# Optional Lab: Python, NumPy and Vectorization

A brief introduction to some of the scientific computing used in this course. In particular the NumPy scientific computing package and its use with python.

# Outline

```
In [1]:  import numpy as np     # it is an unofficial standard to use np for numpy
         import time
```

## 1.1 Goals

In this lab, you will:

- Review the features of NumPy and Python that are used in Course 1

## 1.2 Useful References

- NumPy Documentation including a basic introduction: NumPy.org
- A challenging feature topic: NumPy Broadcasting

# 2 Python and NumPy

Python is the programming language we will be using in this course. It has a set of numeric data types and arithmetic operations. NumPy is a library that extends the base capabilities of python to add a richer data set including more numeric types, vectors, matrices, and many matrix functions. NumPy and python work together fairly seamlessly. Python

arithmetic operators work on NumPy data types and many NumPy functions will accept python data types.

# 3 Vectors

## 3.1 Abstract

Vectors, as you will use them in this course, are ordered arrays of numbers. In notation, vectors are denoted with lower case bold letters such as $\mathbf{x}$. The elements of a vector are all the same type. A vector does not, for example, contain both

$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_{n-1} \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$
$$\text{code} \qquad\qquad \text{lecture}$$

characters and numbers. The number of elements in the array is often referred to as the *dimension* though mathematicians may prefer *rank*. The vector shown has a dimension of $n$. The elements of a vector can be referenced with an index. In math settings, indexes typically run from 1 to n. In computer science and these labs, indexing will typically run from 0 to n-1. In notation, elements of a vector, when referenced individually will indicate the index in a subscript, for example, the $0^{th}$ element, of the vector $\mathbf{x}$ is $x_0$. Note, the x is not bold in this case.

## 3.2 NumPy Arrays

NumPy's basic data structure is an indexable, n-dimensional *array* containing elements of the same type ( `dtype` ). Right away, you may notice we have overloaded the term 'dimension'. Above, it was the number of elements in the vector, here, dimension refers to the number of indexes of an array. A one-dimensional or 1-D array has one index. In Course 1, we will represent vectors as NumPy 1-D arrays.

* 1-D array, shape (n,): n elements indexed [0] through [n-1]

## 3.3 Vector Creation

Data creation routines in NumPy will generally have a first parameter which is the shape of the object. This can either be a single value for a 1-D result or a tuple (n,m,...) specifying the shape of the result. Below are examples of creating vectors using these routines.

```python
In [2]:  # NumPy routines which allocate memory and fill arrays with value
a = np.zeros(4);                print(f"np.zeros(4) :   a = {a}, a shape = {a.
a = np.zeros((4,));             print(f"np.zeros(4,) :  a = {a}, a shape = {a.
a = np.random.random_sample(4); print(f"np.random.random_sample(4): a = {a}, a
```

```
np.zeros(4) :   a = [0. 0. 0. 0.], a shape = (4,), a data type = float64
np.zeros(4,) :  a = [0. 0. 0. 0.], a shape = (4,), a data type = float64
np.random.random_sample(4): a = [0.16054681 0.65264985 0.98302218 0.29864795],
a shape = (4,), a data type = float64
```

Some data creation routines do not take a shape tuple:

In [3]:
```python
# NumPy routines which allocate memory and fill arrays with value but do not a
a = np.arange(4.);                    print(f"np.arange(4.):     a = {a}, a shape =
a = np.random.rand(4);                print(f"np.random.rand(4): a = {a}, a shape =
```

```
np.arange(4.):     a = [0. 1. 2. 3.], a shape = (4,), a data type = float64
np.random.rand(4): a = [0.04467047 0.14854614 0.24949402 0.02417095], a shape
= (4,), a data type = float64
```

values can be specified manually as well.

In [4]:
```python
# NumPy routines which allocate memory and fill with user specified values
a = np.array([5,4,3,2]);  print(f"np.array([5,4,3,2]):  a = {a},     a shape =
a = np.array([5.,4,3,2]); print(f"np.array([5.,4,3,2]): a = {a}, a shape = {a.
```

```
np.array([5,4,3,2]):  a = [5 4 3 2],     a shape = (4,), a data type = int64
np.array([5.,4,3,2]): a = [5. 4. 3. 2.], a shape = (4,), a data type = float64
```

These have all created a one-dimensional vector `a` with four elements. `a.shape` returns the dimensions. Here we see a.shape = `(4,)` indicating a 1-d array with 4 elements.

# 3.4 Operations on Vectors

Let's explore some operations using vectors.

## 3.4.1 Indexing

Elements of vectors can be accessed via indexing and slicing. NumPy provides a very complete set of indexing and slicing capabilities. We will explore only the basics needed for the course here. Reference Slicing and Indexing for more details.

**Indexing** means referring to *an element* of an array by its position within the array.

**Slicing** means getting a *subset* of elements from an array based on their indices.

NumPy starts indexing at zero so the 3rd element of an vector **a** is `a[2]` .

In [5]:
```python
#vector indexing operations on 1-D vectors
a = np.arange(10)
print(a)

#access an element
print(f"a[2].shape: {a[2].shape} a[2]  = {a[2]}, Accessing an element returns

# access the last element, negative indexes count from the end
print(f"a[-1] = {a[-1]}")

#indexes must be within the range of the vector or they will produce and error
try:
    c = a[10]
except Exception as e:
    print("The error message you'll see is:")
    print(e)
```

```
[0 1 2 3 4 5 6 7 8 9]
a[2].shape: () a[2]  = 2, Accessing an element returns a scalar
a[-1] = 9
The error message you'll see is:
index 10 is out of bounds for axis 0 with size 10
```

## 3.4.2 Slicing

Slicing creates an array of indices using a set of three values ( `start:stop:step` ). A subset of values is also valid. Its use is best explained by example:

```
In [6]:  #vector slicing operations
         a = np.arange(10)
         print(f"a         = {a}")

         #access 5 consecutive elements (start:stop:step)
         c = a[2:7:1];    print("a[2:7:1] = ", c)

         # access 3 elements separated by two
         c = a[2:7:2];    print("a[2:7:2] = ", c)

         # access all elements index 3 and above
         c = a[3:];       print("a[3:]    = ", c)

         # access all elements below index 3
         c = a[:3];       print("a[:3]    = ", c)

         # access all elements
         c = a[:];        print("a[:]     = ", c)
```

```
a         = [0 1 2 3 4 5 6 7 8 9]
a[2:7:1] =  [2 3 4 5 6]
a[2:7:2] =  [2 4 6]
a[3:]    =  [3 4 5 6 7 8 9]
a[:3]    =  [0 1 2]
a[:]     =  [0 1 2 3 4 5 6 7 8 9]
```

## 3.4.3 Single vector operations

There are a number of useful operations that involve operations on a single vector.

```
In [7]:  a = np.array([1,2,3,4])
         print(f"a             : {a}")
         # negate elements of a
         b = -a
         print(f"b = -a        : {b}")

         # sum all elements of a, returns a scalar
         b = np.sum(a)
         print(f"b = np.sum(a) : {b}")

         b = np.mean(a)
         print(f"b = np.mean(a): {b}")

         b = a**2
         print(f"b = a**2      : {b}")
```

```
a               : [1 2 3 4]
b = -a          : [-1 -2 -3 -4]
b = np.sum(a) : 10
b = np.mean(a): 2.5
b = a**2        : [ 1  4  9 16]
```

### 3.4.4 Vector Vector element-wise operations

Most of the NumPy arithmetic, logical and comparison operations apply to vectors as well. These operators work on an element-by-element basis. For example

$$\mathbf{a} + \mathbf{b} = \sum_{i=0}^{n-1} a_i + b_i$$

```
In [8]:  a = np.array([ 1, 2, 3, 4])
         b = np.array([-1,-2, 3, 4])
         print(f"Binary operators work element wise: {a + b}")
```

```
Binary operators work element wise: [0 0 6 8]
```

Of course, for this to work correctly, the vectors must be of the same size:

```
In [9]:  #try a mismatched vector operation
         c = np.array([1, 2])
         try:
             d = a + c
         except Exception as e:
             print("The error message you'll see is:")
             print(e)
```

```
The error message you'll see is:
operands could not be broadcast together with shapes (4,) (2,)
```

### 3.4.5 Scalar Vector operations

Vectors can be 'scaled' by scalar values. A scalar value is just a number. The scalar multiplies all the elements of the vector.

```
In [10]:  a = np.array([1, 2, 3, 4])

          # multiply a by a scalar
          b = 5 * a
          print(f"b = 5 * a : {b}")
```

```
b = 5 * a : [ 5 10 15 20]
```

### 3.4.6 Vector Vector dot product

The dot product is a mainstay of Linear Algebra and NumPy. This is an operation used extensively in this course and should be well understood. The dot product is shown below.

# Vector Dot Product

$$a \cdot b = \sum_{i=0}^{n-1} a_i b_i$$

$$a \cdot b = \begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ a_3 \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ b_2 \\ b_3 \end{bmatrix} = [\quad + \quad + \quad + \quad] =$$

The dot product multiplies the values in two vectors element-wise and then sums the result. Vector dot product requires the dimensions of the two vectors to be the same.

Let's implement our own version of the dot product below:

**Using a for loop**, implement a function which returns the dot product of two vectors. The function to return given inputs $a$ and $b$:

$$x = \sum_{i=0}^{n-1} a_i b_i$$

Assume both `a` and `b` are the same shape.

```
In [11]: def my_dot(a, b):
             """
             Compute the dot product of two vectors

             Args:
               a (ndarray (n,)):  input vector
               b (ndarray (n,)):  input vector with same dimension as a

             Returns:
               x (scalar):
             """
             x=0
             for i in range(a.shape[0]):
                 x = x + a[i] * b[i]
             return x
```

```
In [12]: # test 1-D
         a = np.array([1, 2, 3, 4])
         b = np.array([-1, 4, 3, 2])
         print(f"my_dot(a, b) = {my_dot(a, b)}")
```

```
my_dot(a, b) = 24
```

Note, the dot product is expected to return a scalar value.

Let's try the same operations using `np.dot`.

```python
# test 1-D
a = np.array([1, 2, 3, 4])
b = np.array([-1, 4, 3, 2])
c = np.dot(a, b)
print(f"NumPy 1-D np.dot(a, b) = {c}, np.dot(a, b).shape = {c.shape} ")
c = np.dot(b, a)
print(f"NumPy 1-D np.dot(b, a) = {c}, np.dot(a, b).shape = {c.shape} ")
```

```
NumPy 1-D np.dot(a, b) = 24, np.dot(a, b).shape = ()
NumPy 1-D np.dot(b, a) = 24, np.dot(a, b).shape = ()
```

Above, you will note that the results for 1-D matched our implementation.

## 3.4.7 The Need for Speed: vector vs for loop

We utilized the NumPy library because it improves speed memory efficiency. Let's demonstrate:

```python
np.random.seed(1)
a = np.random.rand(10000000)  # very large arrays
b = np.random.rand(10000000)

tic = time.time()  # capture start time
c = np.dot(a, b)
toc = time.time()  # capture end time

print(f"np.dot(a, b) =  {c:.4f}")
print(f"Vectorized version duration: {1000*(toc-tic):.4f} ms ")

tic = time.time()  # capture start time
c = my_dot(a,b)
toc = time.time()  # capture end time

print(f"my_dot(a, b) =  {c:.4f}")
print(f"loop version duration: {1000*(toc-tic):.4f} ms ")

del(a);del(b)  #remove these big arrays from memory
```

```
np.dot(a, b) =  2501072.5817
Vectorized version duration: 29.9268 ms
my_dot(a, b) =  2501072.5817
loop version duration: 2137.4848 ms
```

So, vectorization provides a large speed up in this example. This is because NumPy makes better use of available data parallelism in the underlying hardware. GPU's and modern CPU's implement Single Instruction, Multiple Data (SIMD) pipelines allowing multiple operations to be issued in parallel. This is critical in Machine Learning where the data sets are often very large.

## 3.4.8 Vector Vector operations in Course 1

Vector Vector operations will appear frequently in course 1. Here is why:

- Going forward, our examples will be stored in an array, `X_train` of dimension (m,n). This will be explained more in context, but here it is important to note it is a 2 Dimensional array or matrix (see next section on matrices).
- `w` will be a 1-dimensional vector of shape (n,).

- we will perform operations by looping through the examples, extracting each example to work on individually by indexing X. For example: `X[i]`
- `X[i]` returns a value of shape (n,), a 1-dimensional vector. Consequently, operations involving `X[i]` are often vector-vector.

That is a somewhat lengthy explanation, but aligning and understanding the shapes of your operands is important when performing vector operations.

```
In [15]:  # show common Course 1 example
          X = np.array([[1],[2],[3],[4]])
          w = np.array([2])
          c = np.dot(X[1], w)

          print(f"X[1] has shape {X[1].shape}")
          print(f"w has shape {w.shape}")
          print(f"c has shape {c.shape}")
```
```
X[1] has shape (1,)
w has shape (1,)
c has shape ()
```

# 4 Matrices

## 4.1 Abstract

Matrices, are two dimensional arrays. The elements of a matrix are all of the same type. In notation, matrices are denoted with capitol, bold letter such as $\mathbf{X}$. In this and other labs, `m` is often the number of rows and `n` the number of columns. The elements of a matrix can be referenced with a two dimensional index. In math settings, numbers in the index typically run from 1 to n. In computer science and these labs, indexing will run from 0 to n-1.

$$X = \begin{bmatrix} x_{00} & x_{01} & \cdots & x_{0(n-1)} \\ x_{10} & x_{12} & \cdots & x_{1(n-1)} \\ \vdots & \vdots & \vdots & \vdots \\ x_{(m-1)0} & x_{(m-1)1} & \cdots & x_{(m-1)(n-1)} \end{bmatrix}_{\text{code}} \quad X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{1n} \\ \vdots & \vdots & \vdots & \vdots \\ x_{m1} & x_{m2} & \cdots & x_{mn} \end{bmatrix}_{\text{lecture}}$$

Generic Matrix Notation, 1st index is row, 2nd is column

## 4.2 NumPy Arrays

NumPy's basic data structure is an indexable, n-dimensional *array* containing elements of the same type ( `dtype` ). These were described earlier. Matrices have a two-dimensional (2-D) index [m,n].

In Course 1, 2-D matrices are used to hold training data. Training data is $m$ examples by $n$ features creating an (m,n) array. Course 1 does not do operations directly on matrices but typically extracts an example as a vector and operates on that. Below you will review:

- data creation

- slicing and indexing

## 4.3 Matrix Creation

The same functions that created 1-D vectors will create 2-D or n-D arrays. Here are some examples

Below, the shape tuple is provided to achieve a 2-D result. Notice how NumPy uses brackets to denote each dimension. Notice further than NumPy, when printing, will print one row per line.

```python
In [16]: a = np.zeros((1, 5))
         print(f"a shape = {a.shape}, a = {a}")

         a = np.zeros((2, 1))
         print(f"a shape = {a.shape}, a = {a}")

         a = np.random.random_sample((1, 1))
         print(f"a shape = {a.shape}, a = {a}")
```

```
a shape = (1, 5), a = [[0. 0. 0. 0. 0.]]
a shape = (2, 1), a = [[0.]
 [0.]]
a shape = (1, 1), a = [[0.44236513]]
```

One can also manually specify data. Dimensions are specified with additional brackets matching the format in the printing above.

```python
In [17]: # NumPy routines which allocate memory and fill with user specified values
         a = np.array([[5], [4], [3]]);   print(f" a shape = {a.shape}, np.array: a = {
         a = np.array([[5],    # One can also
                       [4],    # separate values
                       [3]]); #into separate rows
         print(f" a shape = {a.shape}, np.array: a = {a}")
```

```
a shape = (3, 1), np.array: a = [[5]
[4]
[3]]
a shape = (3, 1), np.array: a = [[5]
[4]
[3]]
```

## 4.4 Operations on Matrices

Let's explore some operations using matrices.

### 4.4.1 Indexing

Matrices include a second index. The two indexes describe [row, column]. Access can either return an element or a row/column. See below:

```python
In [18]: #vector indexing operations on matrices
         a = np.arange(6).reshape(-1, 2)   #reshape is a convenient way to create matri
         print(f"a.shape: {a.shape}, \na= {a}")
```

```
#access an element
print(f"\na[2,0].shape:    {a[2, 0].shape}, a[2,0] = {a[2, 0]},     type(a[2,0]

#access a row
print(f"a[2].shape:   {a[2].shape}, a[2]   = {a[2]}, type(a[2])   = {type(a[2]
```

```
a.shape: (3, 2),
a= [[0 1]
 [2 3]
 [4 5]]

a[2,0].shape:   (), a[2,0] = 4,     type(a[2,0]) = <class 'numpy.int64'> Acces
sing an element returns a scalar

a[2].shape:   (2,), a[2]   = [4 5], type(a[2])   = <class 'numpy.ndarray'>
```

It is worth drawing attention to the last example. Accessing a matrix by just specifying the row will return a *1-D vector*.

**Reshape**

The previous example used [reshape](#) to shape the array.

```
a = np.arange(6).reshape(-1, 2)
```

This line of code first created a *1-D Vector* of six elements. It then reshaped that vector into a *2-D* array using the reshape command. This could have been written:

```
a = np.arange(6).reshape(3, 2)
```

To arrive at the same 3 row, 2 column array. The -1 argument tells the routine to compute the number of rows given the size of the array and the number of columns.

## 4.4.2 Slicing

Slicing creates an array of indices using a set of three values ( `start:stop:step` ). A subset of values is also valid. Its use is best explained by example:

In [19]:
```
#vector 2-D slicing operations
a = np.arange(20).reshape(-1, 10)
print(f"a = \n{a}")

#access 5 consecutive elements (start:stop:step)
print("a[0, 2:7:1] = ", a[0, 2:7:1], ",  a[0, 2:7:1].shape =", a[0, 2:7:1].sha

#access 5 consecutive elements (start:stop:step) in two rows
print("a[:, 2:7:1] = \n", a[:, 2:7:1], ",  a[:, 2:7:1].shape =", a[:, 2:7:1].s

# access all elements
print("a[:,:] = \n", a[:,:], ",  a[:,:].shape =", a[:,:].shape)

# access all elements in one row (very common usage)
print("a[1,:] = ", a[1,:], ",  a[1,:].shape =", a[1,:].shape, "a 1-D array")
# same as
print("a[1]   = ", a[1],   ",  a[1].shape   =", a[1].shape, "a 1-D array")
```

```
a =
[[ 0  1  2  3  4  5  6  7  8  9]
 [10 11 12 13 14 15 16 17 18 19]]
a[0, 2:7:1] =  [2 3 4 5 6] ,  a[0, 2:7:1].shape = (5,) a 1-D array
a[:, 2:7:1] =
 [[ 2  3  4  5  6]
 [12 13 14 15 16]] ,  a[:, 2:7:1].shape = (2, 5) a 2-D array
a[:,:] =
 [[ 0  1  2  3  4  5  6  7  8  9]
 [10 11 12 13 14 15 16 17 18 19]] ,  a[:,:].shape = (2, 10)
a[1,:] =  [10 11 12 13 14 15 16 17 18 19] ,  a[1,:].shape = (10,) a 1-D array
a[1]   =  [10 11 12 13 14 15 16 17 18 19] ,  a[1].shape  = (10,) a 1-D array
```

## Congratulations!

In this lab you mastered the features of Python and NumPy that are needed for Course 1.

In [ ]:

# Optional Lab: Multiple Variable Linear Regression

In this lab, you will extend the data structures and previously developed routines to support multiple features. Several routines are updated making the lab appear lengthy, but it makes minor adjustments to previous routines making it quick to review.

# Outline

## 1.1 Goals

- Extend our regression model routines to support multiple features
  - Extend data structures to support multiple features
  - Rewrite prediction, cost and gradient routines to support multiple features
  - Utilize NumPy `np.dot` to vectorize their implementations for speed and simplicity

## 1.2 Tools

In this lab, we will make use of:

- NumPy, a popular library for scientific computing
- Matplotlib, a popular library for plotting data

```python
In [1]: import copy, math
import numpy as np
import matplotlib.pyplot as plt
plt.style.use('./deeplearning.mplstyle')
np.set_printoptions(precision=2)  # reduced display precision on numpy arrays
```

## 1.3 Notation

Here is a summary of some of the notation you will encounter, updated for multiple features.

| General Notation | Description | Python (if applicable) |
|---|---|---|
| $a$ | scalar, non bold | |
| $\mathbf{a}$ | vector, bold | |
| $\mathbf{A}$ | matrix, bold capital | |
| **Regression** | | |
| $\mathbf{X}$ | training example maxtrix | `X_train` |
| $\mathbf{y}$ | training example targets | `y_train` |
| $\mathbf{x}^{(i)}$, $y^{(i)}$ | $i_{th}$ Training Example | `X[i]` , `y[i]` |
| m | number of training examples | `m` |
| n | number of features in each example | `n` |
| $\mathbf{w}$ | parameter: weight, | `w` |
| $b$ | parameter: bias | `b` |
| $f_{\mathbf{w},b}(\mathbf{x}^{(i)})$ | The result of the model evaluation at $\mathbf{x^{(i)}}$ parameterized by $\mathbf{w}, b$: $f_{\mathbf{w},b}(\mathbf{x}^{(i)}) = \mathbf{w} \cdot \mathbf{x}^{(i)} + b$ | `f_wb` |

# 2 Problem Statement

You will use the motivating example of housing price prediction. The training dataset contains three examples with four features (size, bedrooms, floors and, age) shown in the table below. Note that, unlike the earlier labs, size is in sqft rather than 1000 sqft. This causes an issue, which you will solve in the next lab!

| Size (sqft) | Number of Bedrooms | Number of floors | Age of Home | Price (1000s dollars) |
|---|---|---|---|---|
| 2104 | 5 | 1 | 45 | 460 |
| 1416 | 3 | 2 | 40 | 232 |
| 852 | 2 | 1 | 35 | 178 |

You will build a linear regression model using these values so you can then predict the price for other houses. For example, a house with 1200 sqft, 3 bedrooms, 1 floor, 40 years

old.

Please run the following code cell to create your `X_train` and `y_train` variables.

```
In [2]:  X_train = np.array([[2104, 5, 1, 45], [1416, 3, 2, 40], [852, 2, 1, 35]])
         y_train = np.array([460, 232, 178])
```

## 2.1 Matrix X containing our examples

Similar to the table above, examples are stored in a NumPy matrix `X_train`. Each row of the matrix represents one example. When you have $m$ training examples ( $m$ is three in our example), and there are $n$ features (four in our example), $\mathbf{X}$ is a matrix with dimensions ($m$, $n$) (m rows, n columns).

$$\mathbf{X} = \begin{pmatrix} x_0^{(0)} & x_1^{(0)} & \cdots & x_{n-1}^{(0)} \\ x_0^{(1)} & x_1^{(1)} & \cdots & x_{n-1}^{(1)} \\ \cdots \\ x_0^{(m-1)} & x_1^{(m-1)} & \cdots & x_{n-1}^{(m-1)} \end{pmatrix}$$

notation:

- $\mathbf{x}^{(i)}$ is vector containing example i. $\mathbf{x}^{(i)} = (x_0^{(i)}, x_1^{(i)}, \cdots, x_{n-1}^{(i)})$
- $x_j^{(i)}$ is element j in example i. The superscript in parenthesis indicates the example number while the subscript represents an element.

Display the input data.

```
In [3]:  # data is stored in numpy array/matrix
         print(f"X Shape: {X_train.shape}, X Type:{type(X_train)})")
         print(X_train)
         print(f"y Shape: {y_train.shape}, y Type:{type(y_train)})")
         print(y_train)
```

```
X Shape: (3, 4), X Type:<class 'numpy.ndarray'>)
[[2104    5    1   45]
 [1416    3    2   40]
 [ 852    2    1   35]]
y Shape: (3,), y Type:<class 'numpy.ndarray'>)
[460 232 178]
```

## 2.2 Parameter vector w, b

- $\mathbf{w}$ is a vector with $n$ elements.
  - Each element contains the parameter associated with one feature.
  - in our dataset, n is 4.
  - notionally, we draw this as a column vector

$$\mathbf{w} = \begin{pmatrix} w_0 \\ w_1 \\ \ldots \\ w_{n-1} \end{pmatrix}$$

- $b$ is a scalar parameter.

For demonstration, $\mathbf{w}$ and $b$ will be loaded with some initial selected values that are near the optimal. $\mathbf{w}$ is a 1-D NumPy vector.

```
In [4]:  b_init = 785.1811367994083
         w_init = np.array([ 0.39133535, 18.75376741, -53.36032453, -26.42131618])
         print(f"w_init shape: {w_init.shape}, b_init type: {type(b_init)}")
```

```
w_init shape: (4,), b_init type: <class 'float'>
```

# 3 Model Prediction With Multiple Variables

The model's prediction with multiple variables is given by the linear model:

$$f_{\mathbf{w},b}(\mathbf{x}) = w_0 x_0 + w_1 x_1 + \ldots + w_{n-1} x_{n-1} + b \tag{1}$$

or in vector notation:

$$f_{\mathbf{w},b}(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b \tag{2}$$

where $\cdot$ is a vector `dot product`

To demonstrate the dot product, we will implement prediction using (1) and (2).

## 3.1 Single Prediction element by element

Our previous prediction multiplied one feature value by one parameter and added a bias parameter. A direct extension of our previous implementation of prediction to multiple features would be to implement (1) above using loop over each element, performing the multiply with its parameter and then adding the bias parameter at the end.

```
In [5]:  def predict_single_loop(x, w, b):
             """
             single predict using linear regression

             Args:
               x (ndarray): Shape (n,) example with multiple features
               w (ndarray): Shape (n,) model parameters
               b (scalar):  model parameter

             Returns:
               p (scalar):  prediction
             """
             n = x.shape[0]
             p = 0
             for i in range(n):
                 p_i = x[i] * w[i]
                 p = p + p_i
```

```
        p = p + b
    return p
```

```
# get a row from our training data
x_vec = X_train[0,:]
print(f"x_vec shape {x_vec.shape}, x_vec value: {x_vec}")

# make a prediction
f_wb = predict_single_loop(x_vec, w_init, b_init)
print(f"f_wb shape {f_wb.shape}, prediction: {f_wb}")
```

```
x_vec shape (4,), x_vec value: [2104    5    1   45]
f_wb shape (), prediction: 459.9999976194083
```

Note the shape of `x_vec`. It is a 1-D NumPy vector with 4 elements, (4,). The result, `f_wb` is a scalar.

## 3.2 Single Prediction, vector

Noting that equation (1) above can be implemented using the dot product as in (2) above. We can make use of vector operations to speed up predictions.

Recall from the Python/Numpy lab that NumPy `np.dot()` [link] can be used to perform a vector dot product.

```
def predict(x, w, b):
    """
    single predict using linear regression
    Args:
      x (ndarray): Shape (n,) example with multiple features
      w (ndarray): Shape (n,) model parameters
      b (scalar):             model parameter

    Returns:
      p (scalar):  prediction
    """
    p = np.dot(x, w) + b
    return p
```

```
# get a row from our training data
x_vec = X_train[0,:]
print(f"x_vec shape {x_vec.shape}, x_vec value: {x_vec}")

# make a prediction
f_wb = predict(x_vec,w_init, b_init)
print(f"f_wb shape {f_wb.shape}, prediction: {f_wb}")
```

```
x_vec shape (4,), x_vec value: [2104    5    1   45]
f_wb shape (), prediction: 459.99999761940825
```

The results and shapes are the same as the previous version which used looping. Going forward, `np.dot` will be used for these operations. The prediction is now a single statement. Most routines will implement it directly rather than calling a separate predict routine.

# 4 Compute Cost With Multiple Variables

The equation for the cost function with multiple variables $J(\mathbf{w}, b)$ is:

$$J(\mathbf{w}, b) = \frac{1}{2m} \sum_{i=0}^{m-1} (f_{\mathbf{w},b}(\mathbf{x}^{(i)}) - y^{(i)})^2 \tag{3}$$

where:

$$f_{\mathbf{w},b}(\mathbf{x}^{(i)}) = \mathbf{w} \cdot \mathbf{x}^{(i)} + b \tag{4}$$

In contrast to previous labs, $\mathbf{w}$ and $\mathbf{x}^{(i)}$ are vectors rather than scalars supporting multiple features.

Below is an implementation of equations (3) and (4). Note that this uses a *standard pattern for this course* where a for loop over all `m` examples is used.

```
In [9]: def compute_cost(X, y, w, b):
            """
            compute cost
            Args:
              X (ndarray (m,n)): Data, m examples with n features
              y (ndarray (m,)) : target values
              w (ndarray (n,)) : model parameters
              b (scalar)       : model parameter

            Returns:
              cost (scalar): cost
            """
            m = X.shape[0]
            cost = 0.0
            for i in range(m):
                f_wb_i = np.dot(X[i], w) + b           #(n,)(n,) = scalar (see np.dot)
                cost = cost + (f_wb_i - y[i])**2        #scalar
            cost = cost / (2 * m)                       #scalar
            return cost
```

```
In [10]: # Compute and display cost using our pre-chosen optimal parameters.
         cost = compute_cost(X_train, y_train, w_init, b_init)
         print(f'Cost at optimal w : {cost}')
```

Cost at optimal w : 1.5578904880036537e-12

**Expected Result**: Cost at optimal w : 1.5578904045996674e-12

# 5 Gradient Descent With Multiple Variables

Gradient descent for multiple variables:

$$\begin{align*}
&\text{repeat until convergence: } \{ \\
&\quad w_j = w_j - \alpha \frac{\partial J(\mathbf{w}, b)}{\partial w_j} \qquad \text{for } j = 0..n\text{-}1 \tag{5} \\
&\quad b \;\; = b - \alpha \frac{\partial J(\mathbf{w}, b)}{\partial b} \\
&\}
\end{align*}$$

where, n is the number of features, parameters $w_j, b$, are updated simultaneously and where

$$\frac{\partial J(\mathbf{w}, b)}{\partial w_j} = \frac{1}{m} \sum_{i=0}^{m-1} (f_{\mathbf{w},b}(\mathbf{x}^{(i)}) - y^{(i)}) x_j^{(i)} \tag{6}$$

$$\frac{\partial J(\mathbf{w}, b)}{\partial b} = \frac{1}{m} \sum_{i=0}^{m-1} (f_{\mathbf{w},b}(\mathbf{x}^{(i)}) - y^{(i)}) \tag{7}$$

- m is the number of training examples in the data set

- $f_{\mathbf{w},b}(\mathbf{x}^{(i)})$ is the model's prediction, while $y^{(i)}$ is the target value

# 5.1 Compute Gradient with Multiple Variables

An implementation for calculating the equations (6) and (7) is below. There are many ways to implement this. In this version, there is an

- outer loop over all m examples.
    - $\frac{\partial J(\mathbf{w},b)}{\partial b}$ for the example can be computed directly and accumulated
    - in a second loop over all n features:
        - $\frac{\partial J(\mathbf{w},b)}{\partial w_j}$ is computed for each $w_j$.

```
In [11]: def compute_gradient(X, y, w, b):
             """
             Computes the gradient for linear regression
             Args:
               X (ndarray (m,n)): Data, m examples with n features
               y (ndarray (m,)) : target values
               w (ndarray (n,)) : model parameters
               b (scalar)       : model parameter

             Returns:
               dj_dw (ndarray (n,)): The gradient of the cost w.r.t. the parameters w.
               dj_db (scalar):       The gradient of the cost w.r.t. the parameter b.
             """
             m,n = X.shape           #(number of examples, number of features)
             dj_dw = np.zeros((n,))
             dj_db = 0.

             for i in range(m):
                 err = (np.dot(X[i], w) + b) - y[i]
                 for j in range(n):
                     dj_dw[j] = dj_dw[j] + err * X[i, j]
                 dj_db = dj_db + err
             dj_dw = dj_dw / m
             dj_db = dj_db / m

             return dj_db, dj_dw
```

```
In [12]: #Compute and display gradient
         tmp_dj_db, tmp_dj_dw = compute_gradient(X_train, y_train, w_init, b_init)
         print(f'dj_db at initial w,b: {tmp_dj_db}')
         print(f'dj_dw at initial w,b: \n {tmp_dj_dw}')
```

```
dj_db at initial w,b: -1.673925169143331e-06
dj_dw at initial w,b:
 [-2.73e-03 -6.27e-06 -2.22e-06 -6.92e-05]
```

**Expected Result**:

dj_db at initial w,b: -1.6739251122999121e-06

dj_dw at initial w,b:

[-2.73e-03 -6.27e-06 -2.22e-06 -6.92e-05]

## 5.2 Gradient Descent With Multiple Variables

The routine below implements equation (5) above.

```
In [13]:  def gradient_descent(X, y, w_in, b_in, cost_function, gradient_function, alpha
              """
              Performs batch gradient descent to learn theta. Updates theta by taking
              num_iters gradient steps with learning rate alpha

              Args:
                X (ndarray (m,n))   : Data, m examples with n features
                y (ndarray (m,))    : target values
                w_in (ndarray (n,)) : initial model parameters
                b_in (scalar)       : initial model parameter
                cost_function       : function to compute cost
                gradient_function   : function to compute the gradient
                alpha (float)       : Learning rate
                num_iters (int)     : number of iterations to run gradient descent

              Returns:
                w (ndarray (n,)) : Updated values of parameters
                b (scalar)       : Updated value of parameter
              """

              # An array to store cost J and w's at each iteration primarily for graphin
              J_history = []
              w = copy.deepcopy(w_in)  #avoid modifying global w within function
              b = b_in

              for i in range(num_iters):

                  # Calculate the gradient and update the parameters
                  dj_db,dj_dw = gradient_function(X, y, w, b)   ##None

                  # Update Parameters using w, b, alpha and gradient
                  w = w - alpha * dj_dw             ##None
                  b = b - alpha * dj_db             ##None

                  # Save cost J at each iteration
                  if i<100000:      # prevent resource exhaustion
                      J_history.append( cost_function(X, y, w, b))

                  # Print cost every at intervals 10 times or as many iterations if < 10
                  if i% math.ceil(num_iters / 10) == 0:
                      print(f"Iteration {i:4d}: Cost {J_history[-1]:8.2f}   ")

              return w, b, J_history #return final w,b and J history for graphing
```

In the next cell you will test the implementation.

```
# initialize parameters
initial_w = np.zeros_like(w_init)
initial_b = 0.
# some gradient descent settings
iterations = 1000
alpha = 5.0e-7
# run gradient descent
w_final, b_final, J_hist = gradient_descent(X_train, y_train, initial_w, initi
                                            compute_cost, compute_grac
                                            alpha, iterations)
print(f"b,w found by gradient descent: {b_final:0.2f},{w_final} ")
m,_ = X_train.shape
for i in range(m):
    print(f"prediction: {np.dot(X_train[i], w_final) + b_final:0.2f}, target v
```
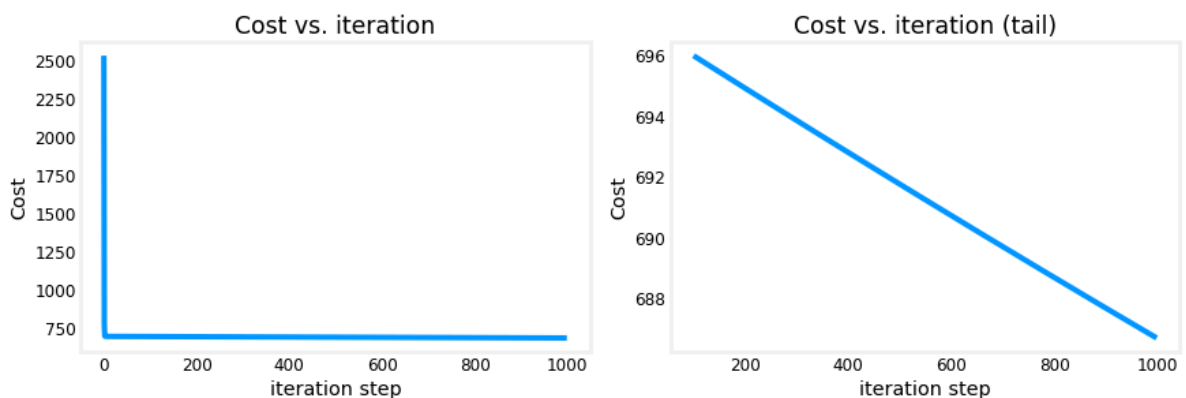
```
Iteration    0: Cost  2529.46
Iteration  100: Cost   695.99
Iteration  200: Cost   694.92
Iteration  300: Cost   693.86
Iteration  400: Cost   692.81
Iteration  500: Cost   691.77
Iteration  600: Cost   690.73
Iteration  700: Cost   689.71
Iteration  800: Cost   688.70
Iteration  900: Cost   687.69
b,w found by gradient descent: -0.00,[ 0.2   0.   -0.01 -0.07]
prediction: 426.19, target value: 460
prediction: 286.17, target value: 232
prediction: 171.47, target value: 178
```

**Expected Result**:

b,w found by gradient descent: -0.00,[ 0.2 0. -0.01 -0.07]

prediction: 426.19, target value: 460

prediction: 286.17, target value: 232

prediction: 171.47, target value: 178

```
# plot cost versus iteration
fig, (ax1, ax2) = plt.subplots(1, 2, constrained_layout=True, figsize=(12, 4))
ax1.plot(J_hist)
ax2.plot(100 + np.arange(len(J_hist[100:])), J_hist[100:])
ax1.set_title("Cost vs. iteration");  ax2.set_title("Cost vs. iteration (tail)
ax1.set_ylabel('Cost')             ;  ax2.set_ylabel('Cost')
ax1.set_xlabel('iteration step')   ;  ax2.set_xlabel('iteration step')
plt.show()
```



*These results are not inspiring*! Cost is still declining and our predictions are not very accurate. The next lab will explore how to improve on this.

# 6 Congratulations!

In this lab you:

- Redeveloped the routines for linear regression, now with multiple variables.
- Utilized NumPy `np.dot` to vectorize the implementations

```
In [ ]:
```

# Optional Lab: Feature scaling and Learning Rate (Multi-variable)

## Goals

In this lab you will:

- Utilize the multiple variables routines developed in the previous lab
- run Gradient Descent on a data set with multiple features
- explore the impact of the *learning rate alpha* on gradient descent
- improve performance of gradient descent by *feature scaling* using z-score normalization

## Tools

You will utilize the functions developed in the last lab as well as matplotlib and NumPy.

```
In [1]: import numpy as np
        import matplotlib.pyplot as plt
        from lab_utils_multi import  load_house_data, run_gradient_descent
        from lab_utils_multi import  norm_plot, plt_equal_scale, plot_cost_i_w
        from lab_utils_common import dlc
        np.set_printoptions(precision=2)
        plt.style.use('./deeplearning.mplstyle')
```

# Notation

| General Notation | Description | Python (if applicable) |
|---|---|---|
| $a$ | scalar, non bold | |
| $\mathbf{a}$ | vector, bold | |
| $\mathbf{A}$ | matrix, bold capital | |
| **Regression** | | |
| $\mathbf{X}$ | training example maxtrix | `X_train` |
| $\mathbf{y}$ | training example targets | `y_train` |
| $\mathbf{x}^{(i)}$, $y^{(i)}$ | $i_{th}$ Training Example | `X[i]` , `y[i]` |
| m | number of training examples | `m` |
| n | number of features in each example | `n` |
| $\mathbf{w}$ | parameter: weight, | `w` |
| $b$ | parameter: bias | `b` |
| $f_{\mathbf{w},b}(\mathbf{x}^{(i)})$ | The result of the model evaluation at $\mathbf{x}^{(i)}$ parameterized by $\mathbf{w}, b$: $f_{\mathbf{w},b}(\mathbf{x}^{(i)}) = \mathbf{w} \cdot \mathbf{x}^{(i)} + b$ | `f_wb` |
| $\frac{\partial J(\mathbf{w},b)}{\partial w_j}$ | the gradient or partial derivative of cost with respect to a parameter $w_j$ | `dj_dw[j]` |
| $\frac{\partial J(\mathbf{w},b)}{\partial b}$ | the gradient or partial derivative of cost with respect to a parameter $b$ | `dj_db` |

# Problem Statement

As in the previous labs, you will use the motivating example of housing price prediction. The training data set contains many examples with 4 features (size, bedrooms, floors and age) shown in the table below. Note, in this lab, the Size feature is in sqft while earlier labs utilized 1000 sqft. This data set is larger than the previous lab.

We would like to build a linear regression model using these values so we can then predict the price for other houses - say, a house with 1200 sqft, 3 bedrooms, 1 floor, 40 years old.
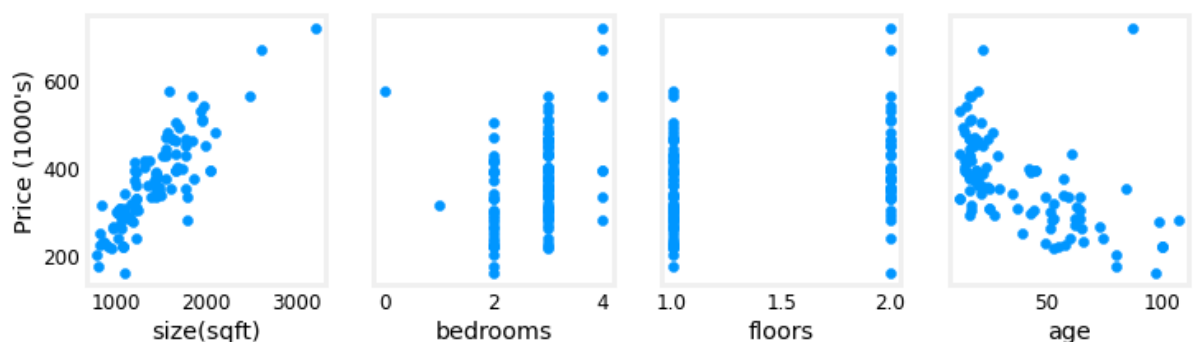
## Dataset:

| Size (sqft) | Number of Bedrooms | Number of floors | Age of Home | Price (1000s dollars) |
|---|---|---|---|---|
| 952 | 2 | 1 | 65 | 271.5 |
| 1244 | 3 | 2 | 64 | 232 |
| 1947 | 3 | 2 | 17 | 509.8 |
| ... | ... | ... | ... | ... |

```
In [2]:  # load the dataset
         X_train, y_train = load_house_data()
         X_features = ['size(sqft)','bedrooms','floors','age']
```

Let's view the dataset and its features by plotting each feature versus price.

```
In [3]:  fig,ax=plt.subplots(1, 4, figsize=(12, 3), sharey=True)
         for i in range(len(ax)):
             ax[i].scatter(X_train[:,i],y_train)
             ax[i].set_xlabel(X_features[i])
         ax[0].set_ylabel("Price (1000's)")
         plt.show()
```



Plotting each feature vs. the target, price, provides some indication of which features have the strongest influence on price. Above, increasing size also increases price. Bedrooms and floors don't seem to have a strong impact on price. Newer houses have higher prices than older houses.

# Gradient Descent With Multiple Variables

Here are the equations you developed in the last lab on gradient descent for multiple variables.:

$$\text{repeat until convergence: } \{$$

$$w_j := w_j - \alpha \frac{\partial J(\mathbf{w}, b)}{\partial w_j} \qquad \text{for j} = 0..\text{n-1} \tag{1}$$

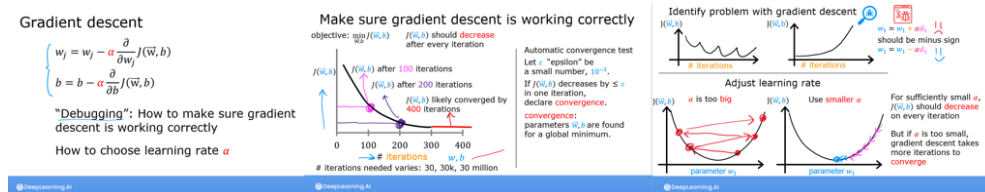$$b := b - \alpha \frac{\partial J(\mathbf{w}, b)}{\partial b}$$

$$\}$$

where, n is the number of features, parameters $w_j$, $b$, are updated simultaneously and where

$$\frac{\partial J(\mathbf{w}, b)}{\partial w_j} = \frac{1}{m} \sum_{i=0}^{m-1} (f_{\mathbf{w}, b}(\mathbf{x}^{(i)}) - y^{(i)}) x_j^{(i)} \tag{2}$$

$$\frac{\partial J(\mathbf{w}, b)}{\partial b} = \frac{1}{m} \sum_{i=0}^{m-1} (f_{\mathbf{w}, b}(\mathbf{x}^{(i)}) - y^{(i)}) \tag{3}$$

- m is the number of training examples in the data set

- $f_{\mathbf{w}, b}(\mathbf{x}^{(i)})$ is the model's prediction, while $y^{(i)}$ is the target value

# Learning Rate



The lectures discussed some of the issues related to setting the learning rate $\alpha$. The learning rate controls the size of the update to the parameters. See equation (1) above. It is shared by all the parameters.

Let's run gradient descent and try a few settings of $\alpha$ on our data set

## $\alpha$ = 9.9e-7

```
In [4]: #set alpha to 9.9e-7
_, _, hist = run_gradient_descent(X_train, y_train, 10, alpha = 9.9e-7)
```
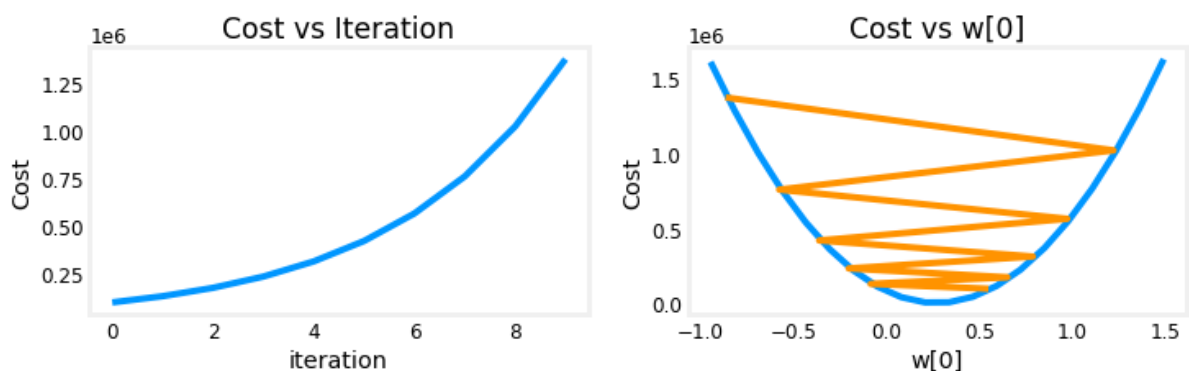
```
Iteration Cost          w0       w1       w2       w3       b        djdw0    djdw1
djdw2    djdw3    djdb
---------------------|--------|--------|--------|--------|--------|--------|-------
-|--------|--------|--------|
        0 9.55884e+04  5.5e-01  1.0e-03  5.1e-04  1.2e-02  3.6e-04 -5.5e+05 -1.0e+0
3 -5.2e+02 -1.2e+04 -3.6e+02
        1 1.28213e+05 -8.8e-02 -1.7e-04 -1.0e-04 -3.4e-03 -4.8e-05  6.4e+05  1.2e+0
3  6.2e+02  1.6e+04  4.1e+02
        2 1.72159e+05  6.5e-01  1.2e-03  5.9e-04  1.3e-02  4.3e-04 -7.4e+05 -1.4e+0
3 -7.0e+02 -1.7e+04 -4.9e+02
        3 2.31358e+05 -2.1e-01 -4.0e-04 -2.3e-04 -7.5e-03 -1.2e-04  8.6e+05  1.6e+0
3  8.3e+02  2.1e+04  5.6e+02
        4 3.11100e+05  7.9e-01  1.4e-03  7.1e-04  1.5e-02  5.3e-04 -1.0e+06 -1.8e+0
3 -9.5e+02 -2.3e+04 -6.6e+02
        5 4.18517e+05 -3.7e-01 -7.1e-04 -4.0e-04 -1.3e-02 -2.1e-04  1.2e+06  2.1e+0
3  1.1e+03  2.8e+04  7.5e+02
        6 5.63212e+05  9.7e-01  1.7e-03  8.7e-04  1.8e-02  6.6e-04 -1.3e+06 -2.5e+0
3 -1.3e+03 -3.1e+04 -8.8e+02
        7 7.58122e+05 -5.8e-01 -1.1e-03 -6.2e-04 -1.9e-02 -3.4e-04  1.6e+06  2.9e+0
3  1.5e+03  3.8e+04  1.0e+03
        8 1.02068e+06  1.2e+00  2.2e-03  1.1e-03  2.3e-02  8.3e-04 -1.8e+06 -3.3e+0
3 -1.7e+03 -4.2e+04 -1.2e+03
        9 1.37435e+06 -8.7e-01 -1.7e-03 -9.1e-04 -2.7e-02 -5.2e-04  2.1e+06  3.9e+0
3  2.0e+03  5.1e+04  1.4e+03
w,b found by gradient descent: w: [-0.87 -0.   -0.   -0.03], b: -0.00
```

It appears the learning rate is too high. The solution does not converge. Cost is *increasing* rather than decreasing. Let's plot the result:

```
In [5]: plot_cost_i_w(X_train, y_train, hist)
```



The plot on the right shows the value of one of the parameters, $w_0$. At each iteration, it is overshooting the optimal value and as a result, cost ends up *increasing* rather than approaching the minimum. Note that this is not a completely accurate picture as there are 4 parameters being modified each pass rather than just one. This plot is only showing $w_0$ with the other parameters fixed at benign values. In this and later plots you may notice the blue and orange lines being slightly off.

## $\alpha$ = 9e-7

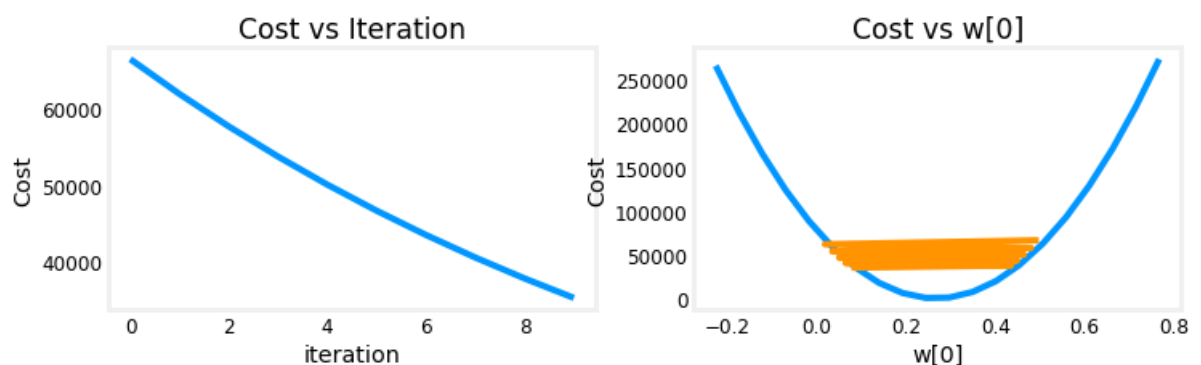Let's try a bit smaller value and see what happens.

```
In [6]: #set alpha to 9e-7
        _,_,hist = run_gradient_descent(X_train, y_train, 10, alpha = 9e-7)
```

```
Iteration Cost          w0        w1       w2       w3       b        djdw0     djdw1
djdw2    djdw3    djdb
---------------------|--------|--------|--------|--------|--------|--------|-------
-|--------|--------|--------|
        0 6.64616e+04  5.0e-01  9.1e-04  4.7e-04  1.1e-02  3.3e-04 -5.5e+05 -1.0e+0
3 -5.2e+02 -1.2e+04 -3.6e+02
        1 6.18990e+04  1.8e-02  2.1e-05  2.0e-06 -7.9e-04  1.9e-05  5.3e+05  9.8e+0
2  5.2e+02  1.3e+04  3.4e+02
        2 5.76572e+04  4.8e-01  8.6e-04  4.4e-04  9.5e-03  3.2e-04 -5.1e+05 -9.3e+0
2 -4.8e+02 -1.1e+04 -3.4e+02
        3 5.37137e+04  3.4e-02  3.9e-05  2.8e-06 -1.6e-03  3.8e-05  4.9e+05  9.1e+0
2  4.8e+02  1.2e+04  3.2e+02
        4 5.00474e+04  4.6e-01  8.2e-04  4.1e-04  8.0e-03  3.2e-04 -4.8e+05 -8.7e+0
2 -4.5e+02 -1.1e+04 -3.1e+02
        5 4.66388e+04  5.0e-02  5.6e-05  2.5e-06 -2.4e-03  5.6e-05  4.6e+05  8.5e+0
2  4.5e+02  1.2e+04  2.9e+02
        6 4.34700e+04  4.5e-01  7.8e-04  3.8e-04  6.4e-03  3.2e-04 -4.4e+05 -8.1e+0
2 -4.2e+02 -9.8e+03 -2.9e+02
        7 4.05239e+04  6.4e-02  7.0e-05  1.2e-06 -3.3e-03  7.3e-05  4.3e+05  7.9e+0
2  4.2e+02  1.1e+04  2.7e+02
        8 3.77849e+04  4.4e-01  7.5e-04  3.5e-04  4.9e-03  3.2e-04 -4.1e+05 -7.5e+0
2 -3.9e+02 -9.1e+03 -2.7e+02
        9 3.52385e+04  7.7e-02  8.3e-05 -1.1e-06 -4.2e-03  8.9e-05  4.0e+05  7.4e+0
2  3.9e+02  1.0e+04  2.5e+02
w,b found by gradient descent: w: [ 7.74e-02  8.27e-05 -1.06e-06 -4.20e-03], b: 0.0
0
```

Cost is decreasing throughout the run showing that alpha is not too large.

In [7]: `plot_cost_i_w(X_train, y_train, hist)`



On the left, you see that cost is decreasing as it should. On the right, you can see that $w_0$ is still oscillating around the minimum, but it is decreasing each iteration rather than increasing. Note above that `dj_dw[0]` changes sign with each iteration as `w[0]` jumps over the optimal value. This alpha value will converge. You can vary the number of iterations to see how it behaves.

## $\alpha$ = 1e-7

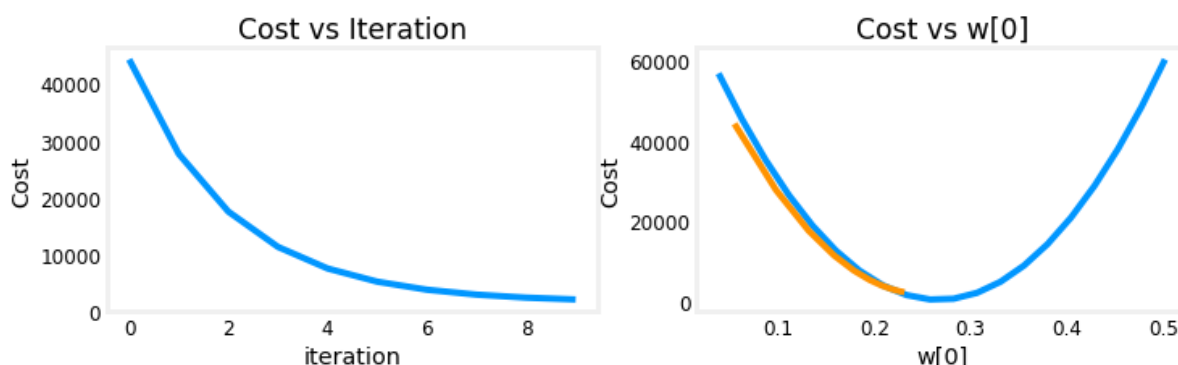Let's try a bit smaller value for $\alpha$ and see what happens.

In [8]:
```
#set alpha to 1e-7
_,_,hist = run_gradient_descent(X_train, y_train, 10, alpha = 1e-7)
```

```
Iteration Cost          w0       w1       w2       w3       b        djdw0    djdw1    djdw2    djdw3    djdb
---------------------|--------|--------|--------|--------|--------|--------|--------|--------|--------|--------|
        0 4.42313e+04  5.5e-02  1.0e-04  5.2e-05  1.2e-03  3.6e-05 -5.5e+05 -1.0e+03 -5.2e+02 -1.2e+04 -3.6e+02
        1 2.76461e+04  9.8e-02  1.8e-04  9.2e-05  2.2e-03  6.5e-05 -4.3e+05 -7.9e+02 -4.0e+02 -9.5e+03 -2.8e+02
        2 1.75102e+04  1.3e-01  2.4e-04  1.2e-04  2.9e-03  8.7e-05 -3.4e+05 -6.1e+02 -3.1e+02 -7.3e+03 -2.2e+02
        3 1.13157e+04  1.6e-01  2.9e-04  1.5e-04  3.5e-03  1.0e-04 -2.6e+05 -4.8e+02 -2.4e+02 -5.6e+03 -1.8e+02
        4 7.53002e+03  1.8e-01  3.3e-04  1.7e-04  3.9e-03  1.2e-04 -2.1e+05 -3.7e+02 -1.9e+02 -4.2e+03 -1.4e+02
        5 5.21639e+03  2.0e-01  3.5e-04  1.8e-04  4.2e-03  1.3e-04 -1.6e+05 -2.9e+02 -1.5e+02 -3.1e+03 -1.1e+02
        6 3.80242e+03  2.1e-01  3.8e-04  1.9e-04  4.5e-03  1.4e-04 -1.3e+05 -2.2e+02 -1.1e+02 -2.3e+03 -8.6e+01
        7 2.93826e+03  2.2e-01  3.9e-04  2.0e-04  4.6e-03  1.4e-04 -9.8e+04 -1.7e+02 -8.6e+01 -1.7e+03 -6.8e+01
        8 2.41013e+03  2.3e-01  4.1e-04  2.1e-04  4.7e-03  1.5e-04 -7.7e+04 -1.3e+02 -6.5e+01 -1.2e+03 -5.4e+01
        9 2.08734e+03  2.3e-01  4.2e-04  2.1e-04  4.8e-03  1.5e-04 -6.0e+04 -1.0e+02 -4.9e+01 -7.5e+02 -4.3e+01
w,b found by gradient descent: w: [2.31e-01 4.18e-04 2.12e-04 4.81e-03], b: 0.00
```
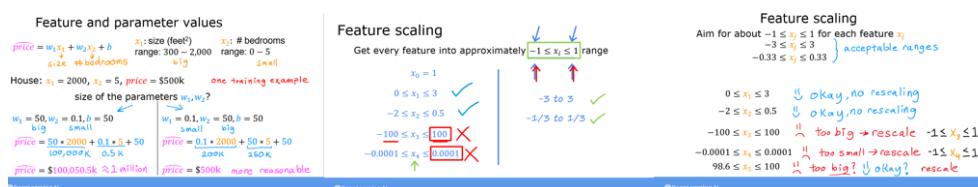
Cost is decreasing throughout the run showing that $\alpha$ is not too large.

```
In [9]:  plot_cost_i_w(X_train,y_train,hist)
```



On the left, you see that cost is decreasing as it should. On the right you can see that $w_0$ is decreasing without crossing the minimum. Note above that `dj_w0` is negative throughout the run. This solution will also converge, though not quite as quickly as the previous example.

## Feature Scaling



The lectures described the importance of rescaling the dataset so the features have a similar range. If you are interested in the details of why this is the case, click on the 'details' header below. If not, the section below will walk through an implementation of how to do feature scaling.

▶ **Details**

The lectures discussed three different techniques:

- Feature scaling, essentially dividing each positive feature by its maximum value, or more generally, rescale each feature by both its minimum and maximum values using (x-min)/(max-min). Both ways normalizes features to the range of -1 and 1, where the former method works for positive features which is simple and serves well for the lecture's example, and the latter method works for any features.
- Mean normalization: $x_i := \dfrac{x_i - \mu_i}{max - min}$
- Z-score normalization which we will explore below.

## z-score normalization

After z-score normalization, all features will have a mean of 0 and a standard deviation of 1.

To implement z-score normalization, adjust your input values as shown in this formula:

$$x_j^{(i)} = \frac{x_j^{(i)} - \mu_j}{\sigma_j} \tag{4}$$

where $j$ selects a feature or a column in the $\mathbf{X}$ matrix. $\mu_j$ is the mean of all the values for feature (j) and $\sigma_j$ is the standard deviation of feature (j).

$$\mu_j = \frac{1}{m} \sum_{i=0}^{m-1} x_j^{(i)} \tag{5}$$

$$\sigma_j^2 = \frac{1}{m} \sum_{i=0}^{m-1} (x_j^{(i)} - \mu_j)^2 \tag{6}$$

> **Implementation Note:** When normalizing the features, it is important

to store the values used for normalization - the mean value and the standard deviation used for the computations. After learning the parameters from the model, we often want to predict the prices of houses we have not seen before. Given a new x value (living room area and number of bed- rooms), we must first normalize x using the mean and standard deviation that we had previously computed from the training set.

**Implementation**

```
In [10]: def zscore_normalize_features(X):
    """
    computes  X, zcore normalized by column

    Args:
      X (ndarray (m,n))     : input data, m examples, n features

    Returns:
      X_norm (ndarray (m,n)): input normalized by column
      mu (ndarray (n,))     : mean of each feature
      sigma (ndarray (n,))  : standard deviation of each feature
    """
    # find the mean of each column/feature
    mu     = np.mean(X, axis=0)                 # mu will have shape (n,)
    # find the standard deviation of each column/feature
    sigma  = np.std(X, axis=0)                  # sigma will have shape (n,)
    # element-wise, subtract mu for that column from each example, divide by std fo
    X_norm = (X - mu) / sigma

    return (X_norm, mu, sigma)
```

```
#check our work
#from sklearn.preprocessing import scale
#scale(X_orig, axis=0, with_mean=True, with_std=True, copy=True)
```
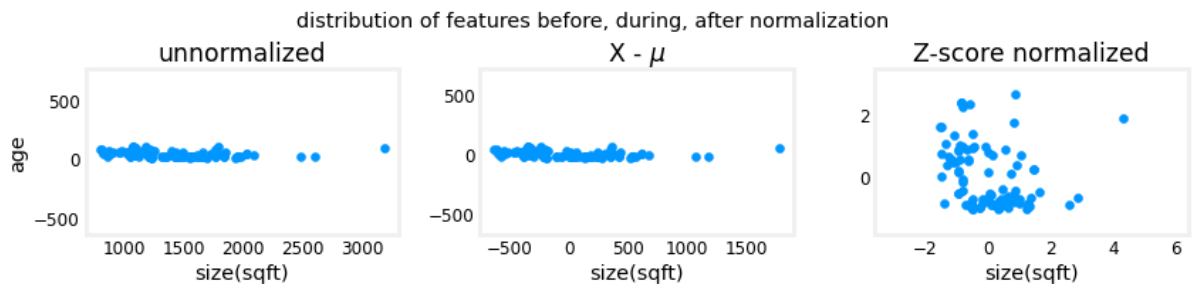
Let's look at the steps involved in Z-score normalization. The plot below shows the transformation step by step.

In [11]:
```
mu     = np.mean(X_train,axis=0)
sigma  = np.std(X_train,axis=0)
X_mean = (X_train - mu)
X_norm = (X_train - mu)/sigma

fig,ax=plt.subplots(1, 3, figsize=(12, 3))
ax[0].scatter(X_train[:,0], X_train[:,3])
ax[0].set_xlabel(X_features[0]); ax[0].set_ylabel(X_features[3]);
ax[0].set_title("unnormalized")
ax[0].axis('equal')

ax[1].scatter(X_mean[:,0], X_mean[:,3])
ax[1].set_xlabel(X_features[0]); ax[0].set_ylabel(X_features[3]);
ax[1].set_title(r"X - $\mu$")
ax[1].axis('equal')

ax[2].scatter(X_norm[:,0], X_norm[:,3])
ax[2].set_xlabel(X_features[0]); ax[0].set_ylabel(X_features[3]);
ax[2].set_title(r"Z-score normalized")
ax[2].axis('equal')
plt.tight_layout(rect=[0, 0.03, 1, 0.95])
fig.suptitle("distribution of features before, during, after normalization")
plt.show()
```



The plot above shows the relationship between two of the training set parameters, "age" and "size(sqft)". *These are plotted with equal scale*.

- Left: Unnormalized: The range of values or the variance of the 'size(sqft)' feature is much larger than that of age
- Middle: The first step removes the mean or average value from each feature. This leaves features that are centered around zero. It's difficult to see the difference for the 'age' feature, but 'size(sqft)' is clearly around zero.
- Right: The second step divides by the variance. This leaves both features centered at zero with a similar scale.

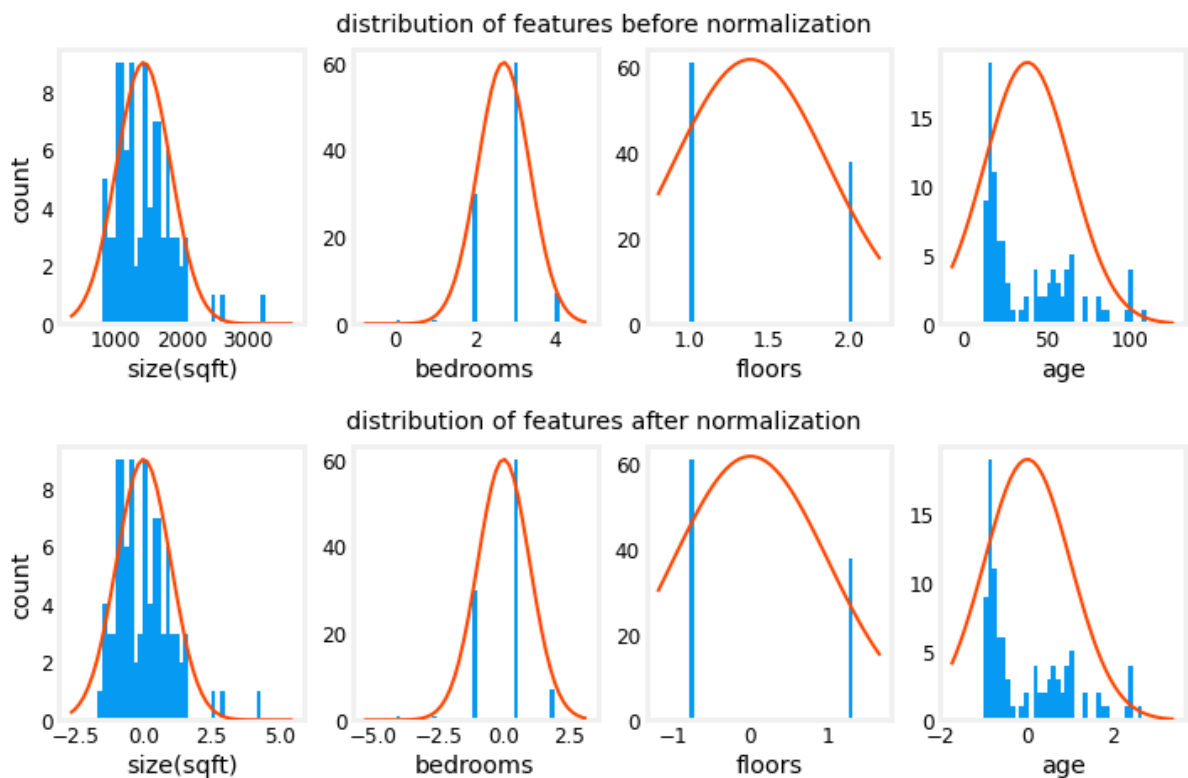Let's normalize the data and compare it to the original data.

In [12]:
```
# normalize the original features
X_norm, X_mu, X_sigma = zscore_normalize_features(X_train)
print(f"X_mu = {X_mu}, \nX_sigma = {X_sigma}")
print(f"Peak to Peak range by column in Raw        X:{np.ptp(X_train,axis=0)}")
print(f"Peak to Peak range by column in Normalized X:{np.ptp(X_norm,axis=0)}")
```

```
X_mu = [1.42e+03 2.72e+00 1.38e+00 3.84e+01],
X_sigma = [411.62   0.65   0.49  25.78]
Peak to Peak range by column in Raw         X:[2.41e+03 4.00e+00 1.00e+00 9.50e+01]
Peak to Peak range by column in Normalized X:[5.85 6.14 2.06 3.69]
```

The peak to peak range of each column is reduced from a factor of thousands to a factor of 2-3 by normalization.

```
In [13]:  fig,ax=plt.subplots(1, 4, figsize=(12, 3))
          for i in range(len(ax)):
              norm_plot(ax[i],X_train[:,i],)
              ax[i].set_xlabel(X_features[i])
          ax[0].set_ylabel("count");
          fig.suptitle("distribution of features before normalization")
          plt.show()
          fig,ax=plt.subplots(1,4,figsize=(12,3))
          for i in range(len(ax)):
              norm_plot(ax[i],X_norm[:,i],)
              ax[i].set_xlabel(X_features[i])
          ax[0].set_ylabel("count");
          fig.suptitle("distribution of features after normalization")

          plt.show()
```



distribution of features before normalization



distribution of features after normalization

Notice, above, the range of the normalized data (x-axis) is centered around zero and roughly +/- 2. Most importantly, the range is similar for each feature.

Let's re-run our gradient descent algorithm with normalized data. Note the **vastly larger value of alpha**. This will speed up gradient descent.

```
In [14]:  w_norm, b_norm, hist = run_gradient_descent(X_norm, y_train, 1000, 1.0e-1, )
```
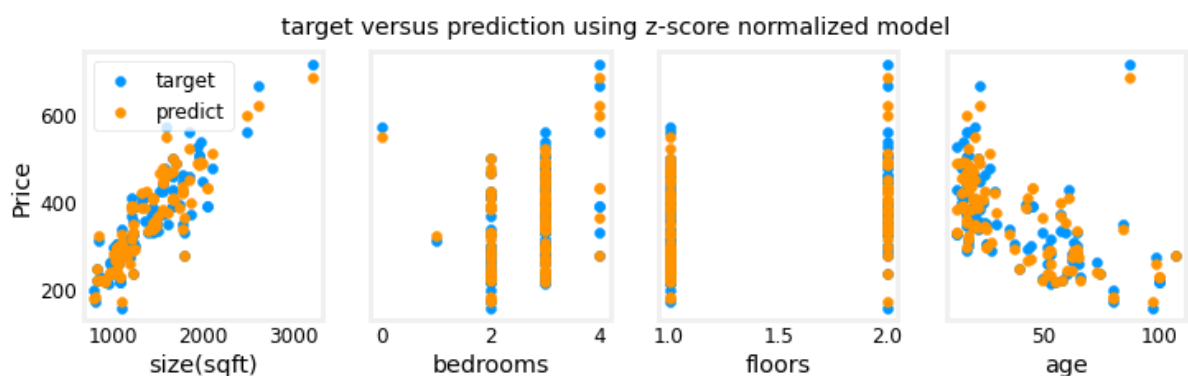
```
Iteration Cost          w0        w1        w2        w3        b         djdw0     djdw1
djdw2     djdw3     djdb
--------------------|--------|--------|--------|--------|--------|--------|-------
-|--------|--------|--------|
        0 5.76170e+04  8.9e+00   3.0e+00   3.3e+00  -6.0e+00   3.6e+01  -8.9e+01  -3.0e+0
1  -3.3e+01   6.0e+01  -3.6e+02
      100 2.21086e+02  1.1e+02  -2.0e+01  -3.1e+01  -3.8e+01   3.6e+02  -9.2e-01   4.5e-0
1   5.3e-01  -1.7e-01  -9.6e-03
      200 2.19209e+02  1.1e+02  -2.1e+01  -3.3e+01  -3.8e+01   3.6e+02  -3.0e-02   1.5e-0
2   1.7e-02  -6.0e-03  -2.6e-07
      300 2.19207e+02  1.1e+02  -2.1e+01  -3.3e+01  -3.8e+01   3.6e+02  -1.0e-03   5.1e-0
4   5.7e-04  -2.0e-04  -6.9e-12
      400 2.19207e+02  1.1e+02  -2.1e+01  -3.3e+01  -3.8e+01   3.6e+02  -3.4e-05   1.7e-0
5   1.9e-05  -6.6e-06  -2.7e-13
      500 2.19207e+02  1.1e+02  -2.1e+01  -3.3e+01  -3.8e+01   3.6e+02  -1.1e-06   5.6e-0
7   6.2e-07  -2.2e-07  -2.7e-13
      600 2.19207e+02  1.1e+02  -2.1e+01  -3.3e+01  -3.8e+01   3.6e+02  -3.7e-08   1.9e-0
8   2.1e-08  -7.3e-09  -2.6e-13
      700 2.19207e+02  1.1e+02  -2.1e+01  -3.3e+01  -3.8e+01   3.6e+02  -1.2e-09   6.2e-1
0   6.9e-10  -2.4e-10  -2.6e-13
      800 2.19207e+02  1.1e+02  -2.1e+01  -3.3e+01  -3.8e+01   3.6e+02  -4.1e-11   2.1e-1
1   2.3e-11  -8.1e-12  -2.6e-13
      900 2.19207e+02  1.1e+02  -2.1e+01  -3.3e+01  -3.8e+01   3.6e+02  -1.4e-12   6.9e-1
3   7.7e-13  -2.7e-13  -2.6e-13
w,b found by gradient descent: w: [110.56 -21.27 -32.71 -37.97], b: 363.16
```

The scaled features get very accurate results **much, much faster!**. Notice the gradient of each parameter is tiny by the end of this fairly short run. A learning rate of 0.1 is a good start for regression with normalized features. Let's plot our predictions versus the target values. Note, the prediction is made using the normalized feature while the plot is shown using the original feature values.

In [15]:
```python
#predict target using normalized features
m = X_norm.shape[0]
yp = np.zeros(m)
for i in range(m):
    yp[i] = np.dot(X_norm[i], w_norm) + b_norm

    # plot predictions and targets versus original features
fig,ax=plt.subplots(1,4,figsize=(12, 3),sharey=True)
for i in range(len(ax)):
    ax[i].scatter(X_train[:,i],y_train, label = 'target')
    ax[i].set_xlabel(X_features[i])
    ax[i].scatter(X_train[:,i],yp,color=dlc["dlorange"], label = 'predict')
ax[0].set_ylabel("Price"); ax[0].legend();
fig.suptitle("target versus prediction using z-score normalized model")
plt.show()
```


target versus prediction using z-score normalized model

The results look good. A few points to note:

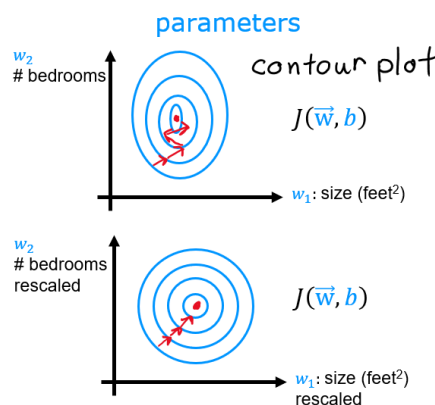- with multiple features, we can no longer have a single plot showing results versus features.

- when generating the plot, the normalized features were used. Any predictions using the parameters learned from a normalized training set must also be normalized.

**Prediction** The point of generating our model is to use it to predict housing prices that are not in the data set. Let's predict the price of a house with 1200 sqft, 3 bedrooms, 1 floor, 40 years old. Recall, that you must normalize the data with the mean and standard deviation derived when the training data was normalized.

```
In [16]: # First, normalize out example.
         x_house = np.array([1200, 3, 1, 40])
         x_house_norm = (x_house - X_mu) / X_sigma
         print(x_house_norm)
         x_house_predict = np.dot(x_house_norm, w_norm) + b_norm
         print(f" predicted price of a house with 1200 sqft, 3 bedrooms, 1 floor, 40 years o
```

```
[-0.53  0.43 -0.79  0.06]
 predicted price of a house with 1200 sqft, 3 bedrooms, 1 floor, 40 years old = $31
8709
```

### Cost Contours



Another way to view feature scaling is in terms of the cost contours. When feature scales do not match, the plot of cost versus parameters in a contour plot is asymmetric.

In the plot below, the scale of the parameters is matched. The left plot is the cost contour plot of w[0], the square feet versus w[1], the number of bedrooms before normalizing the features. The plot is so asymmetric, the curves completing the contours are not visible. In contrast, when the features are normalized, the cost contour is much more symmetric. The result is that updates to parameters during gradient descent can make equal progress for each parameter.

```
In [17]: plt_equal_scale(X_train, X_norm, y_train)
```



Cost contour with equal scale

## Congratulations!

In this lab you:

- utilized the routines for linear regression with multiple features you developed in previous labs
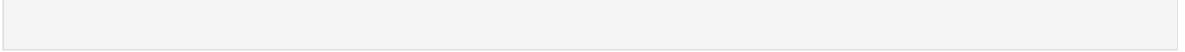- explored the impact of the learning rate $\alpha$ on convergence

- discovered the value of feature scaling using z-score normalization in speeding convergence
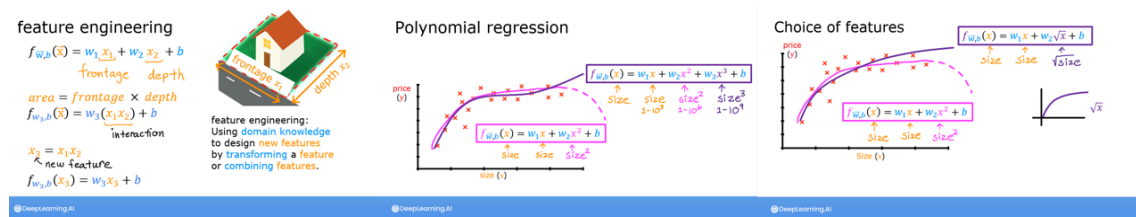
## Acknowledgments

The housing data was derived from the Ames Housing dataset compiled by Dean De Cock for use in data science education.

In [17]:

# Optional Lab: Feature Engineering and Polynomial Regression



## Goals

In this lab you will:

- explore feature engineering and polynomial regression which allows you to use the machinery of linear regression to fit very complicated, even very non-linear functions.

## Tools

You will utilize the function developed in previous labs as well as matplotlib and NumPy.

```
In [1]:  import numpy as np
         import matplotlib.pyplot as plt
         from lab_utils_multi import zscore_normalize_features, run_gradient_descent_feng
         np.set_printoptions(precision=2)  # reduced display precision on numpy arrays
```

# Feature Engineering and Polynomial Regression Overview

Out of the box, linear regression provides a means of building models of the form:

$$f_{\mathbf{w},b} = w_0 x_0 + w_1 x_1 + \ldots + w_{n-1} x_{n-1} + b \tag{1}$$

What if your features/data are non-linear or are combinations of features? For example, Housing prices do not tend to be linear with living area but penalize very small or very large houses resulting in the curves shown in the graphic above. How can we use the machinery of linear regression to fit this curve? Recall, the 'machinery' we have is the ability to modify the parameters $\mathbf{w}$, $\mathbf{b}$ in (1) to 'fit' the equation to the training data. However, no amount of adjusting of $\mathbf{w}$,$\mathbf{b}$ in (1) will achieve a fit to a non-linear curve.

## Polynomial Features

Above we were considering a scenario where the data was non-linear. Let's try using what we know so far to fit a non-linear curve. We'll start with a simple quadratic: $y = 1 + x^2$

You're familiar with all the routines we're using. They are available in the lab_utils.py file for review. We'll use `np.c_[..]` which is a NumPy routine to concatenate along the column
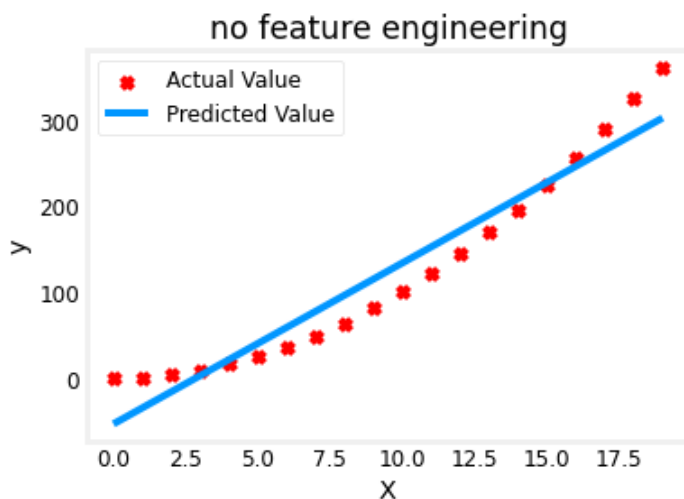
boundary.

```
In [2]:  # create target data
         x = np.arange(0, 20, 1)
         y = 1 + x**2
         X = x.reshape(-1, 1)

         model_w,model_b = run_gradient_descent_feng(X,y,iterations=1000, alpha = 1e-2)

         plt.scatter(x, y, marker='x', c='r', label="Actual Value"); plt.title("no feature e
         plt.plot(x,X@model_w + model_b, label="Predicted Value");  plt.xlabel("X"); plt.yla
```

```
Iteration        0, Cost: 1.65756e+03
Iteration      100, Cost: 6.94549e+02
Iteration      200, Cost: 5.88475e+02
Iteration      300, Cost: 5.26414e+02
Iteration      400, Cost: 4.90103e+02
Iteration      500, Cost: 4.68858e+02
Iteration      600, Cost: 4.56428e+02
Iteration      700, Cost: 4.49155e+02
Iteration      800, Cost: 4.44900e+02
Iteration      900, Cost: 4.42411e+02
w,b found by gradient descent: w: [18.7], b: -52.0834
```



Well, as expected, not a great fit. What is needed is something like $y = w_0x_0^2 + b$, or a **polynomial feature**. To accomplish this, you can modify the *input data* to *engineer* the needed features. If you swap the original data with a version that squares the $x$ value, then you can achieve $y = w_0x_0^2 + b$. Let's try it. Swap `X` for `X**2` below:

```
In [3]:  # create target data
         x = np.arange(0, 20, 1)
         y = 1 + x**2

         # Engineer features
         X = x**2       #<-- added engineered feature
```
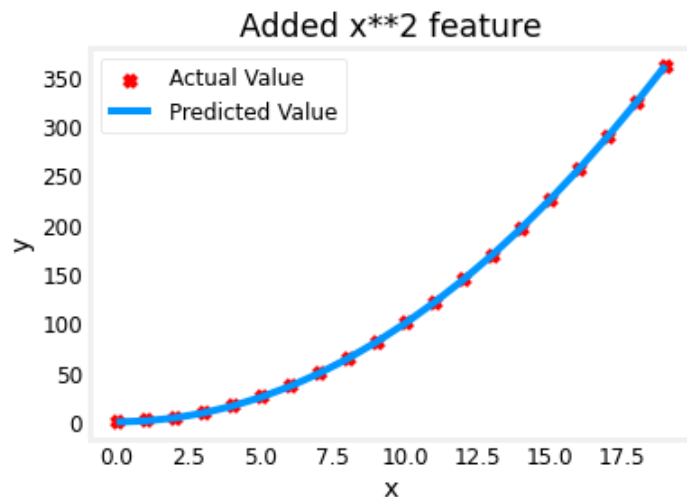
```
In [4]:  X = X.reshape(-1, 1)  #X should be a 2-D Matrix
         model_w,model_b = run_gradient_descent_feng(X, y, iterations=10000, alpha = 1e-5)

         plt.scatter(x, y, marker='x', c='r', label="Actual Value"); plt.title("Added x**2 f
         plt.plot(x, np.dot(X,model_w) + model_b, label="Predicted Value"); plt.xlabel("x");
```

```
Iteration         0, Cost: 7.32922e+03
Iteration      1000, Cost: 2.24844e-01
Iteration      2000, Cost: 2.22795e-01
Iteration      3000, Cost: 2.20764e-01
Iteration      4000, Cost: 2.18752e-01
Iteration      5000, Cost: 2.16758e-01
Iteration      6000, Cost: 2.14782e-01
Iteration      7000, Cost: 2.12824e-01
Iteration      8000, Cost: 2.10884e-01
Iteration      9000, Cost: 2.08962e-01
w,b found by gradient descent: w: [1.], b: 0.0490
```



Great! near perfect fit. Notice the values of $\mathbf{w}$ and b printed right above the graph: `w,b found by gradient descent: w: [1.], b: 0.0490`. Gradient descent modified our initial values of $\mathbf{w}, b$ to be (1.0,0.049) or a model of $y = 1 * x_0^2 + 0.049$, very close to our target of $y = 1 * x_0^2 + 1$. If you ran it longer, it could be a better match.

## Selecting Features

Above, we knew that an $x^2$ term was required. It may not always be obvious which features are required. One could add a variety of potential features to try and find the most useful. For example, what if we had instead tried : $y = w_0 x_0 + w_1 x_1^2 + w_2 x_2^3 + b$ ?

Run the next cells.

In [5]:
```python
# create target data
x = np.arange(0, 20, 1)
y = x**2

# engineer features .
X = np.c_[x, x**2, x**3]   #<-- added engineered feature
```

In [6]:
```python
model_w,model_b = run_gradient_descent_feng(X, y, iterations=10000, alpha=1e-7)

plt.scatter(x, y, marker='x', c='r', label="Actual Value"); plt.title("x, x**2, x**
plt.plot(x, X@model_w + model_b, label="Predicted Value"); plt.xlabel("x"); plt.yla
```
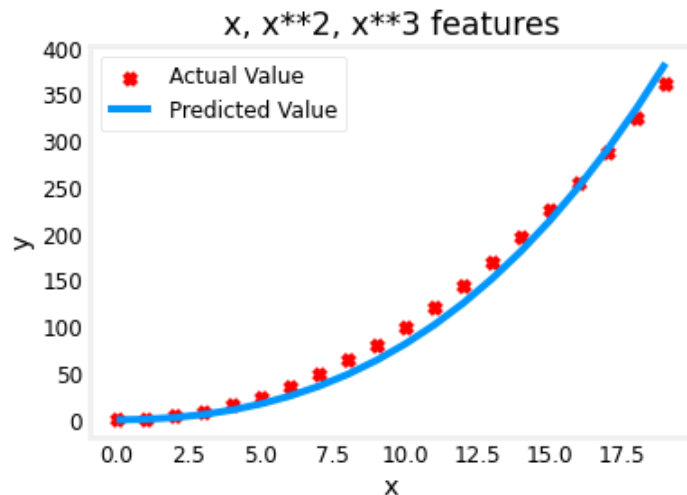
```
Iteration          0, Cost: 1.14029e+03
Iteration       1000, Cost: 3.28539e+02
Iteration       2000, Cost: 2.80443e+02
Iteration       3000, Cost: 2.39389e+02
Iteration       4000, Cost: 2.04344e+02
Iteration       5000, Cost: 1.74430e+02
Iteration       6000, Cost: 1.48896e+02
Iteration       7000, Cost: 1.27100e+02
Iteration       8000, Cost: 1.08495e+02
Iteration       9000, Cost: 9.26132e+01
w,b found by gradient descent: w: [0.08 0.54 0.03], b: 0.0106
```



Note the value of $\mathbf{w}$, `[0.08 0.54 0.03]` and b is `0.0106` .This implies the model after fitting/training is:

$$0.08x + 0.54x^2 + 0.03x^3 + 0.0106$$

Gradient descent has emphasized the data that is the best fit to the $x^2$ data by increasing the $w_1$ term relative to the others. If you were to run for a very long time, it would continue to reduce the impact of the other terms.

> Gradient descent is picking the 'correct' features for us by emphasizing its associated parameter

Let's review this idea:

- Intially, the features were re-scaled so they are comparable to each other
- less weight value implies less important/correct feature, and in extreme, when the weight becomes zero or very close to zero, the associated feature useful in fitting the model to the data.
- above, after fitting, the weight associated with the $x^2$ feature is much larger than the weights for $x$ or $x^3$ as it is the most useful in fitting the data.

## An Alternate View

Above, polynomial features were chosen based on how well they matched the target data. Another way to think about this is to note that we are still using linear regression once we have created new features. Given that, the best features will be linear relative to the target. This is best understood with an example.
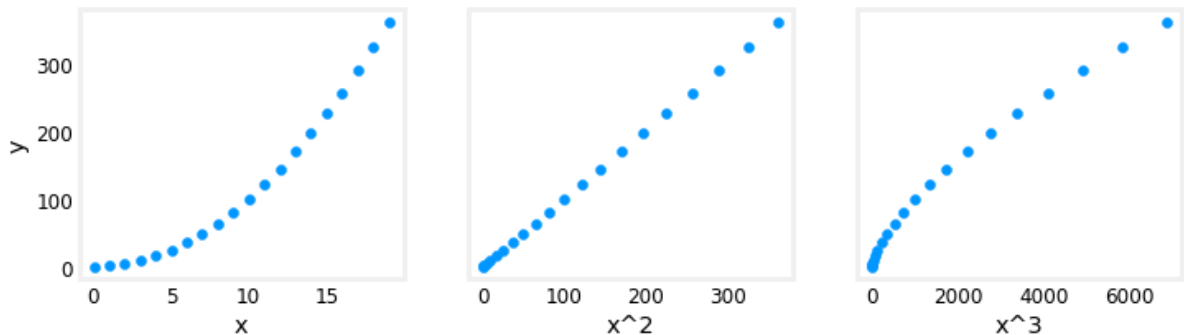
```
In [7]: # create target data
        x = np.arange(0, 20, 1)
        y = x**2
```

```
# engineer features .
X = np.c_[x, x**2, x**3]    #<-- added engineered feature
X_features = ['x','x^2','x^3']
```

In [8]:
```
fig,ax=plt.subplots(1, 3, figsize=(12, 3), sharey=True)
for i in range(len(ax)):
    ax[i].scatter(X[:,i],y)
    ax[i].set_xlabel(X_features[i])
ax[0].set_ylabel("y")
plt.show()
```



Above, it is clear that the $x^2$ feature mapped against the target value $y$ is linear. Linear regression can then easily generate a model using that feature.

## Scaling features

As described in the last lab, if the data set has features with significantly different scales, one should apply feature scaling to speed gradient descent. In the example above, there is $x$, $x^2$ and $x^3$ which will naturally have very different scales. Let's apply Z-score normalization to our example.

In [9]:
```
# create target data
x = np.arange(0,20,1)
X = np.c_[x, x**2, x**3]
print(f"Peak to Peak range by column in Raw        X:{np.ptp(X,axis=0)}")

# add mean_normalization
X = zscore_normalize_features(X)
print(f"Peak to Peak range by column in Normalized X:{np.ptp(X,axis=0)}")
```

```
Peak to Peak range by column in Raw        X:[  19   361 6859]
Peak to Peak range by column in Normalized X:[3.3  3.18 3.28]
```

Now we can try again with a more aggressive value of alpha:

In [10]:
```
x = np.arange(0,20,1)
y = x**2

X = np.c_[x, x**2, x**3]
X = zscore_normalize_features(X)

model_w, model_b = run_gradient_descent_feng(X, y, iterations=100000, alpha=1e-1)

plt.scatter(x, y, marker='x', c='r', label="Actual Value"); plt.title("Normalized x
plt.plot(x,X@model_w + model_b, label="Predicted Value"); plt.xlabel("x"); plt.ylab
```
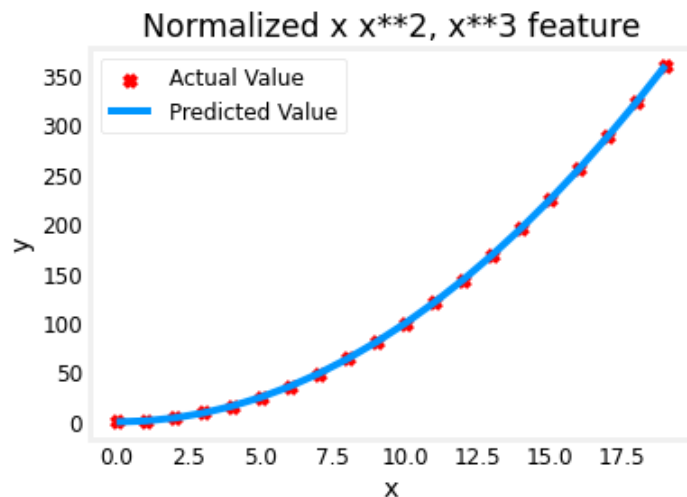
```
Iteration         0, Cost: 9.42147e+03
Iteration     10000, Cost: 3.90938e-01
Iteration     20000, Cost: 2.78389e-02
Iteration     30000, Cost: 1.98242e-03
Iteration     40000, Cost: 1.41169e-04
Iteration     50000, Cost: 1.00527e-05
Iteration     60000, Cost: 7.15855e-07
Iteration     70000, Cost: 5.09763e-08
Iteration     80000, Cost: 3.63004e-09
Iteration     90000, Cost: 2.58497e-10
w,b found by gradient descent: w: [5.27e-05 1.13e+02 8.43e-05], b: 123.5000
```

Normalized x x**2, x**3 feature



Feature scaling allows this to converge much faster.

Note again the values of $\mathbf{w}$. The $w_1$ term, which is the $x^2$ term is the most emphasized. Gradient descent has all but eliminated the $x^3$ term.

## Complex Functions

With feature engineering, even quite complex functions can be modeled:

```
In [11]: x = np.arange(0,20,1)
         y = np.cos(x/2)

         X = np.c_[x, x**2, x**3,x**4, x**5, x**6, x**7, x**8, x**9, x**10, x**11, x**12, x*
         X = zscore_normalize_features(X)

         model_w,model_b = run_gradient_descent_feng(X, y, iterations=1000000, alpha = 1e-1)

         plt.scatter(x, y, marker='x', c='r', label="Actual Value"); plt.title("Normalized x
         plt.plot(x,X@model_w + model_b, label="Predicted Value"); plt.xlabel("x"); plt.ylab
```
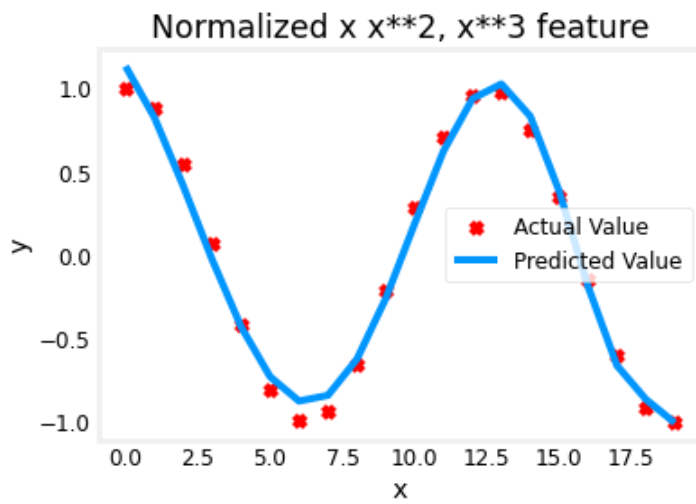```
Iteration         0, Cost: 2.20188e-01
Iteration    100000, Cost: 1.70074e-02
Iteration    200000, Cost: 1.27603e-02
Iteration    300000, Cost: 9.73032e-03
Iteration    400000, Cost: 7.56440e-03
Iteration    500000, Cost: 6.01412e-03
Iteration    600000, Cost: 4.90251e-03
Iteration    700000, Cost: 4.10351e-03
Iteration    800000, Cost: 3.52730e-03
Iteration    900000, Cost: 3.10989e-03
w,b found by gradient descent: w: [ -1.34 -10.    24.78   5.96 -12.49 -16.26  -9.51
  0.59    8.7   11.94
  9.27   0.79 -12.82], b: -0.0073
```

## Congratulations!

In this lab you:

- learned how linear regression can model complex, even highly non-linear functions using feature engineering
- recognized that it is important to apply feature scaling when doing feature engineering

In [11]:

# Optional Lab: Linear Regression using Scikit-Learn

There is an open-source, commercially usable machine learning toolkit called scikit-learn. This toolkit contains implementations of many of the algorithms that you will work with in this course.

## Goals

In this lab you will:

- Utilize scikit-learn to implement linear regression using Gradient Descent

## Tools

You will utilize functions from scikit-learn as well as matplotlib and NumPy.

```
In [1]: import numpy as np
        import matplotlib.pyplot as plt
        from sklearn.linear_model import SGDRegressor
        from sklearn.preprocessing import StandardScaler
        from lab_utils_multi import  load_house_data
        from lab_utils_common import dlc
        np.set_printoptions(precision=2)
        plt.style.use('./deeplearning.mplstyle')
```

# Gradient Descent

Scikit-learn has a gradient descent regression model sklearn.linear_model.SGDRegressor. Like your previous implementation of gradient descent, this model performs best with normalized inputs. sklearn.preprocessing.StandardScaler will perform z-score normalization as in a previous lab. Here it is referred to as 'standard score'.

## Load the data set

```
In [2]: X_train, y_train = load_house_data()
        X_features = ['size(sqft)','bedrooms','floors','age']
```

## Scale/normalize the training data

```
In [3]: scaler = StandardScaler()
        X_norm = scaler.fit_transform(X_train)
        print(f"Peak to Peak range by column in Raw        X:{np.ptp(X_train,axis=0)}")
        print(f"Peak to Peak range by column in Normalized X:{np.ptp(X_norm,axis=0)}")
```

```
Peak to Peak range by column in Raw        X:[2.41e+03 4.00e+00 1.00e+00 9.50e+01]
Peak to Peak range by column in Normalized X:[5.85 6.14 2.06 3.69]
```

## Create and fit the regression model

```
In [4]: sgdr = SGDRegressor(max_iter=1000)
        sgdr.fit(X_norm, y_train)
```

```
print(sgdr)
print(f"number of iterations completed: {sgdr.n_iter_}, number of weight updates: {
```

```
SGDRegressor()
number of iterations completed: 111, number of weight updates: 10990.0
```

## View parameters

Note, the parameters are associated with the *normalized* input data. The fit parameters are very close to those found in the previous lab with this data.

In [5]:
```
b_norm = sgdr.intercept_
w_norm = sgdr.coef_
print(f"model parameters:                  w: {w_norm}, b:{b_norm}")
print( "model parameters from previous lab: w: [110.56 -21.27 -32.71 -37.97], b: 36
```

```
model parameters:                  w: [109.95 -20.97 -32.35 -38.07], b:[363.15]
model parameters from previous lab: w: [110.56 -21.27 -32.71 -37.97], b: 363.16
```

## Make predictions

Predict the targets of the training data. Use both the `predict` routine and compute using $w$ and $b$.

In [6]:
```
# make a prediction using sgdr.predict()
y_pred_sgd = sgdr.predict(X_norm)
# make a prediction using w,b.
y_pred = np.dot(X_norm, w_norm) + b_norm
print(f"prediction using np.dot() and sgdr.predict match: {(y_pred == y_pred_sgd).a

print(f"Prediction on training set:\n{y_pred[:4]}" )
print(f"Target values \n{y_train[:4]}")
```

```
prediction using np.dot() and sgdr.predict match: True
Prediction on training set:
[295.17 485.84 389.62 492.  ]
Target values
[300.  509.8 394.  540. ]
```
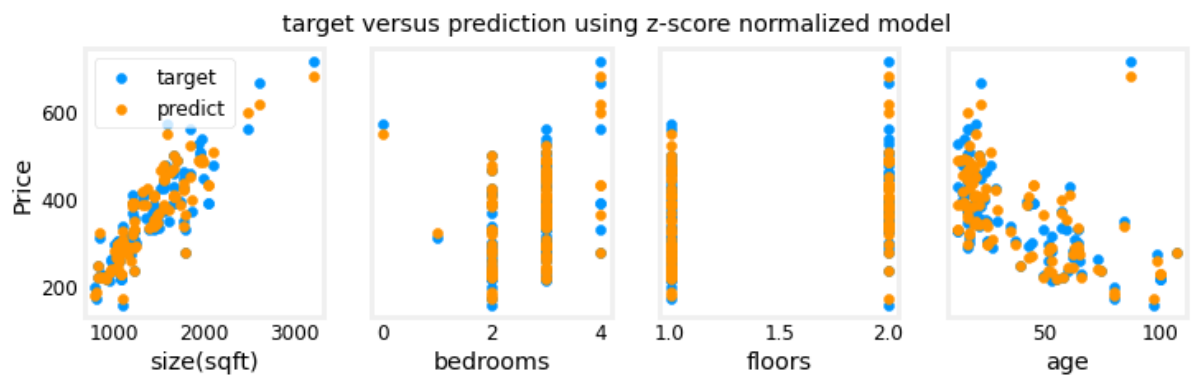
## Plot Results

Let's plot the predictions versus the target values.

In [7]:
```
# plot predictions and targets vs original features
fig,ax=plt.subplots(1,4,figsize=(12,3),sharey=True)
for i in range(len(ax)):
    ax[i].scatter(X_train[:,i],y_train, label = 'target')
    ax[i].set_xlabel(X_features[i])
    ax[i].scatter(X_train[:,i],y_pred,color=dlc["dlorange"], label = 'predict')
ax[0].set_ylabel("Price"); ax[0].legend();
fig.suptitle("target versus prediction using z-score normalized model")
plt.show()
```

target versus prediction using z-score normalized model

## Congratulations!

In this lab you:

- utilized an open-source machine learning toolkit, scikit-learn
- implemented linear regression using gradient descent and feature normalization from that toolkit

In [7]: