

# CPSC 340 Assignment 3 (due Monday, Feb 10 at 11:55pm)

Chia Xiang Rong  
53513537  
Zachary Chua  
48995336

2020-02-09

## Instructions

Rubric: {mechanics:5}

**IMPORTANT!!! Before proceeding, please carefully read the general homework instructions at <https://www.cs.ubc.ca/~fwood/CS340/homework/>.** The above 5 points are for following the submission instructions. You can ignore the words “mechanics”, “reasoning”, etc.

We use [blue](#) to highlight the deliverables that you must answer/do/submit with the assignment.

## 1 Finding Similar Items

For this question we’ll use the Amazon product data set<sup>1</sup> from <http://jmcauley.ucsd.edu/data/amazon/>. We will focus on the “Patio, Lawn, and Garden” section. You should start by downloading the ratings at <https://stanford.io/2Q7QTVu> and place the file in your **data** directory with the original filename. Once you do that, running `python main.py -q 1` should do the following steps:

- Load the raw ratings data set into a Pandas dataframe.
- Construct the user-product matrix as a sparse matrix (to be precise, a `scipy.sparse.csr_matrix`).
- Create bi-directional mappings from the user ID (e.g. “A2VNYWOPJ13AFP”) to the integer index into the rows of **X**.
- Create bi-directional mappings from the item ID (e.g. “0981850006”) to the integer index into the columns of **X**.

### 1.1 Exploratory data analysis

#### 1.1.1 Most popular item

Rubric: {code:1}

Find the item with the most total stars. [Submit the product name and the number of stars.](#)

---

<sup>1</sup>The author of the data set has asked for the following citations: (1) Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. R. He, J. McAuley. WWW, 2016, and (2) Image-based recommendations on styles and substitutes. J. McAuley, C. Targett, J. Shi, A. van den Hengel. SIGIR, 2015.

Note: once you find the ID of the item, you can look up the name by going to the url [https://www.amazon.com/dp/ITEM\\_ID](https://www.amazon.com/dp/ITEM_ID), where ITEM\_ID is the ID of the item. For example, the URL for item ID “B00CFM0P7Y” is <https://www.amazon.com/dp/B00CFM0P7Y>.

Answer:

Product name: Classic Accessories Veranda Grill Cover, X-Large

Number of stars: 14454

### 1.1.2 User with most reviews

Rubric: {code:1}

Find the user who has rated the most items, and the number of items they rated.

Answer:

User ID: A09532783OZ58XW6TJSEP

Number of reviews: 161

### 1.1.3 Histograms

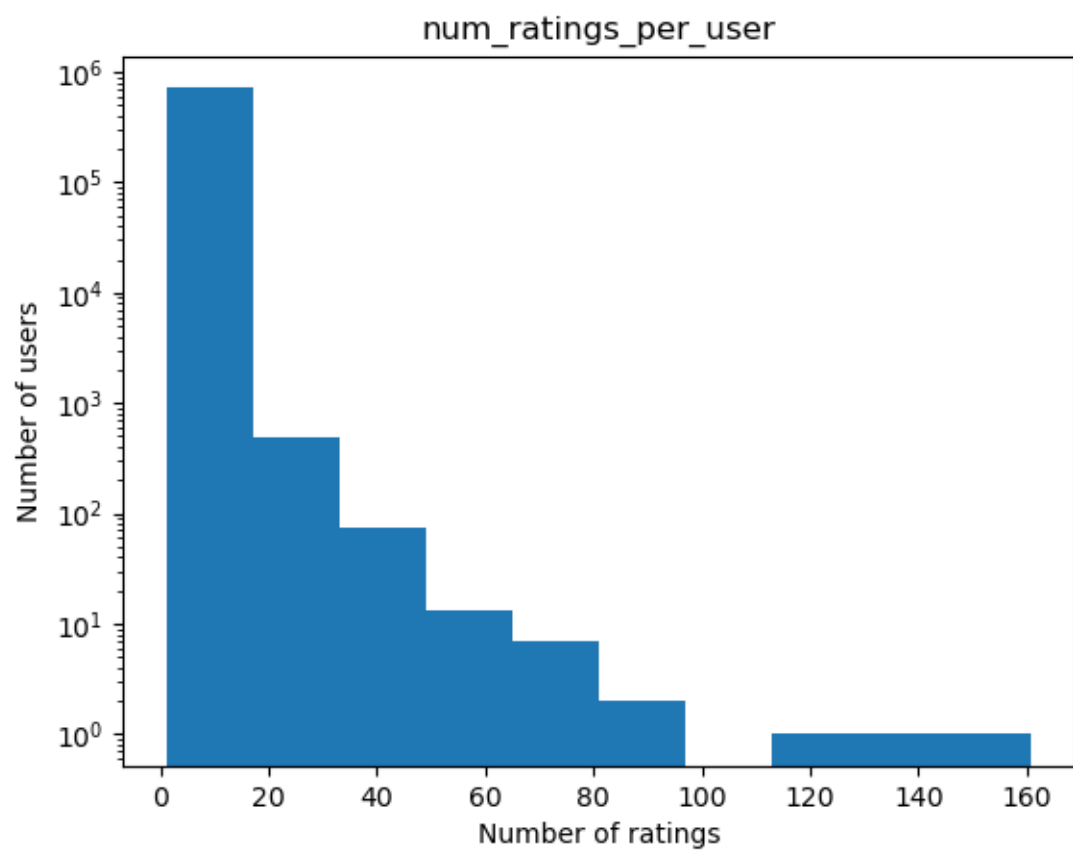
Rubric: {code:2}

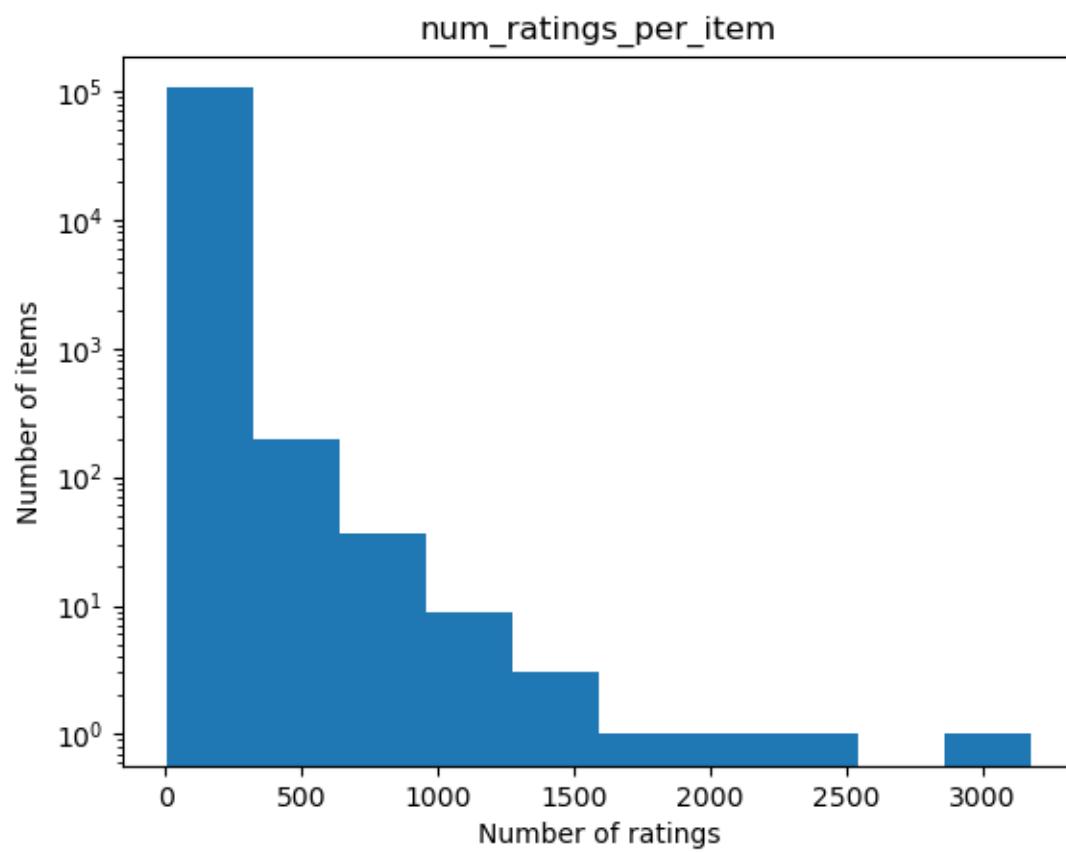
Make the following histograms:

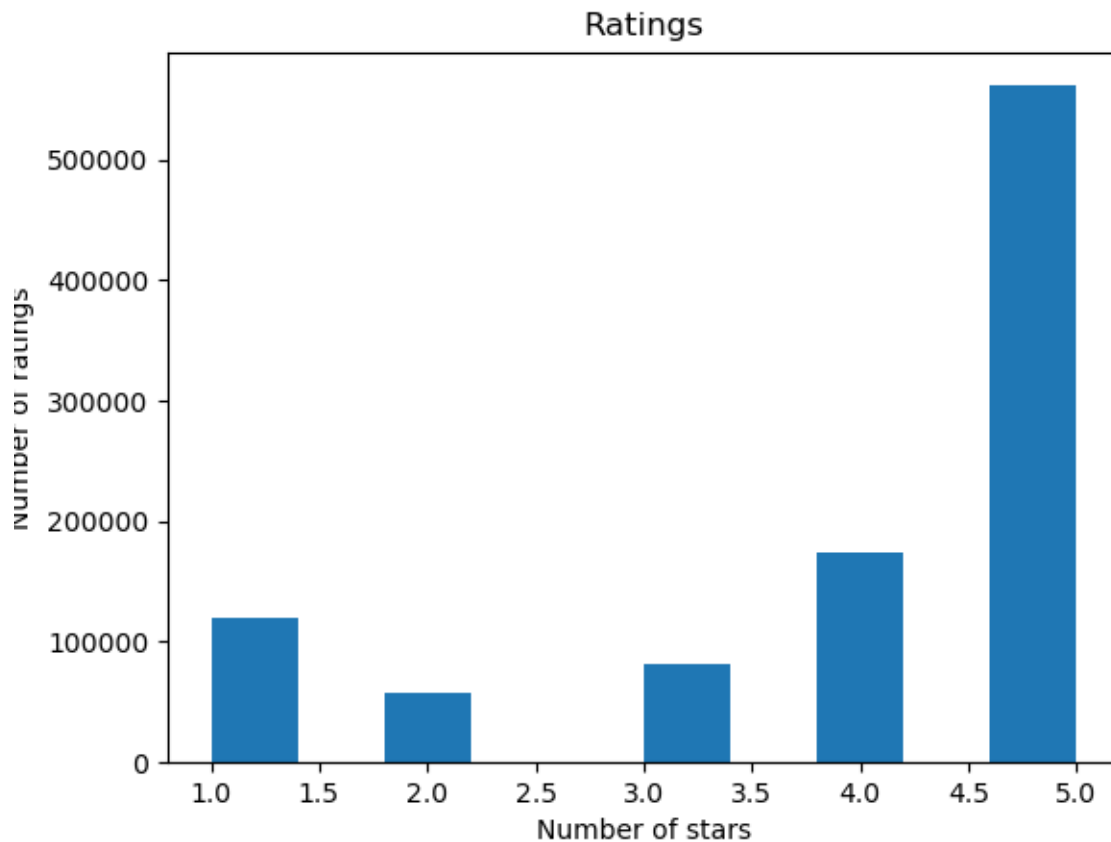
1. The number of ratings per user
2. The number of ratings per item
3. The ratings themselves

Note: for the first two, use `plt.yscale('log', nonposy='clip')` to put the histograms on a log-scale. Also, you can use `X.getnnz` to get the total number of nonzero elements along a specific axis.

Answer:







## 1.2 Finding similar items with nearest neighbours

Rubric: {code:6}

We'll use scikit-learn's `neighbors.NearestNeighbors` object to find the items most similar to the example item above, namely the Brass Grill Brush 18 Inch Heavy Duty and Extra Strong, Solid Oak Handle, at URL <https://www.amazon.com/dp/B00CFM0P7Y>.

Find the 5 most similar items to the Grill Brush using the following metrics (Make sure to include the code you have written in your pdf GradeScope submission.):

1. Euclidean distance (the `NearestNeighbors` default)
2. Normalized Euclidean distance (you'll need to do the normalization)
3. Cosine similarity (by setting `metric='cosine'`)

Some notes/hints...

- If you run `python main.py -q 1.2`, it will grab the row of `X` associated with the grill brush. The mappers take care of going back and forth between the IDs (like "B00CFM0P7Y") and the indices of the sparse array (0, 1, 2, ...).

- Keep in mind that scikit-learn's `NearestNeighbors` is for taking neighbors across rows, but here we're working across columns.
- Keep in mind that scikit-learn's `NearestNeighbors` will include the query item itself as one of the nearest neighbours if the query item is in the "training set".
- Normalizing the columns of a matrix would usually be reasonable to implement, but because  $X$  is stored as a sparse matrix it's a bit more of a mess. Therefore, use `sklearn.preprocessing.normalize` to help you with the normalization in part 2.

Answer:

```
# YOUR CODE HERE FOR Q1.2
#tuning matrix
d = len(set(ratings["item"]))
to_keep = list(set(range(d))-{grill_brush_ind})
X_wo_query = X[:,to_keep]
new_X = X_wo_query.transpose()
nbrs = NearestNeighbors(n_neighbors=5, algorithm='brute', metric='euclidean').fit(new_X)
distances, indices = nbrs.kneighbors(grill_brush_vec.transpose())
print("Closest 5 items using Euclidean distance:")
for i in indices[0]:
    print(item_inverse_mapper[i])

print('Euclidean number of ratings', X[:,indices[0]].getnnz(axis=0))

normalized_X = normalize(new_X)
nbrs = NearestNeighbors(n_neighbors=5, algorithm='brute', metric='euclidean').fit(normalized_X)
distances, indices = nbrs.kneighbors(grill_brush_vec.transpose())
print("Closest 5 items using normalized euclidean distance:")
for i in indices[0]:
    print(item_inverse_mapper[i])

nbrs = NearestNeighbors(n_neighbors=5, algorithm='brute', metric='cosine').fit(new_X)
distances, indices = nbrs.kneighbors(grill_brush_vec.transpose())
print("Closest 5 items using cosine distance:")
for i in indices[0]:
    print(item_inverse_mapper[i])

print('Cosine number of ratings:', X[:,indices[0]].getnnz(axis=0))
```

Closest 5 items using Euclidean distance:

B00IJB4MLA  
B00IJ8951K  
B00HVXLFZS  
B00EXDNA3O  
B00743MZCM

Closest 5 items using normalized euclidean distance:

B00IJB4MLA  
B00IJB5MCS  
B00IJ8951K  
B00EF4206E  
B00EF0I5Q2

Closest 5 items using cosine similarity:

B00IJB4MLA  
B00IJB5MCS  
B00IJ8951K  
B00EF4206E  
B00EF0I5Q2

Did normalized Euclidean distance and cosine similarity yields the same similar items, as expected?

Answer: Yes, the closest items for normalized euclidean distance is the same as cosine similarity.

### 1.3 Total popularity

Rubric: {reasoning:2}

For both Euclidean distance and cosine similarity, find the number of reviews for each of the 5 recommended items and report it. Do the results make sense given what we discussed in class about Euclidean distance vs. cosine similarity and popular items?

Note: in `main.py` you are welcome to combine this code with your code from the previous part, so that you don't have to copy/paste all that code in another section of `main.py`.

Answer:

Number of reviews for euclidean distance: [45 1 6 22 1]

Number of reviews for cosine similarity: [45 55 1 1 2]

Yes it makes sense as euclidean distance favours the unpopular items with less reviews.

## 2 Matrix Notation and Minimizing Quadratics

### 2.1 Converting to Matrix/Vector/Norm Notation

Rubric: {reasoning:3}

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.  $\max_{i \in \{1, 2, \dots, n\}} |w^T x_i - y_i|$ .
2.  $\sum_{i=1}^n v_i (w^T x_i - y_i)^2 + \frac{\lambda}{2} \sum_{j=1}^d w_j^2$ .
3.  $(\sum_{i=1}^n |w^T x_i - y_i|)^2 + \frac{1}{2} \sum_{j=1}^d \lambda_j |w_j|$ .

Note that in part 2 we give a *weight*  $v_i$  to each training example and the value  $\lambda$  is a non-negative scalar, whereas in part 3 we are regularizing the parameters with different weights  $\lambda_j$ . You can use  $V$  to denote a diagonal matrix that has the values  $v_i$  along the diagonal, and  $\Lambda$  as a diagonal matrix that has the  $\lambda_j$  values along the diagonal. You can assume that all the  $v_i$  and  $\lambda_i$  values are non-negative.

Answer:

1.  $\max_{i \in \{1, 2, \dots, n\}} |w^T x_i - y_i| = \|Xw - y\|_\infty$

2.

$$\begin{aligned} & \sum_{i=1}^n v_i (w^T x_i - y_i)^2 + \frac{\lambda}{2} \sum_{j=1}^d w_j^2 \\ &= \sum_{i=1}^n (w^T x_i v_i^{1/2} - y_i v_i^{1/2})^2 + \frac{\lambda}{2} \|w\|^2 \\ &= \|V^{1/2}(Xw - y)\|^2 + \frac{\lambda}{2} \|w\|^2 \end{aligned}$$

3.

$$\begin{aligned} & (\sum_{i=1}^n |w^T x_i - y_i|)^2 + \frac{1}{2} \sum_{j=1}^d \lambda_j |w_j| \\ &= \|Xw - y\|_1^2 + \frac{1}{2} \|\Lambda w\|_1 \end{aligned}$$

## 2.2 Minimizing Quadratic Functions as Linear Systems

Rubric: {reasoning:3}

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).

1.  $f(w) = \frac{1}{2}\|w - v\|^2$  (projection of  $v$  onto real space).
2.  $f(w) = \frac{1}{2}\|Xw - y\|^2 + \frac{1}{2}w^T \Lambda w$  (least squares with weighted regularization).
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  $d$  by 1 vectors (in part 3  $v$  is a vector of length  $n$  by 1), and  $\Lambda$  is a  $d$  by  $d$  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 sanity check for your derivation, make sure that your results have the right dimensions. **As a sanity check, make that the dimensions match for all quantities/operations: in order to make the dimensions match for some parts you may need to introduce an identity matrix.** For example,  $X^T X w + \lambda w$  can be re-written as  $(X^T X + \lambda I)w$ .

Answer:

1.

$$\begin{aligned} f(w) &= \frac{1}{2} \|w - v\|^2 \\ &= \frac{1}{2} (w - v)^T (w - v) \\ &= \frac{1}{2} (w^T w - 2w^T v + v^T v) \end{aligned}$$

$$\nabla f(w) = w - v$$

$$w - v = 0$$

$$w = v$$

2.

$$\begin{aligned} f(w) &= \frac{1}{2} \|Xw - y\|^2 + \frac{1}{2} w^T \Lambda w \\ &= \frac{1}{2} (Xw - y)^T (Xw - y) + \frac{1}{2} w^T \Lambda w \\ &= \frac{1}{2} (w^T X^T X w - w^T X^T y + y^T y + \frac{1}{2} w^T \Lambda w) \end{aligned}$$

$$\nabla f(w) = X^T X w - X^T y + \Lambda w$$

$$X^T X w - X^T y + \Lambda w = 0$$

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

3.

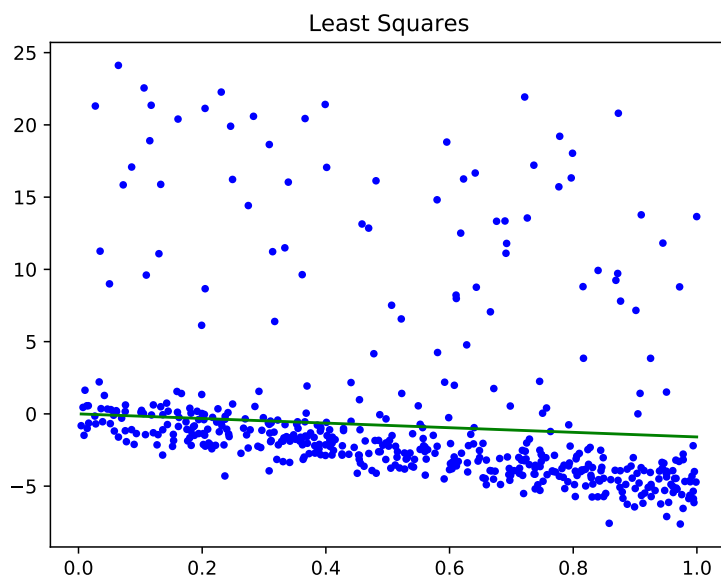
$$f(w) = \frac{1}{2} \sum_{i=1}^n v_i (w^T x_i - y_i)^2 + \frac{\lambda}{2} \|w - w^0\|^2$$



$$\begin{aligned}
&= \|V^{1/2}(Xw - y)\|^2 + \frac{\lambda}{2}(w - w^0)^T(w - w^0) \\
&= V^{1/2}((Xw - y)^T(Xw - y)) + \frac{\lambda}{2}(w - w^0)^T(w - w^0) \\
&= V^{1/2}(w^T X^T Xw - w^T X^T y + y^T y) + \frac{\lambda}{2}w^T w^0 \\
\nabla f(w) &= V^{1/2}(2X^T Xw - X^T y) + \lambda(w - w^0) \\
V^{1/2}(2X^T Xw - X^T y) + \lambda(w - w^0) &= 0 \\
2V^{1/2}X^T Xw + \lambda(w - w^0) &= V^{1/2}X^T y
\end{aligned}$$

### 3 Robust Regression and Gradient Descent

If you run `python main.py -q 3`, 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 Q4.1 we’ll address this explicitly.

#### 3.1 Weighted Least Squares in One Dimension

Rubric: {code:3}

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,

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

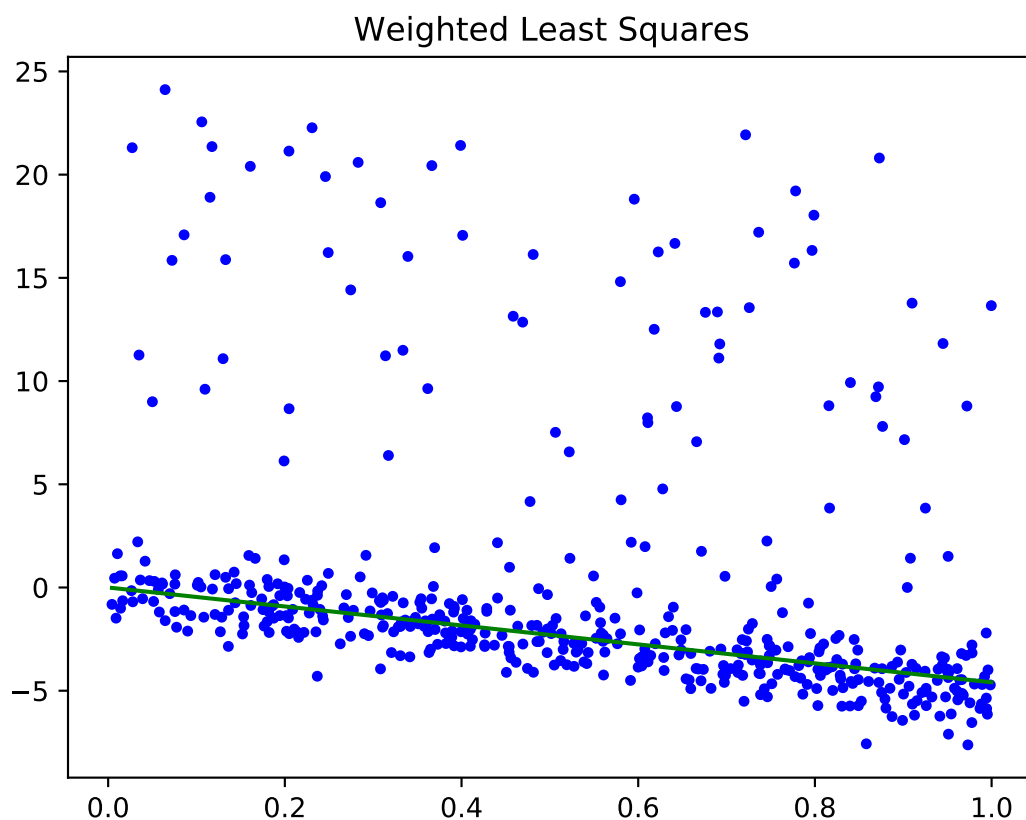
In this formulation, the model focuses on making the error small for examples  $i$  where  $v_i$  is high. Similarly, if  $v_i$  is low then the model allows a larger error. Note: these weights  $v_i$  (one per training example) are completely different from the model parameters  $w_j$  (one per feature), which, confusingly, we sometimes also call "weights".

Complete the model class, `WeightedLeastSquares`, that implements this model (note that Q2.2.3 asks you to show how a few 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). [Hand in your code and the updated plot.](#)

Answer:

```

15 # Least squares where each sample point x has a weight associated with it.
16 class WeightedLeastSquares(LeastSquares): # inherits the predict() function f
17     def fit(self, X, y, z):
18         self.w = solve(X.T@(X*z), X.T@(y*z))
19 
```



### 3.2 Smooth Approximation to the L1-Norm

Rubric: {reasoning:3}

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 function<sup>2</sup>:

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

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  $\nabla f$  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:

---

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

$$\begin{aligned}
g(\omega) &= f'(\omega) \\
&= \sum_{i=1}^n \frac{d}{d\omega} (\log(\exp(\omega^T x_i - y_i) + \exp(y_i - \omega^T x_i))) \\
&= \sum_{i=1}^n \frac{\frac{d}{d\omega} (\omega^T x_i - y_i) \exp(\omega^T x_i - y_i) + \frac{d}{d\omega} (y_i - \omega^T x_i) \exp(y_i - \omega^T x_i)}{\exp(\omega^T x_i - y_i) + \exp(y_i - \omega^T x_i)} \\
&= \sum_{i=1}^n \frac{\frac{d}{d\omega} (\omega^T x_i - y_i) \times (\exp(\omega^T x_i - y_i) - \exp(y_i - \omega^T x_i))}{\exp(\omega^T x_i - y_i) + \exp(y_i - \omega^T x_i)} \\
&= \begin{bmatrix} \sum_{i=1}^n \frac{x_{i1} (\exp(\omega^T x_i - y_i) - \exp(y_i - \omega^T x_i))}{\exp(\omega^T x_i - y_i) + \exp(y_i - \omega^T x_i)} \\ \sum_{i=1}^n \frac{x_{i2} (\exp(\omega^T x_i - y_i) - \exp(y_i - \omega^T x_i))}{\exp(\omega^T x_i - y_i) + \exp(y_i - \omega^T x_i)} \\ \vdots \\ \sum_{i=1}^n \frac{x_{id} (\exp(\omega^T x_i - y_i) - \exp(y_i - \omega^T x_i))}{\exp(\omega^T x_i - y_i) + \exp(y_i - \omega^T x_i)} \end{bmatrix}
\end{aligned}$$

### 3.3 Robust Regression

Rubric: {code:3}

The class `LinearModelGradient` 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 3.3` 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 `LinearModelGradient` for an example. Note that the `fit` function of `LinearModelGradient` also has a numerical check that the gradient code is approximately correct, since implementing gradients is often error-prone.<sup>3</sup>

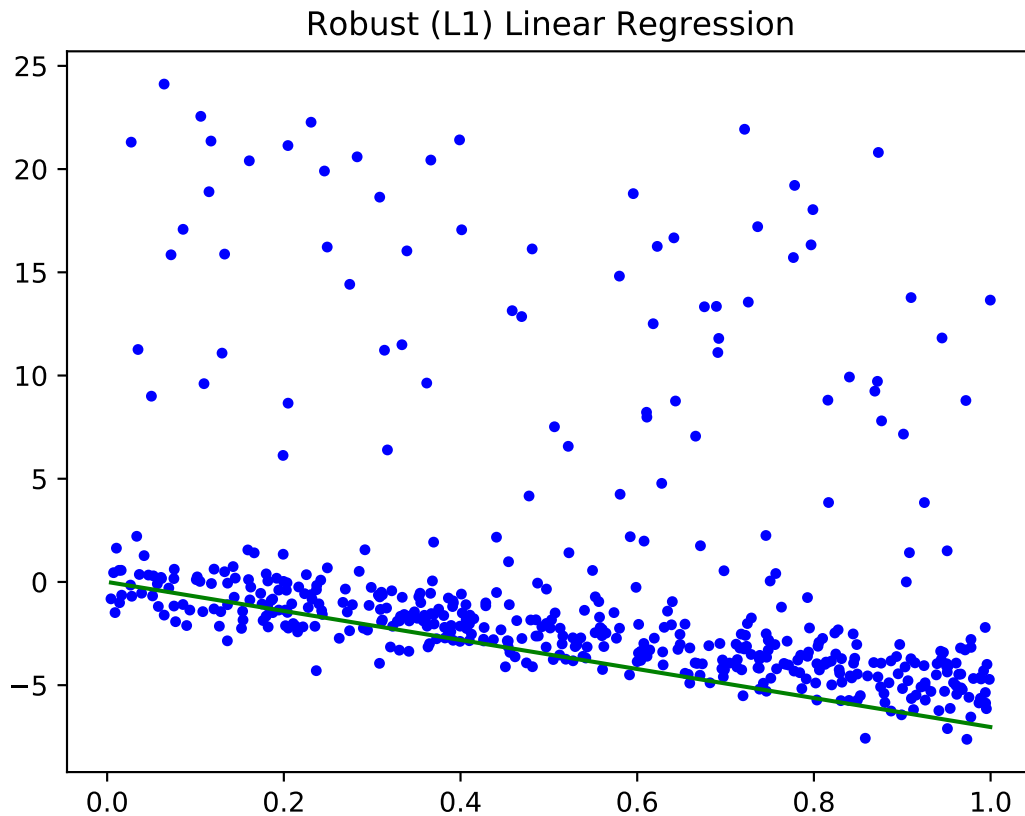
An advantage of gradient-based strategies is that they are able to solve problems that do not have closed-form solutions, such as the formulation from the previous section. The class `LinearModelGradient` has most of the implementation of a gradient-based strategy for fitting the robust regression model under the log-sum-exp approximation. The only part missing is the function and gradient calculation inside the `funObj` code. [Modify `funObj` 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

<sup>3</sup>Sometimes the numerical gradient checker itself can be wrong. See CPSC 303 for a lot more on numerical differentiation.

using this robust regression approach.

Answer:

```
def funObj(self,w,X,y):  
  
    # Calculate the function value  
    f = np.sum(np.log(np.exp(X@w - y) + np.exp(y - X@w)))  
  
    # Calculate the gradient value  
    g = np.zeros((1, X.shape[1]))  
    for j in range(X.shape[1]):  
        g_i = 0  
        for i in range(X.shape[0]):  
            e1 = np.exp(X[i]@w - y[i])  
            e2 = np.exp(X[i]@w + y[i])  
            g_i += X[i,j] * (e1-e2) / (e1+e2)  
        g[j] = g_i  
    return (f,g)
```



## 4 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 test error.

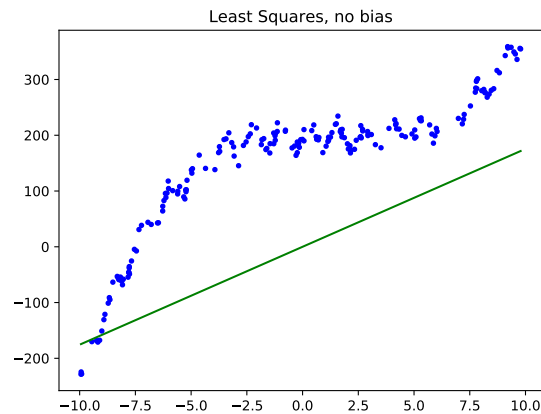
### 4.1 Adding a Bias Variable

Rubric: {code:3}

If you run `python main.py -q 4`, 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 test error (on a dataset not used for training).
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 test 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_model.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  $\beta$  in lecture). Hand in your new class, the updated plot, and the updated training/test 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:

Updated training error: 3551.3

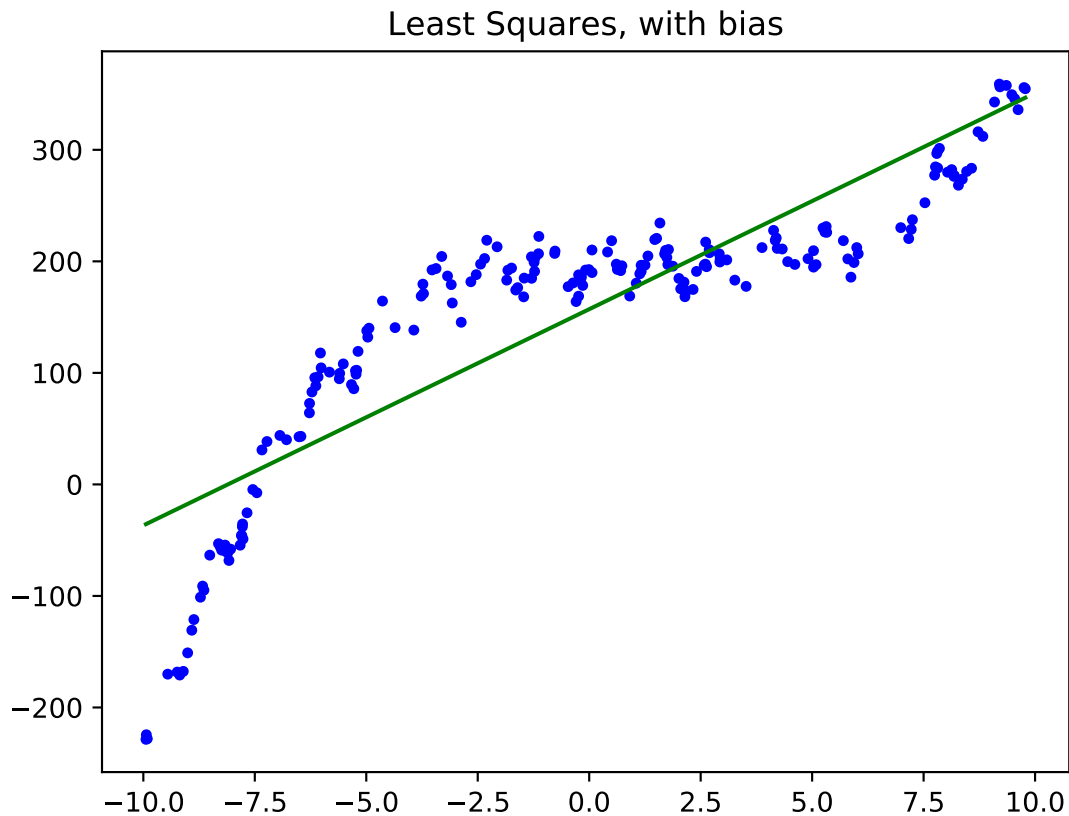
Updated test error: 3393.9

*# Least Squares with a bias added*

**class** `LeastSquaresBias`:

```
def fit(self, X, y):
    Z = np.concatenate((np.ones((X.shape[0], 1)), X), axis=1)
    self.w = solve(Z.T@Z, Z.T@y)

def predict(self, X):
    Z = np.concatenate((np.ones((X.shape[0], 1)), X), axis=1)
    return Z@self.w
```



## 4.2 Polynomial Basis

Rubric: {code:4}

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. [Hand in the new class](#), and report the training and test error for  $p = 0$  through  $p = 10$ . Explain the effect of  $p$  on the training error and on the test error.

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

Answer:

p=0



Training error = 15480.5  
Test error = 14390.8

$p=1$   
Training error = 3551.3  
Test error = 3393.9

$p=2$   
Training error = 2168.0  
Test error = 2480.7

$p=3$   
Training error = 252.0  
Test error = 242.8

$p=4$   
Training error = 251.5  
Test error = 242.1

$p=5$   
Training error = 251.1  
Test error = 239.5

$p=6$   
Training error = 248.6  
Test error = 246.0

$p=7$   
Training error = 247.0  
Test error = 242.9

$p=8$   
Training error = 241.3  
Test error = 246.0

$p=9$   
Training error = 235.8  
Test error = 259.3

$p=10$   
Training error = 235.1  
Test error = 256.3

At low values of  $p$  ( $p < 3$ ), there is high training and test error. As  $p$  increases, there is a decrease in both the training and test error. However, as  $p$  increases beyond a certain point ( $p = 6$ ), the test error starts increasing even though the training error continues to fall. This is due to the model overfitting to the training data.

```

# Least Squares with polynomial basis
class LeastSquaresPoly:
    def __init__(self, p):
        self.leastSquares = LeastSquares()
        self.p = p

    def fit(self, X, y):
        Z = self.__polyBasis(X)
        self.w = solve(Z.T@Z, Z.T@y)

    def predict(self, X):
        Z = self.__polyBasis(X)
        return Z@self.w

    # A private helper function to transform any matrix X into
    # the polynomial basis defined by this class at initialization
    # Returns the matrix Z that is the polynomial basis of X.
    def __polyBasis(self, X):
        Z = np.ones((X.shape[0], self.p+1))
        for i in range(1, self.p+1):
            Z[:, i] = X[:, 0]**i
        return Z

```

## 5 Very-Short Answer Questions

Rubric: {reasoning:7}

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: It is assigned to the cluster with the closest mean. For density based clustering they are not assigned.
2. Why do need random restarts for  $k$ -means but not for density-based clustering?  
Answer: kmeans is sensitive to initialization or the selection of the starting mean whereas density-based clustering is not as it does not matter where it starts from.
3. Can hierarchical clustering find non-convex clusters?  
Answer: Yes it can produce a tree of clusterings which can represent non-convex clusters
4. For model-based outlier detection, list an example method and problem with identifying outliers using this method.  
Answer: KNN outlier detection. It is not good for local outliers, for example an outlier far from the mean of a cluster but still relatively close to the data points of the data set.

5. For graphical-based outlier detection, list an example method and problem with identifying outliers using this method.  
 Answer: Scatterplot. It can only handle 2 variables at a time.
6. For supervised outlier detection, list an example method and problem with identifying outliers using this method.  
 Answer: Decision trees. It cannot detect new types of outliers and we need to know what the outliers look like.
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: No it would not make sense. The time complexity of linear regression is  $O(ndt)$  and the time complexity of normal equations is  $O(nd^2 + d^3)$  since  $d=1$  and  $t$  needs to be reasonably large in order for gradient descent to be accurate, it is better to use normal equations.
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: So that we can have a non-zero y-intercept(bias) by changing features and ignore the y-intercept in derivations, which is clearer.
9. If a function is convex, what does that say about stationary points of the function? Does convexity imply that a stationary points exists?  
 Answer: A stationary point is the global minima. Yes it must exist for the function to be convex.
10. 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: For a large number of features, gradient descent is faster because of its time complexity and if  $t$  is less than  $\text{minimum}(d, d^2/n)$  and it is good enough, then gradient descent should be used.
11. What is the problem with having too small of a learning rate in gradient descent?  
 Answer:  $t(\text{number of iterations})$  needs to be large to find a minimum, so it will take a long time.
12. What is the problem with having too large of a learning rate in gradient descent?  
 Answer: Going down the negative gradient may result in the function actually increasing as the minimum point is overshoot because the learning rate is too large.
13. What is the purpose of the log-sum-exp function and how is this related to gradient descent?  
 Answer:  $L_\infty - \text{norm}$  is convex but non-smooth. Log-sum-exp function is a smooth approximation to max function and hence can fit  $L_\infty - \text{norm}$  for gradient descent.
14. What type of non-linear transform might be suitable if we had a periodic function?  
 Answer: A sine transform as it is periodic.