluster of points that don't belong to any other group.

For density-based clustering, the outliers do not get assigned to clusters.

2. Why do we need random restarts for k -means but not for density-based clustering?

Answer: K -means is affected by its initializing values, so it can get "stuck" at a certain cluster configuration even though it may not be ideal. For density-based clustering, we end up at the same clusters no matter how we start.

3. Can hierarchical clustering find non-convex clusters?

Answer: Yes, because it finds clusters by comparing which points are closest to each other.

4. For model-based outlier detection, list an example method and problem with identifying outliers using this method.

Answer: One method is to fit a probabilistic model and consider the outliers to be the values with low probabilities. The issue with this is that the z -score is also skewed by outliers, which defeats the whole purpose of using it.

5. For graphical-based outlier detection, list an example method and problem with identifying outliers using this method.

Answer: An example would be to generate a boxplot to visualize the data. The problem with this method is that you can only use it for 1 variable.

6. For supervised outlier detection, list an example method and problem with identifying outliers using this method.

Answer: We could label a dataset as $y_i = 0$ if it's not an outlier, and $y_i = 1$ if it is an outlier. The problem with this is that it needs supervision, and that it may not detect new types of outliers that it hasn't seen before.

7. If we want to do linear regression with 1 feature, explain why it would or would not make sense to use gradient descent to compute the least squares solution.

Answer: It wouldn't make sense to use gradient descent for linear regressions with 1 feature because gradient descent is used to find the global minimum of convex functions, but for linear regression, there is no global minimum so the algorithm would run infinitely.

8. Why do we typically add a column of 1 values to X when we do linear regression? Should we do this if we're using decision trees?

Answer: We add a column of 1 values to account for the bias, or the intercept.

No, we shouldn't do this with decision trees, because there is no algebra done with it so it would be meaningless.

9. Why do we need gradient descent for the robust regression problem, as opposed to just using the normal equations? Hint: it is NOT because of the non-differentiability. Recall that we used gradient descent even after smoothing away the non-differentiable part of the loss.

Answer: Because we don't have normal equations with robust regression, so it's more difficult to compute the gradient with the Huber-loss.

10. What is the problem with having too small of a learning rate in gradient descent? What is the problem with having too large of a learning rate in gradient descent?

Answer: The problem of having a learning rate that's too small is that our algorithm would have to run many iterations in order to find the minimum, so it would cost a lot of time. The problem of having a learning rate that's too large is that it's possible that the algorithm could overshoot and miss the minimum, which means that it would also cost more time to find the minimum.

11. What is the purpose of the log-sum-exp function and how is this related to gradient descent?

Answer: The purpose of the log-sum-exp function is to do a smooth approximation when using brittle regression. This is so that we can perform gradient descent to minimize the maximum error (comparable to using Huber-loss for robust regression).

12. What type of non-linear transform might be suitable if we had a periodic function?

Answer: If we had a periodic function, we could add a column with $\sin(x_i)$ values instead of 1's.