

Linear Regression*

*Based on Chapter 5, “Machine Learning Refined: Foundations, Algorithms, and Applications”, Jeremy Watt, Reza Borhani, Aggelos K. Katsaggelos, 2nd Ed., Cambridge University Press, 2020

Least Squares Linear Regression

- In this lecture, we formally describe the problem of *linear regression*.
- That is, *fitting of a representative line or hyperplane in higher dimensions to a set of input/output data points*.
- Regression, in general, may be performed for a variety of reasons to:
 - Produce a so-called trend line or - more generally - a curve that can be used to help visually summarize the data.
 - Emphasize a particular point about the data under study.
 - Learn a model so that precise predictions can be made regarding output values in the future.

Notation and Modeling

- Data for regression problems comes in the form of a set of P input/output observation pairs:

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_P, y_P)$$

- More compactly, this is $\{(\mathbf{x}_p, y_p)\}_{p=1}^P$, where \mathbf{x}_p and y_p denote the p^{th} input and output, respectively.
- In simple instances, the input is scalar-valued.
- Then, the linear regression problem is geometrically equivalent of fitting a line to the associated scatter of data points in 2-dimensional space.
- However, each input \mathbf{x}_p generally may be a column vector of length N :

$$\mathbf{x}_p = [x_{1,p} \ x_{2,p} \ \cdots \ x_{N,p}]^T$$

- Now, the linear regression problem is analogous to fitting a hyperplane to a scatter of points in $(N + 1)$ -dimensional space.

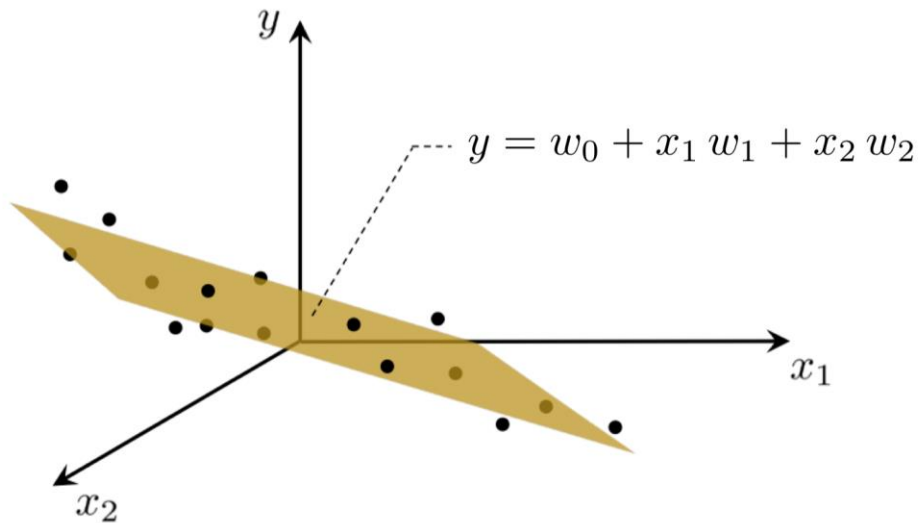
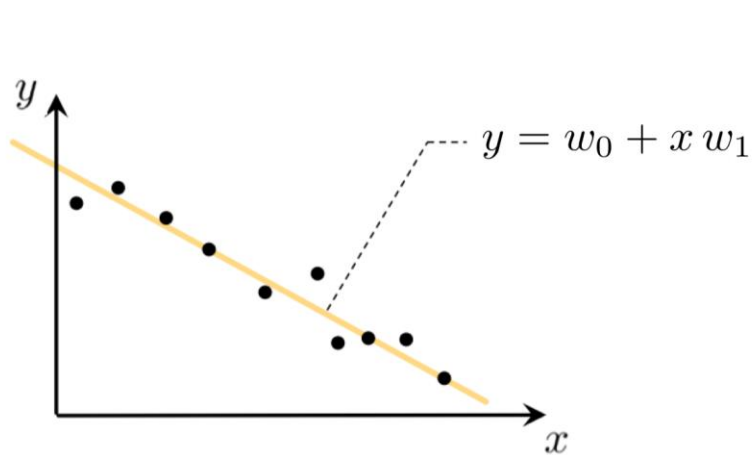
Notation and Modeling (cont'd)

- With scalar input, we must determine appropriate vertical (y-) intercept w_0 and slope w_1 so that the following holds:

$$\hat{y}_p = w_0 + x_p w_1 \approx y_p, \quad p = 1, \dots, P.$$

- The approximately equal sign (\approx) is used *because we cannot be sure that all data lies completely on a line*.
- When dealing with N dimensional input, we have a bias w_0 and N associated slope weights w_i to tune appropriately:

$$\hat{y}_p = w_0 + x_{1,p} w_1 + x_{2,p} w_2 + \dots + x_{N,p} w_N \approx y_p, \quad p = 1, \dots, P.$$



Notation and Modeling (cont'd)

- Note:
 - Each dimension of the *input* is referred to in the jargon of machine learning as a *feature* or *input feature*.
 - So, $w_{1,p}$, $w_{2,p}$, ..., $w_{N,p}$ are often referred to as the *feature-touching weights*.
 - The only weight *not* touching a feature is the bias w_0 .
- For any N , we can write the previous equation more compactly as:

$$\dot{\mathbf{x}}_p^T \mathbf{w} \approx y_p, \quad p = 1, \dots, P$$

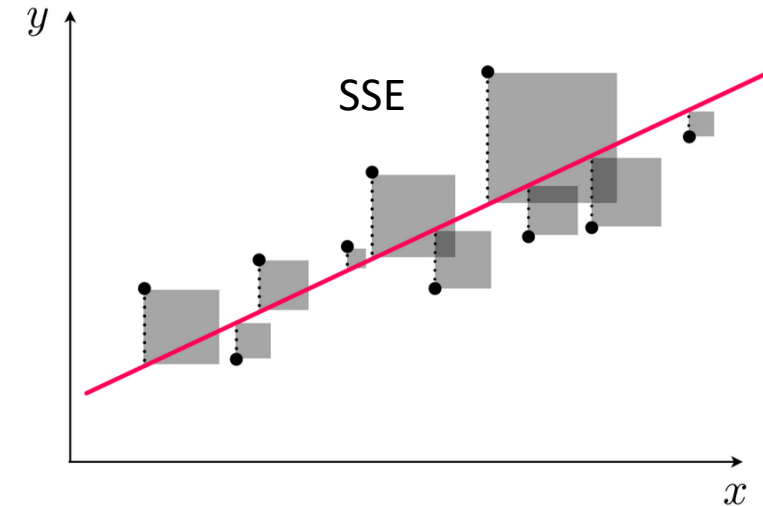
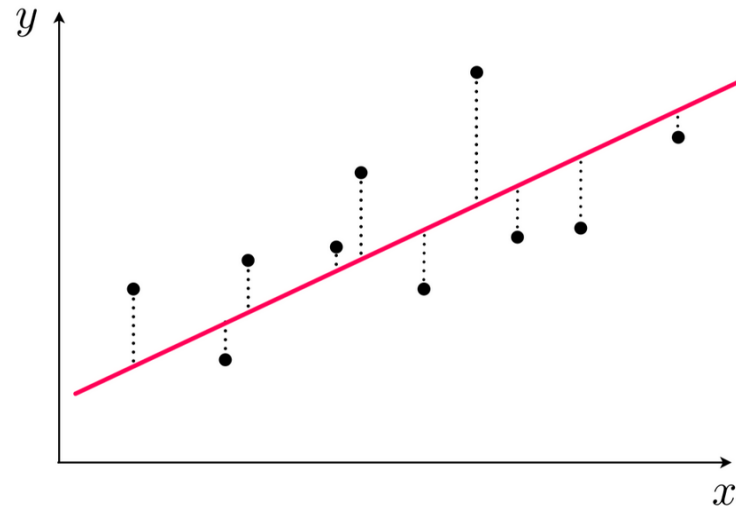
where:

$$\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_{N-1} \\ w_N \end{bmatrix} \quad \dot{\mathbf{x}}_p = \begin{bmatrix} 1 \\ x_{1,p} \\ x_{2,p} \\ \vdots \\ x_{N,p} \end{bmatrix}$$

and notation $\dot{\mathbf{x}}$ denotes an input \mathbf{x} with a 1 placed on top of it.

The Least Squares cost function

- To find the parameters of the hyperplane that best fits a regression dataset, we must use a cost function to measure how well a linear model fits it.
- A popular choice is to employ the *Least Squares cost function*.
- For a given set of parameters \mathbf{w} , this cost function computes the total squared error SSE between the associated hyperplane and the data, as illustrated in the figure below.
- Then, the best fitting hyperplane is the one whose parameters minimize this error.



The Least Squares cost function (cont'd)

- Where does this “Least Squares” cost come from?
- Remember: We want to find a weight vector \mathbf{w} so that each of the P approximate equalities below holds as tightly as possible:

$$\hat{y}_p = \dot{\mathbf{x}}_p^T \mathbf{w} \approx y_p, \quad p = 1, \dots, P.$$

- Equivalently, the *error* (i.e., *difference*):

$$\epsilon_p = \hat{y}_p - y_p$$

between $\dot{\mathbf{x}}_p^T \mathbf{w}$ (approximation) and y_p (true value) must be small.

- One natural way to measure the error between two quantities like this is to measure its *square* as:

$$g_p(\mathbf{w}) = (\dot{\mathbf{x}}_p^T \mathbf{w} - y_p)^2$$

- This means that *both negative and positive errors are treating equally*.

The Least Squares cost function (cont'd)

- This squared error $g_p(\cdot)$ is one example of a *point-wise cost* or *loss* that measures the error of a model on the point $\{\mathbf{x}_p, y_p\}$.
- Since we want all P such values to be small, we can take their average*, forming the *Least Squares* cost function for linear regression:

$$g(\mathbf{w}) = \frac{1}{P} \sum_{p=1}^P g_p(\mathbf{w}) = \frac{1}{P} \underbrace{\sum_{p=1}^P \overbrace{(\dot{\mathbf{x}}_p^T \mathbf{w} - y_p)^2}^{\text{loss}}}_{\text{SSE}}$$

- Observation: The Least Squares cost function is a function of both the weights \mathbf{w} and the data.
- However, when we express the function in mathematical shorthand as $g(\mathbf{w})$ (as we do on the left-hand side above), we only show dependency on the weights \mathbf{w} .

*Another option followed: apply minimization to SSE, which is equivalent.

The Least Squares cost function (cont'd)

- Why do we do this?
 1. Foremost, for notational simplicity and convenience.
 2. For a given dataset, the weights \mathbf{w} are the important input, *since this is what we need to tune to produce a good fit.*
 3. From an optimization perspective, *the dataset itself is considered fixed.*
- In overall:

We want to tune our parameters \mathbf{w} to *minimize* the Least Squares cost,...

...since the larger this value becomes ...

...the larger the squared error between the corresponding linear model and the data,...

...and so, the poorer we represent the given dataset using a linear model.
- In other words:

We want to determine a value for weights \mathbf{w} that *minimizes* $g(\mathbf{w})$, written formally as:

$$\text{minimize}_{\mathbf{w}} \frac{1}{P} \sum_{p=1}^P (\mathbf{x}_p^T \mathbf{w} - y_p)^2$$

Least Squares method: Simple or bivariate regression

- Taking the partial derivatives:

$$\frac{\partial \sum (y_i - (w_0 + w_1 x_i))^2}{\partial w_0} = -2 \sum (y_i - (w_0 + w_1 x_i)) = 0$$

$$\frac{\partial \sum (y_i - (w_0 + w_1 x_i))^2}{\partial w_1} = -2 \sum x_i (y_i - (w_0 + w_1 x_i)) = 0$$

we find

(a) the slope w_1 :

$$w_1 = \frac{\sum (x_i - \mu_X)(y_i - \mu_Y)}{\sum (x_i - \mu_X)^2} = \frac{\sum x_i y_i - P \mu_X \mu_Y}{\sum x_i^2 - P \mu_X^2} = \frac{S_{XY}}{S_{XX}} = \frac{Cov(X, Y)}{Var(X)}$$

where:

$$\mu_X = \frac{\sum x_i}{n} \quad \text{and} \quad \mu_Y = \frac{\sum y_i}{n}$$

Least Squares method: Simple or bivariate regression (cont'd)

and

(b) the y -intercept w_0 :

$$w_0 = \mu_Y - w_1\mu_X$$

- So:

$$\hat{y} = \mu_Y + w_1(x - \mu_X)$$

- Observation:

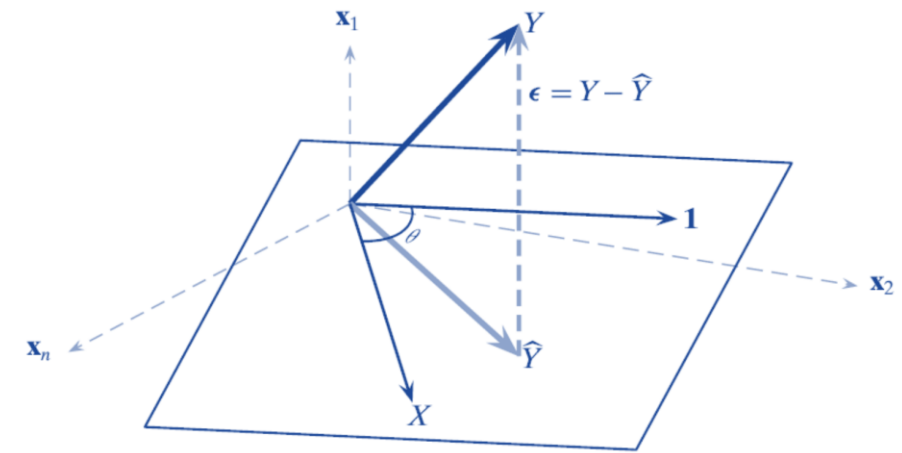
The point (μ_X, μ_Y) *always belongs to the regression line.*

Least Squares method: Geometry of simple or bivariate regression

- Note that the predictions can be written as:

$$\hat{Y} = w_0 \mathbf{1} + w_1 X$$

- This equation indicates that the predicted vector \hat{Y} is a *linear combination* of $\mathbf{1}$ and X .
- That is, it must lie in the column space spanned by $\mathbf{1}$ and X , given as $\text{span}(\{\mathbf{1}, X\})$.
- The figure on the right depicts that the optimal \hat{Y} minimizing the error ϵ is the orthogonal projection of Y onto the subspace $\text{span}(\{\mathbf{1}, X\})$
- The residual error vector ϵ is *thus orthogonal* to the subspace $\text{span}(\{\mathbf{1}, X\})$, and its squared length (or magnitude) equals the SSE value.



Least Squares method: Geometry of simple or bivariate regression (cont'd)

- We can create an orthogonal basis by decomposing \mathbf{X} into a component along $\mathbf{1}$, and a component orthogonal to $\mathbf{1}$, as shown in the figure below, where:

- The projection of \mathbf{X} onto $\mathbf{1}$ equals:

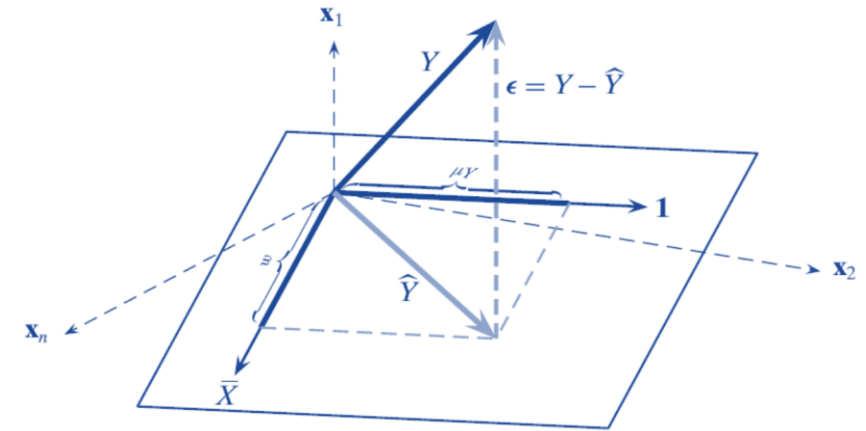
$$\left(\frac{\mathbf{X}^T \mathbf{1}}{\mathbf{1}^T \mathbf{1}} \right) \mathbf{1} = \left(\frac{\sum_{i=1}^n x_i}{n} \right) \mathbf{1} = \mu_X \mathbf{1}$$

- The orthogonal (demeaned) component $\bar{\mathbf{X}}$ equals:
 $\mathbf{X} - \mu_X \mathbf{1}$

- Since $\mathbf{1}$ and $\bar{\mathbf{X}}$ form an orthogonal basis, we can obtain the predicted vector $\hat{\mathbf{Y}}$ by projecting \mathbf{Y} onto them.
- Summing up these two components, we have:

$$\hat{\mathbf{Y}} = \mu_Y \mathbf{1} + \left(\frac{\mathbf{Y}^T \bar{\mathbf{X}}}{\bar{\mathbf{X}}^T \bar{\mathbf{X}}} \right) \bar{\mathbf{X}} = \left(\mu_Y - \left(\frac{\mathbf{Y}^T \bar{\mathbf{X}}}{\bar{\mathbf{X}}^T \bar{\mathbf{X}}} \right) \mu_X \right) \mathbf{1} + \left(\frac{\mathbf{Y}^T \bar{\mathbf{X}}}{\bar{\mathbf{X}}^T \bar{\mathbf{X}}} \right) \mathbf{X}$$

- Since $\hat{\mathbf{Y}} = w_0 \mathbf{1} + w_1 \mathbf{X}$, equating the two expressions, we arrive at the same w_0, w_1 values.



Least Squares method: Multiple regression

- Let $\dot{\mathbf{X}}$ be the matrix, where:
 1. The first column $\dot{\mathbf{X}}_0$ is equal to $[1 \ 1 \cdots 1]^T$.
 2. The other columns $\dot{\mathbf{X}}_j$ contain the data of the j th variable (factor).
- $\dot{\mathbf{X}}$ is also known as the “*model matrix*”.
- Setting:

$$\nabla_{\mathbf{w}} \left(\|\mathbf{Y} - \dot{\mathbf{X}}\mathbf{w}\|^2 \right) = -2\dot{\mathbf{X}}^T (\mathbf{Y} - \dot{\mathbf{X}}\mathbf{w}) = \mathbf{0}$$

we have:

$$\dot{\mathbf{X}}^T \mathbf{Y} = \dot{\mathbf{X}}^T \dot{\mathbf{X}} \mathbf{w}$$

- This expression defines a set of equations known as the *normal equations*.
- Given that $\dot{\mathbf{X}}$ has linear independent columns (and so does the Gram matrix $\dot{\mathbf{X}}^T \dot{\mathbf{X}}$), we have:

$$\mathbf{w} = \underbrace{(\dot{\mathbf{X}}^T \dot{\mathbf{X}})^{-1} \dot{\mathbf{X}}^T}_{\text{Pseudoinverse of } \dot{\mathbf{X}}} \mathbf{Y} \quad \text{and} \quad \hat{\mathbf{Y}} = \dot{\mathbf{X}} \mathbf{w} = \underbrace{\dot{\mathbf{X}} (\dot{\mathbf{X}}^T \dot{\mathbf{X}})^{-1} \dot{\mathbf{X}}^T}_{\mathbf{H}} \mathbf{Y}$$

- \mathbf{H} is called the “*hat matrix*”, since it ‘puts’ a hat ($\hat{}$) on \mathbf{Y} .

Least Squares method: Geometry of multiple regression

- As in the simple case:
 - \hat{Y} is the orthogonal projection of Y onto the subspace spanned by the column vectors \dot{X}_j .
 - The error $\epsilon = Y - \dot{X}w$ is orthogonal to this subspace and so to each \dot{X}_j .
- Thus:

$$\dot{X}_j^T \epsilon = 0, \quad \forall j$$

- From the above equations, we can derive again the *normal equations*.

Least Squares method: Geometry of multiple regression (cont'd)

- Since \dot{X}_j are linear independent, we can use the *Gram–Schmidt* orthogonalization to construct a set of orthogonal basis vectors U_0, U_1, \dots, U_d :

$$U_0 = \dot{X}_0$$

$$U_1 = \dot{X}_1 - p_{01}U_0$$

$$U_2 = \dot{X}_2 - p_{02}U_0 - p_{12}U_1$$

$$\vdots$$

$$U_d = \dot{X}_d - p_{0d}U_0 - p_{1d}U_1 - \dots - p_{d-1d}U_{d-1}$$

where:

$$p_{ij} = \frac{\dot{X}_j^T U_i}{U_i^T U_i} \text{ is the scalar projection of attribute } \dot{X}_j \text{ onto the basis vector } U_i.$$

- From the above expression, with rearrangement, we can write \dot{X} , also known as *QR-factorization*, as:

$$\dot{X} = QR,$$

where Q is the orthogonal matrix of orthogonal basis vectors U_0, U_1, \dots, U_d (note: $Q^T Q = \text{diag}(\|U_i\|^2) = \Delta$) and R the *upper triangular* matrix, with 1 in the main diagonal and the scalar projections p_{ij} in the rest cells.

- Then:

$$\dot{X}^T Y = \dot{X}^T \dot{X} w \Rightarrow R^T Q^T Y = R^T Q^T Q R w \text{ or } R^T Q^T Y = R^T \Delta R w \text{ or } w = R^{-1} \Delta^{-1} Q^T Y$$

Least Squares method: Geometry of multiple regression (cont'd)

- The above discussion can be summarized as follows:

MULTIPLE-REGRESSION (\mathbf{X}, \mathbf{Y})

- $\dot{\mathbf{X}} = [\mathbf{1} \ \mathbf{X}]$ // augmented data with $\mathbf{X}_0 = \mathbf{1} \in \mathbb{R}^n$
- $\{\mathbf{Q}, \mathbf{R}\} = \text{QR-factorization}(\dot{\mathbf{X}})$ // $\mathbf{Q} = [\mathbf{U}_0 \ \mathbf{U}_1 \cdots \mathbf{U}_d]$
- $\Delta^{-1} = \begin{bmatrix} \frac{1}{\|\mathbf{U}_0\|^2} & 0 & \cdots & 0 \\ 0 & \frac{1}{\|\mathbf{U}_1\|^2} & \cdots & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & \cdots & \frac{1}{\|\mathbf{U}_d\|^2} \end{bmatrix}$ // reciprocal squared norm
- Calculate $\Delta^{-1} \mathbf{Q}^T \mathbf{Y}$
- Solve for \mathbf{w} by back-substitution, using the fact that $\mathbf{R}\mathbf{w} = \Delta^{-1} \mathbf{Q}^T \mathbf{Y}$
- $\hat{\mathbf{Y}} = \mathbf{Q} \Delta^{-1} \mathbf{Q}^T \mathbf{Y}$ // $\hat{\mathbf{Y}} = \mathbf{Q} \mathbf{R} \mathbf{w}$

Least Squares method: Gradient descent

- Instead of using the QR-factorization approach to exactly solve the multiple regression problem, we can also employ the *simpler* gradient descent algorithm.
- Since:

$$\nabla_{\mathbf{w}} \left(\frac{1}{p} \|\mathbf{Y} - \mathbf{X}\mathbf{w}\|^2 \right) = -\frac{2}{p} \mathbf{X}^T (\mathbf{Y} - \mathbf{X}\mathbf{w})$$

we have the following update rule:

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \alpha \nabla_{\mathbf{w}} = \mathbf{w}^k + \alpha \underbrace{\frac{2}{p} \mathbf{X}^T (\mathbf{Y} - \mathbf{X}\mathbf{w}^k)}_{\tilde{\nabla}_{\mathbf{w}}} = \mathbf{w}^k + \alpha' \tilde{\nabla}_{\mathbf{w}}$$

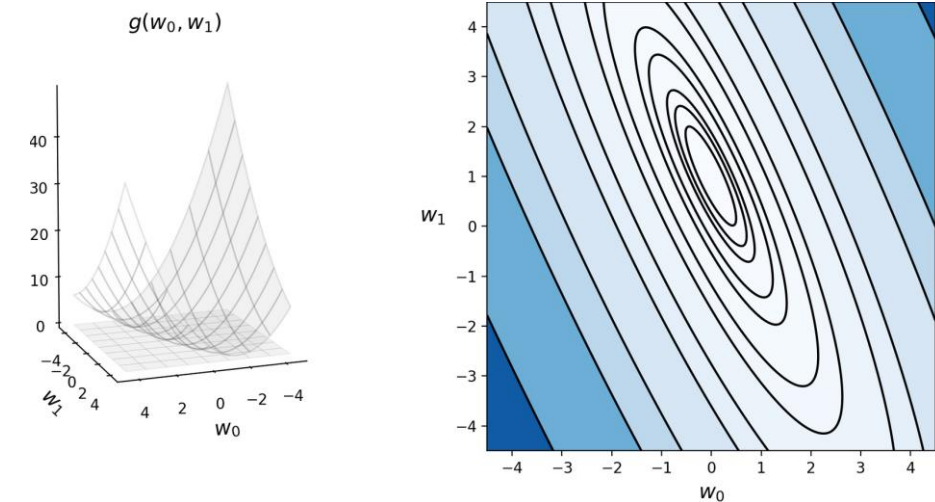
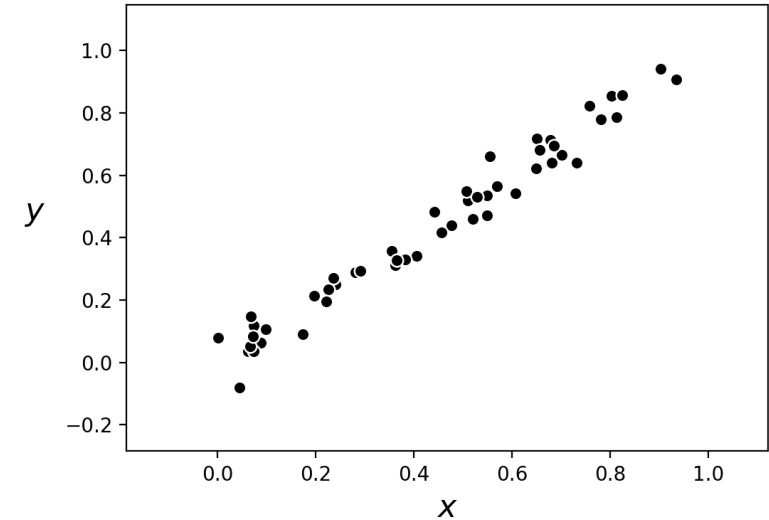
Minimization of the Least Squares cost function

- Remember:
Determining the overall shape of a function (i.e., the function is convex)...
...helps determine the appropriate optimization method(s) we can apply to efficiently determine the ideal parameters.
- In the case of the Least Squares cost function for linear regression, it is easy to check that *it is always a convex quadratic function regardless of the dataset**.
- For small input dimensions (e.g., $N = 1$), we can empirically verify this claim for any given dataset by simply plotting the function g - as a surface and/or contour plot - as the following example demonstrates.

* See Section 5.9 of the textbook.

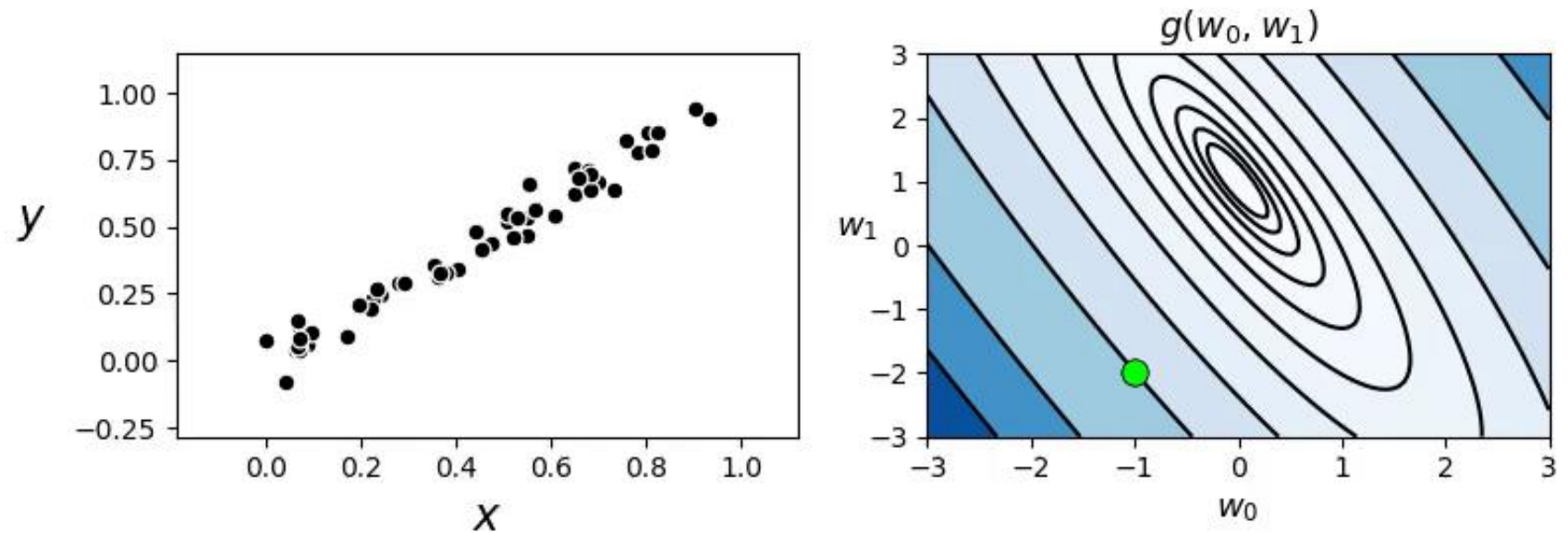
Example. Visually verifying the convexity of the cost function for a toy dataset

- In this example, we plot the contour and surface plot for the Least Squares cost function for linear regression for a toy dataset.
- This toy dataset consists of 50 points randomly selected from the line $y = x$, adding a small amount of Gaussian noise to each.
- The contour plot and corresponding surface generated by the Least Squares cost function using this data are shown in the bottom panel of the figure.



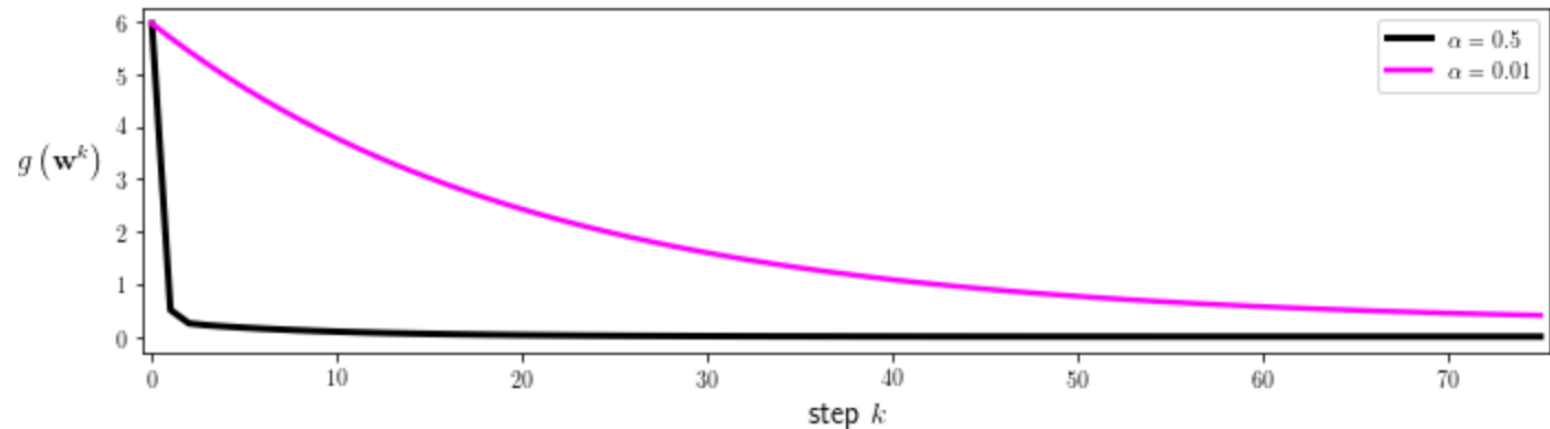
Example. Using gradient descent to minimize the Least Squares cost on a toy dataset

- Below, the minimization of the Least Squares cost is animated with the toy dataset presented in the previous example.
- More specifically, gradient descent is used, employing a fixed steplength value $\alpha = 0.5$ for all 75 steps, until approximately reaching the minimum of the function.



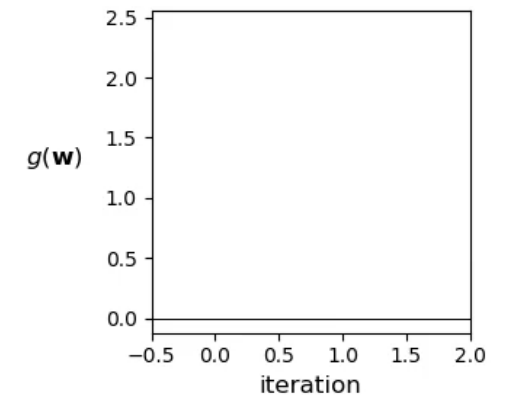
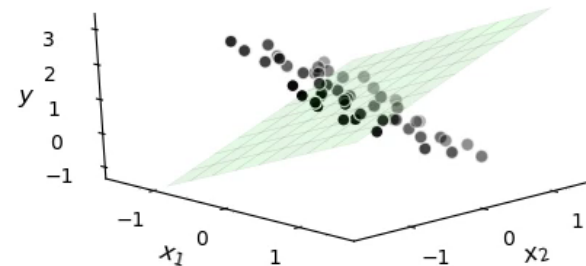
Example. Using gradient descent to minimize the Least Squares cost on a toy dataset (cont'd)

- Remember: Whenever we use a local optimization method like gradient descent, we must properly tune the steplength parameter α .
- This was done for the previous run, with several fixed steplength values.
- Below, those runs are re-created, using $\alpha = 0.5$, $\alpha = 0.01$ and showing the cost function history plot for each steplength value choice.
- We can see from the plot that indeed the first steplength value works considerably better.



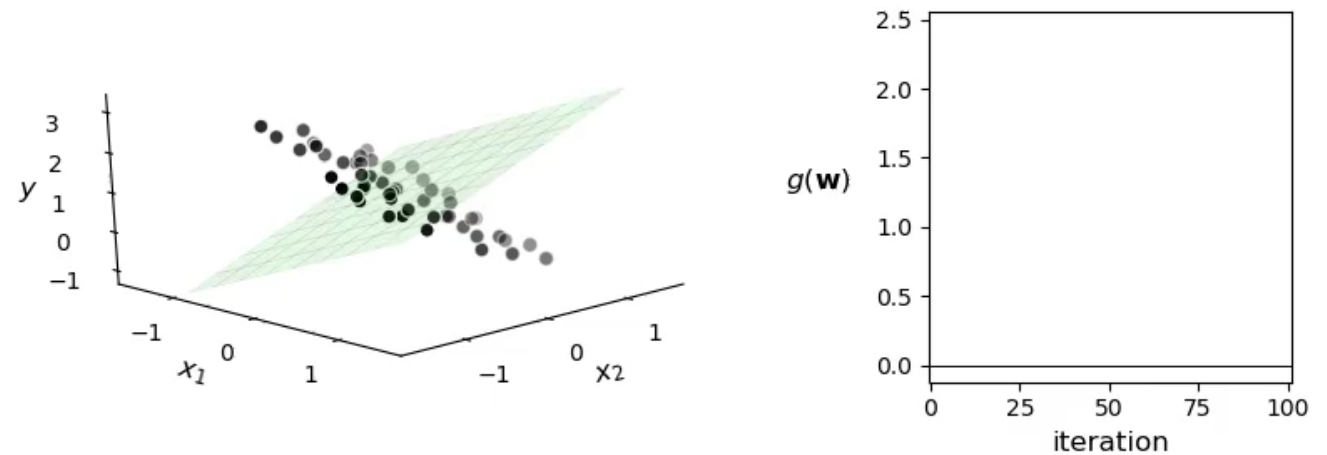
Example. Completely minimizing the Least Squares cost function using a single Newton step

- In the optional reading of the previous topic, we described how Newton's method perfectly minimizes any convex quadratic function in a single step.
- As mentioned, one can easily show that the Least Squares cost function is *a convex quadratic function for any dataset*.
- Thus, a *single Newton step* will completely minimize it, *regardless of the dataset used*.
- We animate this for a toy dataset with $N = 2$ inputs.



Example. Completely minimizing the Least Squares cost function using a single Newton step

- The same minimization using gradient descent



Implementing the Least Squares cost in Python

- When implementing a cost function like Least squares, it is helpful to think modularly, with the aim lightening the amount of mental 'bookkeeping' required, breaking down the cost into a few distinct pieces.
- Here, we can really break things down into two chunks: we have our model - a linear combination of input - and the cost - squared error - itself.
- Our linear model - a function of input and weights - is a function worthy enough of its own notation.
- We will write it as:

$$\text{model}(\mathbf{x}_p, \mathbf{w}) = \mathbf{x}_p^T \mathbf{w}.$$

- If we were to go back then, and use this modeling notation, we could equally well express our ideal settings of the weights our original approximation as

$$\text{model}(\mathbf{x}_p, \mathbf{w}) \approx y_p$$

and, likewise, our Least Squares cost function itself as

$$g(\mathbf{w}) = \frac{1}{P} \sum_{p=1}^P (\text{model}(\mathbf{x}_p, \mathbf{w}) - y_p)^2.$$

Implementing the Least Squares cost in Python (cont'd)

- To write this in Python we un-ravel the bias from our weights, using slicing and express the model computation as `a = w[0] + np.dot(x_p.T, w[1:])`
- The linear combination $\mathbf{x}_p^T \mathbf{w}_{[1:]}$ - implemented using `np.dot` - is very efficient (far more than an explicit for loop in Python).

```
# compute linear combination of input point
def model(x_p, w):
    # compute linear combination and return
    a = w[0] + np.dot(x_p.T, w[1:])
    return a.T
```

- Next, in writing out the Least Squares function, we could explicitly loop over our P points, one at-a-time.

```
# a least squares function for linear regression
def least_squares(w, x, y):
    # loop over points and compute cost contribution from each input/output pair
    cost = 0
    for p in range(y.size):
        # get pth input/output pair
        x_p = x[:, p][:, np.newaxis]
        y_p = y[p]

        ## add to current cost
        cost += (model(x_p, w) - y_p)**2

    # return average least squares error
    return cost/float(y.size)
```

Implementing the Least Squares cost in Python (cont'd)

- However, when such `for` loops can be equivalently written using `numpy` operations, this is typically far more efficient.
- *Explicit* `for` loops (including list comprehensions) written in Python are rather slow due to the very nature of the language (e.g., it being a dynamically typed interpreted language).
- It is easy to get around most of this inefficiency by replacing explicit `for` loops with numerically equivalent operations performed using operations from the numpy library.
- `numpy` is an API for some very efficient vector/matrix manipulation libraries written in C.
- In fact, Python code, employing heavy use of `numpy` functions, can often execute almost as fast a raw C implementation itself.

Implementing the Least Squares cost in Python (cont'd)

- Below, we show compact implementations of both our model and Least Squares cost in Python leveraging `numpy`.

```
# compute linear combination of input points
def model(x,w):
    a = w[0] + np.dot(x.T,w[1:])
    return a.T

# an implementation of the least squares cost function for linear regression
def least_squares(w):
    # compute the least squares cost
    cost = np.sum((model(x,w) - y)**2)
    return cost/float(y.size)
```

- Notice that, for simplicity, we write the Pythonic Least Squares cost function `least_squares(w)` instead of `least_squares(w, x, y)`, where in the latter case we explicitly list its other two arguments: the input `x` and output `y` data.
- This is done for notational simplicity, as `autograd` will correctly differentiate both forms (since by default it computes the gradient of a Python function with respect to its first input only).
- We will use this kind of simplified Pythonic notation when introducing future machine learning cost functions as well.

Least Absolute Deviations

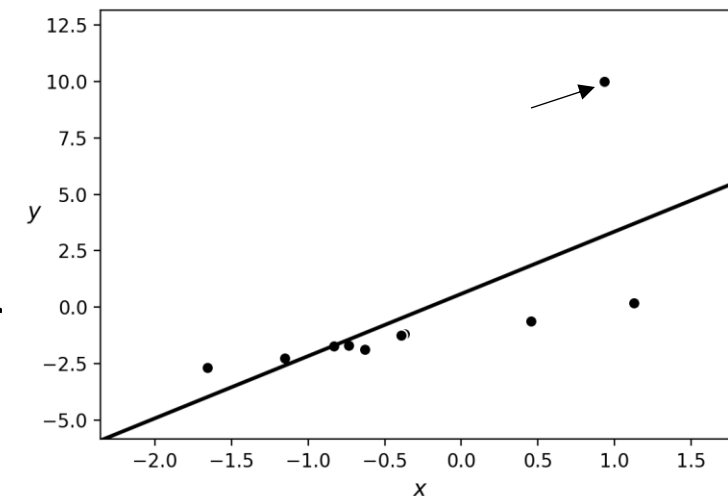
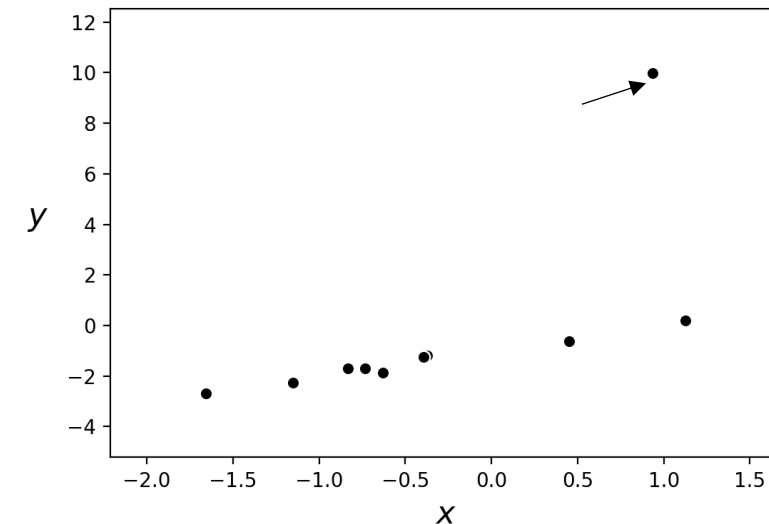
- In the sequel, we consider a slight twist on the derivation of the Least Squares cost function that leads to an alternative cost for linear regression.
- This cost is called *Least Absolute Deviations*.
- This alternative cost function is much more *robust* to outliers in a dataset than the original Least Squares.

The susceptibility of the Least Squares cost to outliers

- One downside of using the squared error in the Least Squares cost is that *squaring the error increases the importance of larger errors.*
- Specifically, *squaring errors of length greater than 1 makes these values considerably larger.*
- This forces weights learned via the Least Squares cost to produce a linear fit that is especially focused on trying to minimize these large errors.
- This is often due to *outliers* in a dataset.
- In other words: *the Least Squares cost produces linear models that tend to overfit to outliers in a dataset.*
- We illustrate this fact via a simple dataset in the following example.

Example. Least Squares overfits to outliers

- In this example, we use the dataset plotted on the right.
- This dataset can largely be represented by a proper linear model, except for a single large outlier.
- This example shows how the Least Squares cost function for linear regression tends to create linear models that *overfit* to outliers.
- First, the parameters of a linear regressor to this dataset are tuned by minimizing the Least squares cost via gradient descent.
- Next, the linear model associated with those weights achieving the smallest cost function values is plotted.
- This fit, shown in black, does not fit most of the data points well, bending upward.
- This occurred with the aim of minimizing the large squared error on the singleton outlier point.



Replacing Squared Error with Absolute Error

- How can we make our linear regression framework more robust to outliers?
- An alternative to *squared* error for our point-wise cost is the *absolute error*, defined as:

$$g_p(\mathbf{w}) = |\dot{\mathbf{x}}_p^T \mathbf{w} - y_p|, \quad p = 1, \dots, P.$$

- By using absolute error, *we still treat negative and positive errors equally.*
- *However, we do not exaggerate the importance of large errors greater than 1.*
- Why?
- *Because we do not square them.*

Replacing Squared Error with Absolute Error (cont'd)

- If we form the average of these absolute error point-wise costs (losses), we have the so-called *Least Absolute Deviations* cost function:

$$g(\mathbf{w}) = \frac{1}{P} \sum_{p=1}^P g_p(\mathbf{w}) = \frac{1}{P} \sum_{p=1}^P |\dot{\mathbf{x}}_p^T \mathbf{w} - y_p|$$

- This cost function is also *always convex, regardless of the input dataset* (like Least Squares).
- But we cannot use Newton's method to minimize this in a single step, since it doesn't have a useful second derivative as it is zero almost everywhere.
- Even though it is not differentiable everywhere, we can use gradient descent to minimize it.
- Technically speaking when doing this, we are using 'sub-gradient* descent', but it works all the same.

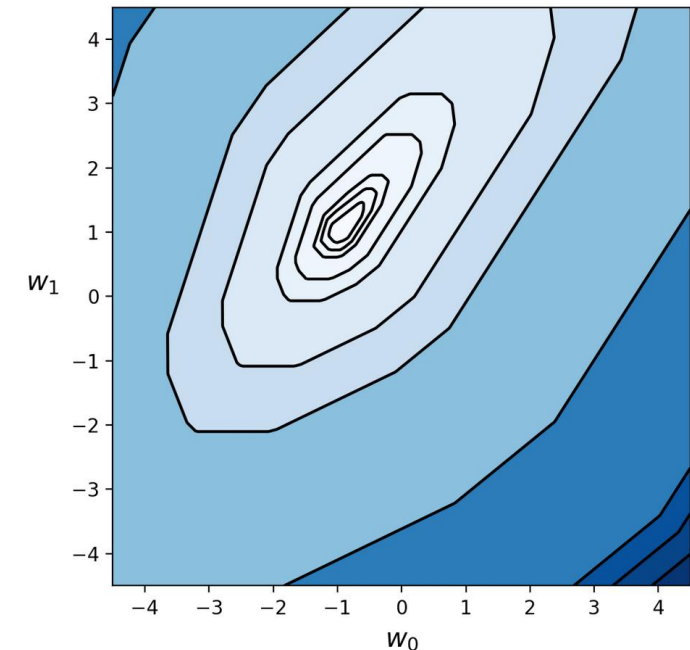
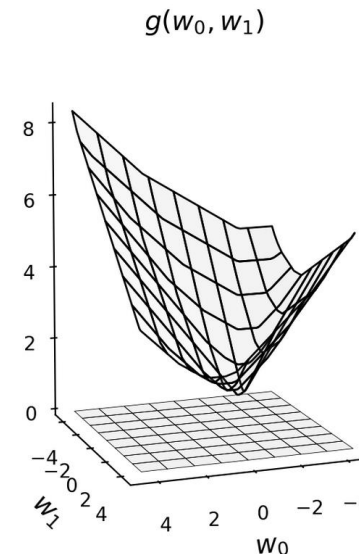
* v : $f(x) - f(x_0) \geq v \cdot (x - x_0)$

Least absolute deviations versus Least Squares

- In implementing the cost in Python, we can employ the model function we used with our Least Squares implementation shown previously.
- All we then need to do is slightly alter the cost function itself to get our desired implementation, as shown below.

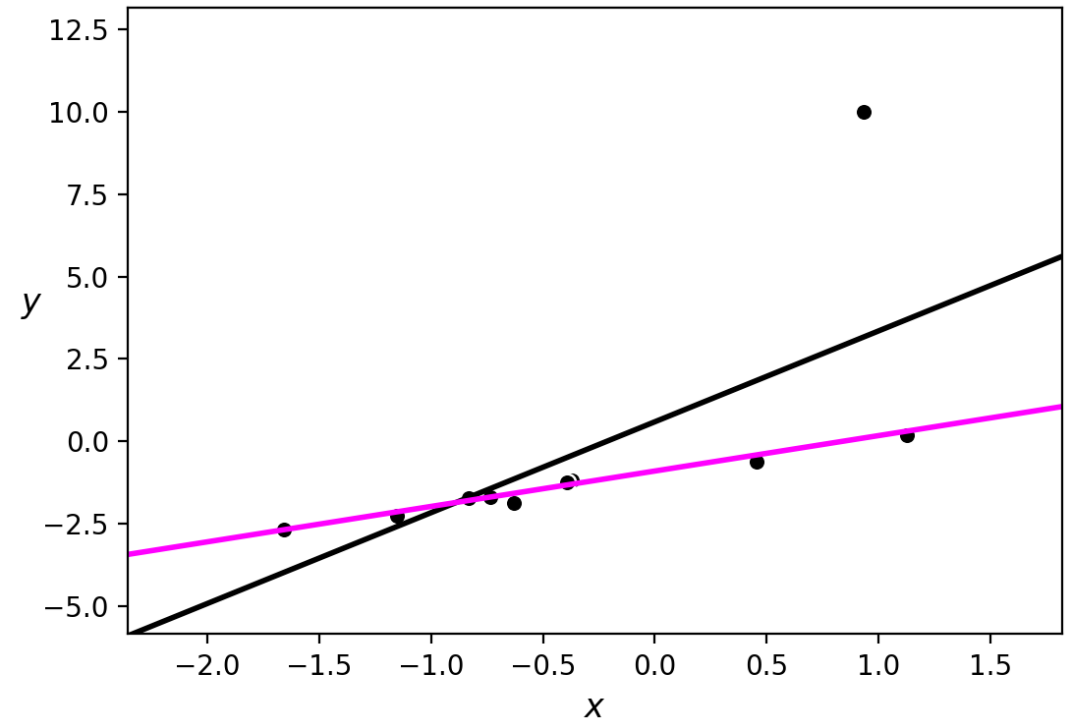
```
# a compact least absolute deviations cost function
def least_absolute_deviations(w):
    cost = np.sum(np.abs(model(x,w) - y))
    return cost/float(np.size(y))
```

- On the right, we plot the surface/ contour plot of this cost function using the previously shown dataset - indeed it is convex.



Example. Least absolute deviations versus Least Squares

- Let us plot and compare the best fit found via gradient descent for both Least Squares and Least Absolute Deviations cost functions on the previous dataset.
- The Least Squares fit is shown in black, while the Least Absolute Deviation fit is depicted in magenta.
- The latter fit is considerably better, since it does not exaggerate the large error produced by the single outlier.



Regression Quality Metrics

- Next, we discuss how to make predictions using a trained model.
- We also describe simple metrics for judging the quality of a trained regression model.

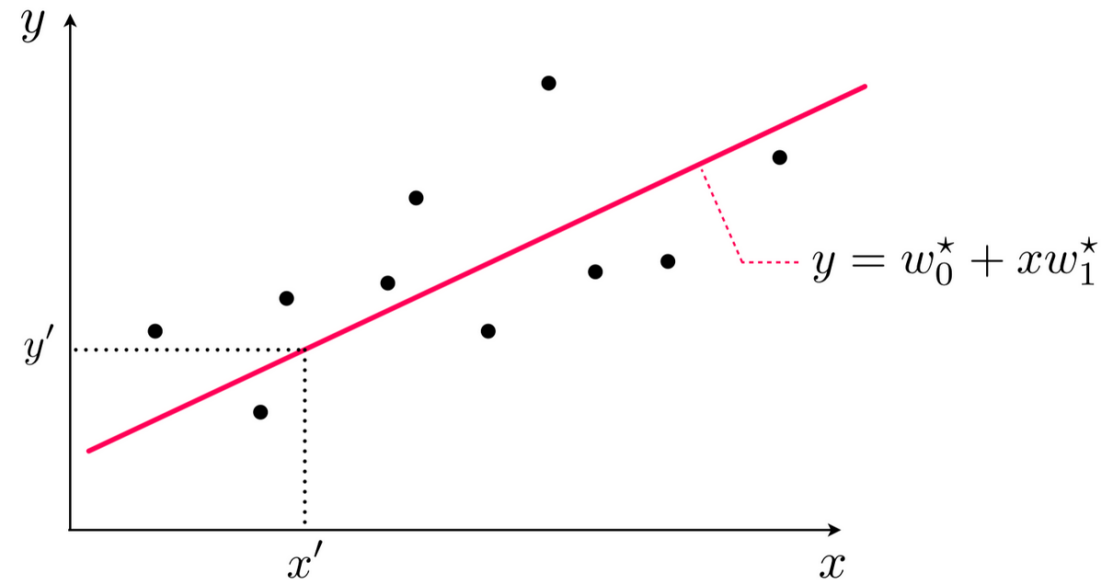
Making predictions using a trained model

- If we denote the optimal set of weights found by minimizing a regression cost function by \mathbf{w}^* , then we can write our fully tuned linear model as:

$$\text{model}(\mathbf{x}, \mathbf{w}^*) = \mathbf{x}^T \mathbf{w}^* = w_0^* + x_1 w_1^* + x_2 w_2^* + \cdots + x_N w_N^*.$$

- Given an input \mathbf{x}' , whether one from our training dataset or a new input, we predict its output y' by passing it along with our trained weights into our model as:

$$\text{model}(\mathbf{x}', \mathbf{w}^*) = y'$$



Judging the quality of a trained model

- Once we have successfully minimized the adopted cost function for linear regression, it is an easy matter to determine the quality of our regression model.
- We must simply evaluate a cost (error) function using our optimal weights.
- For example, we can evaluate the quality of this trained model using a Least Squares cost.
- This is especially natural to use, when we employ this cost in training.
- To do this, we plug in our trained model and dataset into the Least Squares cost, giving the *Mean Squared Error* (or MSE for short) of our trained model:

$$\text{MSE} = \frac{1}{P} \sum_{p=1}^P (\text{model}(\mathbf{x}_p, \mathbf{w}^*) - y_p)^2$$

- The name for this quality measurement describes precisely what the Least Squares cost computes: the *average* or *mean squared error*.

Judging the quality of a trained model (cont'd)

- In the same way, we can employ the Least Absolute Deviations cost to determine the quality of our trained model.
- Plugging in our trained model and dataset into this cost computes the *Mean Absolute Deviations* (or MAD for short).
- This is precisely what this cost function computes:

$$\text{MAD} = \frac{1}{P} \sum_{p=1}^P |\text{model}(\mathbf{x}_p, \mathbf{w}^*) - y_p|$$

Judging the quality of a trained model (cont'd)

- These two measurements differ in precisely the ways we have seen their respective cost functions differ:
the MSE measure is far more sensitive to outliers.
- Which measure one employs in practice can therefore depend on personal and/or occupational preference, or the nature of a problem at hand.
- Of course, in general, the *lower* one can make a quality measure, by proper tuning of model weights, the *better* the quality of the corresponding trained model and vice versa.
- However, the threshold for what one considers 'good' or 'great' performance can depend on: personal preference, an occupational or institutionally set benchmark, or some other problem-dependent concern.

Weighted Regression

- Regression cost functions are *summable losses over individual points*.
- Hence, we can weight individual points in order to emphasize or de-emphasize their importance to a regression model.
- This is called *weighted regression*.

Dealing with duplicates

- Imagine we have a linear regression dataset $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_P, y_P)$ that contains multiple copies of the same point.
- This can appear in a variety of contexts including:
 - Experimental data (e.g., in physics, medicine, etc.): If a repeated experiment produces the same result.
 - Metadata-type datasets (e.g., census, customer databases): Due to necessary/helpful pre-processing that quantizes (bins) input features to make human-in-the-loop analysis of the data and/or modeling easier.
- In such instances, *'duplicate' datapoints should not be thrown away*:
They accurately represent the true phenomenon under study.
- If we examine any regression cost function over such a dataset (i.e., one with repeated entries), we can see that it naturally *collapses* into a weighted version itself.

Dealing with duplicates (cont'd)

- For example, let us examine the Least Squares cost.
- Suppose that our first two datapoints (\mathbf{x}_1, y_1) and (\mathbf{x}_2, y_2) are *identical*.
- In this instance, using our model notation, the first two summands of our cost function (in the first two datapoints) can be combined since they will always take on the same value:

$$\begin{aligned} \frac{1}{P} \sum_{p=1}^P (\text{model}(\mathbf{x}_p, \mathbf{w}) - y_p)^2 &= \\ &= \frac{1}{P} (\text{model}(\mathbf{x}_1, \mathbf{w}) - y_1)^2 + \frac{1}{P} (\text{model}(\mathbf{x}_2, \mathbf{w}) - y_2)^2 + \frac{1}{P} \sum_{p=3}^P (\text{model}(\mathbf{x}_p, \mathbf{w}) - y_p)^2 \\ &= \frac{1}{P} 2 (\text{model}(\mathbf{x}_1, \mathbf{w}) - y_1)^2 + \frac{1}{P} \sum_{p=3}^P (\text{model}(\mathbf{x}_p, \mathbf{w}) - y_p)^2 \end{aligned}$$

- Here, we can see that the cost function naturally collapses, so that a repeated point in a dataset is represented *in the cost function* by a single weighted summand.

Dealing with duplicates (cont'd)

- Of course, this holds more generally as well.
- If a dataset has any number of identical points, then we can collapse the summands of the regression cost for each set of identical points just as we have seen previously.
- In general, this leads to the notion that each term in a regression cost can be *weighted* to reflect repeated points.

Dealing with duplicates (cont'd)

- We can write such a *weighted regression* Least Squares as:

$$g(\mathbf{w}) = \frac{1}{\sum_{p=1}^P \beta_p} \sum_{p=1}^P \beta_p (\text{model}(\mathbf{x}_p, \mathbf{w}) - y_p)^2,$$

where $\beta_1, \beta_2, \dots, \beta_P$ are *point-wise* weights.

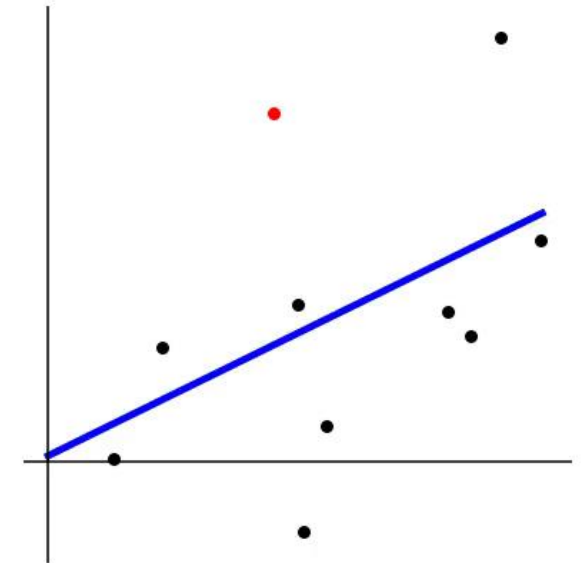
- That is, a unique point (\mathbf{x}_p, y_p) in the dataset has weight $\beta_p = 1$.
- If a point is repeated R times in the dataset, then one instance of it will have weight $\beta_p = R$, while the others will have weight $\beta_p = 0$.
- Since these weights are fixed (i.e., they are *not* parameters that need to be tuned, like \mathbf{w}), *we can minimize a weighted regression cost precisely as we would any other.*

Weighting points by confidence

- Weighted regression can also be employed when we want to weight each point based on our *confidence on the trustworthiness* of each datapoint.
- That assumes knowledge of the process generating the dataset.
- For example, if our dataset came in two batches, one from a trustworthy and another from a less trustworthy source, where some datapoints could be noisy or fallacious, we would want to weight datapoints from the trustworthy source more in our final regression.
- We can do this very easily, using precisely the weighted regression paradigm introduced above.
- So, now we set the weights $\beta_1, \beta_2, \dots, \beta_p$ ourselves, based on our *confidence of each point*.
- If we believe that a point is very trustworthy, we can set its corresponding weight β_p closer to 1.
- The less trustworthy a point is the smaller we set β_p in the range $[0,1]$, where the zero value implies that *we do not trust the point at all*.
- Making these weight selections, effectively we determine *how important each datapoint is* in the training of the model.

Weighting points by confidence (cont'd)

- Below, the result of increasing the confidence/weight β_p on a *single point* in a toy dataset and how this effects a fully trained regression model on a toy linear regression dataset is animated.
- This single point is colored red, and its diameter is shown increasing as the corresponding weight β_p increases.
- With each weighting, a weighted Least Squares cost is completely minimized over the entire dataset, and the resulting line fit to data is shown in blue.
- Observe: The higher the weighting of this single point, the more a linear regressor is incentivized to fit it.
- If the weight is increased enough, the fully trained regression model naturally starts fitting to this single datapoint alone, *disregarding all other points*.



Multi-Output Regression

- Thus far we have assumed that datapoints for linear regression consist of N dimensional vector valued inputs and scalar-valued outputs.
- So, a prototypical datapoint takes the form (\mathbf{x}_p, y_p) , where \mathbf{x}_p is an N dimensional input vector and a y_p is a scalar output.
- While this configuration covers most of the regression cases one may well encounter in practice,...
...it is possible to perform (linear) regression where both input and output are vector-valued.
- This is often called *multiple-output regression*.

Notation and Modeling

- Suppose now that the output of a regression dataset is *vector-valued*.
- That is, our datapoints $(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_P, \mathbf{y}_P)$ have both *vector-valued* input and output.
- The p th point $(\mathbf{x}_p, \mathbf{y}_p)$ has N dimensional input \mathbf{x}_p and associated C dimensional output \mathbf{y}_p .
- While, in principle, we can treat \mathbf{y}_p as a $C \times 1$ column vector,...
...to keep the formulas that follow looking similar to what we have already seen in the scalar case,...
...the input is treated as a $N \times 1$ column vector and the output as a $1 \times C$ row vector.

Notation and Modeling (cont'd)

- So, if we assume that a linear relationship holds between the input \mathbf{x}_p and just the c th dimension of the output $y_{c,p}$, we are using precisely the sort of regression framework discussed thus far.
- So, we can write:

$$\dot{\mathbf{x}}_p^T \mathbf{w}_c \approx y_{c,p}, \quad p = 1, \dots, P$$

where \mathbf{w}_c is a vector of weights and $\dot{\mathbf{x}}_p$ is the vector formed by stacking 1 on top of \mathbf{x}_p .

- If we then further assume that a linear relationship holds between the input and all C entries of the output, we can place each weight vector \mathbf{w}_c into the c th column of an $(N + 1) \times C$ weight matrix \mathbf{W} like this:

$$\mathbf{W} = \begin{bmatrix} w_{0,0} & w_{0,1} & w_{0,2} & \cdots & w_{0,C-1} \\ w_{1,0} & w_{1,1} & w_{1,2} & \cdots & w_{1,C-1} \\ w_{2,0} & w_{2,1} & w_{2,2} & \cdots & w_{2,C-1} \\ \vdots & \vdots & \vdots & \cdots & \vdots \\ w_{N,0} & w_{N,1} & w_{N,2} & \cdots & w_{N,C-1} \end{bmatrix}$$

Notation and Modeling (cont'd)

- Then, the entire set of C linear models can be written compactly as $1 \times C$ vector-matrix product:

$$\dot{\mathbf{x}}_p^T \mathbf{W} = [\dot{\mathbf{x}}_p^T \mathbf{w}_0 \quad \dot{\mathbf{x}}_p^T \mathbf{w}_1 \quad \cdots \quad \dot{\mathbf{x}}_p^T \mathbf{w}_{c-1}]$$

- This allows us to write the entire set of C linear relationships very compactly as:

$$\hat{\mathbf{y}}_p = \dot{\mathbf{x}}_p^T \mathbf{W} \approx \mathbf{y}_p, \quad p = 1, \dots, P$$

Cost Functions

- The thought process involved in deriving a regression cost function for the case of multi-output regression mirrors almost exactly the scalar-output.
- For example, to derive a Least Squares cost function, we begin by taking the difference of both sides of $\mathbf{x}_p^T \mathbf{W} \approx \mathbf{y}_p$.
- However, now the error $\boldsymbol{\epsilon}_p = \mathbf{x}_p^T \mathbf{W} - \mathbf{y}_p$ associated with the p th point has C values.
- The Least Squares cost function in this case is then the *average squared ℓ_2 norm* of each point's error:

$$g(\mathbf{W}) = \frac{1}{P} \sum_{p=1}^P \|\mathbf{x}_p^T \mathbf{W} - \mathbf{y}_p\|_2^2 = \frac{1}{P} \sum_{p=1}^P \sum_{c=0}^{C-1} (\mathbf{x}_p^T \mathbf{w}_c - y_{c,p})^2$$

- Likewise, the Least Absolute Deviations cost for our present case, using the ℓ_1 norm of each point's error, takes the form:

$$g(\mathbf{W}) = \frac{1}{P} \sum_{p=1}^P \|\mathbf{x}_p^T \mathbf{W} - \mathbf{y}_p\|_1 = \frac{1}{P} \sum_{p=1}^P \sum_{c=0}^{C-1} |\mathbf{x}_p^T \mathbf{w}_c - y_{c,p}|$$

Cost Functions (cont'd)

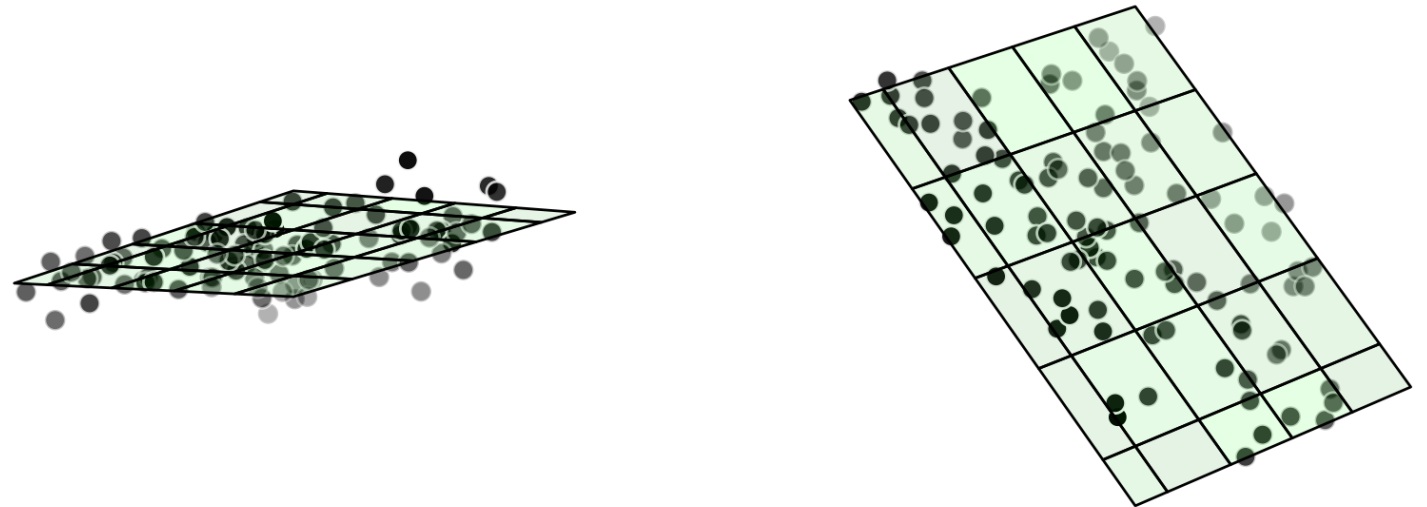
- Just like their scalar-valued versions, these cost functions *are always convex regardless of the dataset used*.
- They also decompose over the weights \mathbf{w}_c associated with each output dimension.
- For example, we can rewrite the Least Absolute Deviations cost by swapping the summands over P and C , as:

$$g(\mathbf{W}) = \sum_{c=0}^{C-1} \left(\frac{1}{P} \sum_{p=1}^P |\dot{\mathbf{x}}_p^T \mathbf{w}_c - y_{c,p}| \right) = \sum_{c=0}^{C-1} g_c(\mathbf{w}_c)$$

- Since the weights from each of the C subproblems do not interact,...
...if desired, we can minimize each $g_c(\cdot)$ for an optimal setting of \mathbf{w}_c *independently*,...
...and then take their sum to form the full cost function g .

Example. Fitting a linear model to a multi-output regression dataset

- Below, an example of multi-output linear regression is shown, using a toy dataset with input dimension $N = 2$ and output dimension $C = 2$.
- In each of the two panels of the figure, the input and one output value are plotted.
- The parameters of an appropriate linear model are tuned via minimizing the Least Squares cost, using gradient descent, with 200 steps and steplength $\alpha = 1$.
- The fully trained model is shown in green in each panel, by evaluating a fine mesh of points in the input region of the dataset.



Python implementation

- Multi-output regression cost functions can also be implemented in Python precisely as we have seen previously.
- For example, our linear model and Least Squares cost can be written as shown below.

```
# linear model
def model(x, w):
    a = w[0] + np.dot(x.T, w[1:])
    return a.T

# least squares cost
def least_squares(w):
    cost = np.sum((model(x, w) - y)**2)
    return cost/float(np.size(y))
```

- Implementing the model using `np.dot` as `np.dot(x.T, w[1:])` is especially efficient since NumPy's `np.dot` operation is far more efficient than constructing a linear combination in Python via an explicit `for` loop.

Resources

- Ch. 5: 5.1 – 5.7. Optional sections: 5.2.4, 5.6.3, 5.9.1.
- Animations
- [Optional Reading: Zaki and Meira \(2020\) Ch. 23: 23.1 – 23.3](#)