Make sure you fill in any place that says `YOUR CODE HERE`.

+ Code      + Text

# ▾ Homework 7

*This* is a Python Notebook homework. It consists of various types of cells:

- Text: you can read them :-)
- Code: you should run them, as they may set up the problems that you are asked to solve.
- **Solution:** These are cells where you should enter a solution. You will see a marker in these cells that indicates where your work should be inserted.

```
# YOUR CODE HERE
```

- Test: These cells contains some tests, and are worth some points. You should run the cells as a way to debug your code, and to see if you understood the question, and whether the output of your code is produced in the correct format. The notebook contains both the tests you see, and some secret ones that you cannot see. This prevents you from using the simple trick of hard-coding the desired output.

## Questions

There are several groups of questions in this notebook:

- Implementation of `Expr.compute`
- Implementation of `Multiply`, `Minus`, `Divide`, `Power`, `Negative`
- Implementation of `Expr.compute_gradient`
- Implementation of `fit`
- Implementation of fitting for various expressions.

There are other pieces of text called "exercises", but you only have to do those that are explicitly marked with a place in the code for you to write the answer.

## Working on Your Notebook

To work on your notebook:

- Click on *File > Save a copy in Drive* : this will create a copy of this file in your Google Drive; you will find the notebook in your *Colab Notebooks* folder.
- Work on that notebook.

## Submitting Your Notebook

Submit your work as follows:

- Download the notebook from Colab, clicking on "File > Download .ipynb".
- Upload the resulting file to this Google form.
- **Deadline: see home page**

You can submit multiple times, and the last submittion before the deadline will be used to assign you a grade.

## ML in a nutshell

Optimization, and machine learning, are intimately connected. At a very coarse level, ML works as follows.

First, you come up somehow with a very complicated model $\vec{y} = M(\vec{x}, \vec{\theta})$, which computes an output $\vec{y}$ as a function of an input $\vec{x}$ and of a vector of parameters $\vec{\theta}$. In general, $\vec{x}$, $\vec{y}$, and $\vec{\theta}$ are vectors, as the model has multiple inputs, multiple outputs, and several parameters. The model $M$ needs to be complicated, because only complicated models can represent complicated phenomena; for instance, $M$ can be a multi-layer neural net with parameters $\vec{\theta} = [\theta_1, \ldots, \theta_k]$, where $k$ is the number of parameters of the model.

Second, you come up with a notion of *loss* $L$, that is, how badly the model is doing. For instance, if you have a list of inputs $\vec{x}_1, \ldots, \vec{x}_n$, and a set of desired outputs $\vec{y}_1, \ldots, \vec{y}_m$, you can use as loss:

$$L(\vec{\theta}) = \sum_{i=1}^{n} \|\vec{y}_i - M(\vec{x}_i, \vec{\theta})\| \ .$$

Here, we wrote $L(\vec{\theta})$ because, once the inputs $\vec{x}_1, \ldots, \vec{x}_n$ and the desired outputs $\vec{y}_1, \ldots, \vec{y}_n$ are chosen, the loss $L$ depends on $\vec{\theta}$.

Once the loss is chosen, you decrease it, by computing its *gradient* with respect to $\vec{\theta}$. Remembering that $\vec{\theta} = [\theta_1, \ldots, \theta_k]$,

$$\nabla_{\vec{\theta}} L = \left[ \frac{\partial L}{\partial \theta_1}, \ldots, \frac{\partial L}{\partial \theta_k} \right] \ .$$

The gradient is a vector that indicates how to tweak $\vec{\theta}$ to decrease the loss. You then choose a small *step size* $\delta$, and you update $\vec{\theta}$ via $\vec{\theta} := \vec{\theta} - \delta \nabla_{\vec{\theta}} L$. This makes the loss a little bit smaller, and the model a little bit better. If you repeat this step many times, the model will hopefully get (a good bit) better.

## Autogradient

The key to *pleasant* ML is to focus on building the model $M$ in a way that is sufficiently expressive, and on choosing a loss $L$ that is helpful in guiding the optimization. The computation of the gradient is done automatically for you. This capability, called *autogradient*, is implemented in ML frameworks such as [Tensorflow](#), [Keras](#), and [PyTorch](#).

It is possible to use these advanced ML libraries without ever knowing what is under the hood, and how autogradient works. Here, we will insted dive in, and implement autogradient.

Building a model $M$ corresponds to building an expression with inputs $\vec{x}, \vec{\theta}$. We will provide a representaton for expressions that enables both the calculation of the expression value, and the differentiation with respect to any of the inputs. This will enable us to implement autogradient. On the basis of this, we will be able to implement a simple ML framework.

We say we, but we mean you. *You* will implement it; we will just provide guidance.

## ▾ Expressions with autogradient

Our main task will be to implement a class `Expr` that represents expressions with autogradient.

### Implementing expressions

We will have `Expr` be the abstract class of a generic expression, and `Plus`, `Multiply`, and so on, be derived classes representing expression with given top-level operators. The constructor takes the children node. The code for the constructor, and the code to create addition expressions, is as follows.

```
1 class Expr(object):
2
3     def __init__(self, *args):
4         """Initializes an expression node, with a given list of children
5         expressions."""
6         self.children = args
7         self.value = None # The value of the expression.
8         self.values = None # The values of the child expressions.
9
10    def __add__(self, other):
11        """Constructs the sum of two expressions."""
12        return Plus(self, other)
```

The code for the `Plus` class, initially, is empty; no `Expr` methods are over-ridden.

```
1 class Plus(Expr):
2     """An addition expression."""
3
```

of values `self.values` of the children.

Let's implement the `compute` method for an expression. This method will:

1. Loop over the children, and computes the list `self.values` of children values as follows:

   ○ If the child is an expression (an instance of `Expr`, obtain its value by calling `compute` on it.

   ○ If the child is not an instance of `Expr`, then the child must be a number, and we can use its value directly.

2. Call the method `op` of the expression, to compute `self.value` from `self.values`.
3. return `self.value`.

We will let you implement the `compute` method. Hint: it takes just a couple of lines of code.

```
 1 class Expr(object):
 2
 3     def __init__(self, *args):
 4         """Initializes an expression node, with a given list of children
 5         expressions."""
 6         self.children = args
 7         self.value = None # The value of the expression.
 8         self.values = None # The values of the child expressions.
 9         self.gradient = 0 # The value of the gradient.
10
11     def op(self):
12         """This operator must be implemented in subclasses; it should
13         compute self.value from self.values, thus implementing the
14         operator at the expression node."""
15         raise NotImplementedError()
16
17     def compute(self):
18         """This method computes the value of the expression.
19         It first computes the value of the children expressions,
20         and then uses self.op to compute the value of the expression."""
21         self.value = None ### INSERT YOUR SOLUTION HERE
22         return self.value
23
24     def __repr__(self):
25         return ("%s:%r %r (g: %r)" % (
26             self.__class__.__name__, self.children, self.value, self.gradient))
27
28     # Expression constructors
29
30     def __add__(self, other):
31         return Plus(self, other)
32
33     def __radd__(self, other):
34         return Plus(self, other)
```

```
4      pass
```

To construct expressions, we need one more thing. So far, if we write things like `2 + 3`, Python will just consider these as expressions involving numbers, and compute their value. To write *symbolic* expressions, we need symbols, or variables. A variable is a type of expression that just contains a value as child, and that has an `assign` method to assign a value to the variable. The `assign` method can be used to modify the variable's content (without `assign`, our variables would be constants!).

```
1 class V(Expr):
2     """This class represents a variable.  The derivative rule corresponds
3     to d/dx x = 1, but note that it will not be called, since the children
4     of a variable are just numbers."""
5
6     def assign(self, v):
7         """Assigns a value to the variable."""
8         self.children = [v]
```

This suffices for creating expressions. Let's create one.

```
1 e = V(3) + 4
2 e
```

⤷   <__main__.Plus at 0x7f7d9049b080>

```
1 # Let us ensure that nose is installed.
2 try:
3     from nose.tools import assert_equal, assert_true
4     from nose.tools import assert_false, assert_almost_equal
5 except:
6     !pip install nose
7     from nose.tools import assert_equal, assert_true
8     from nose.tools import assert_false, assert_almost_equal
```

⤷   Collecting nose
      Downloading https://files.pythonhosted.org/packages/15/d8/dd071918c040f50fa1cf:
          |████████████████████████████████| 163kB 4.7MB/s
      Installing collected packages: nose
      Successfully installed nose-1.3.7

## ▾ Computing the value of expressions

We now have our first expression. To compute the expression value, we endow each expression with a method `op`, whose task is to compute the value `self.value` of the expression from the list

```
35
36    def __sub__(self, other):
37        return Minus(self, other)
38
39    def __rsub__(self, other):
40        return Minus(other, self)
41
42    def __mul__(self, other):
43        return Multiply(self, other)
44
45    def __rmul__(self, other):
46        return Multiply(other, self)
47
48    def __truediv__(self, other):
49        return Divide(self, other)
50
51    def __rtruediv__(self, other):
52        return Divide(other, self)
53
54    def __pow__(self, other):
55        return Power(self, other)
56
57    def __rpow__(self, other):
58        return Power(other, self)
59
60    def __neg__(self):
61        return Negative(self)
```

Let us give `op` for `Plus`, `Multiply`, and for variables via `v`, so you can see how it works.

```
 1 class V(Expr):
 2     """This class represents a variable."""
 3
 4     def assign(self, v):
 5         """Assigns a value to the variable.  Used to fit a model, so we
 6         can assign the various input values to the variable."""
 7         self.children = [v]
 8
 9     def op(self):
10         self.value = self.values[0]
11
12     def __repr__(self):
13         return "Variable: " + str(self.children[0])
14
15
16 class Plus(Expr):
17
18     def op(self):
19         self.value = self.values[0] + self.values[1]
20
```

```
21
22 class Multiply(Expr):
23
24     def op(self):
25         self.value = self.values[0] * self.values[1]
```

Here you can write your implementation of the `compute` method.

```
1 ### Exercise: Implementation of `compute` method
2
3 def expr_compute(self):
4     """This method computes the value of the expression.
5     It first computes the value of the children expressions,
6     and then uses self.op to compute the value of the expression."""
7     # YOUR CODE HERE
8     self.values = []
9     for c in self.children:
10        if isinstance(c, Expr):
11            self.values.append(c.compute())
12        else:
13            self.values.append(c)
14    self.op()
15    return self.value
16
17 Expr.compute = expr_compute
```

```
1 ### Tests for compute
2
3 from nose.tools import assert_equal, assert_true, assert_false
4
5 # First, an expression consisting only of one variable.
6 e = V(3)
7 assert_equal(e.compute(), 3)
8 assert_equal(e.value, 3)
9
10 # Then, an expression involving plus.
11 e = V(3) + 4
12 assert_equal(e.compute(), 7)
13 assert_equal(e.value, 7)
14
15 # And finally, a more complex expression.
16 e = (V(3) + 4) + V(2)
17 assert_equal(e.compute(), 9)
18 assert_equal(e.value, 9)
19
```

We will have you implement also multiplication.

```
1 ### Exercise: Implement `Multiply`
2
3 class Multiply(Expr):
4     """A multiplication expression."""
5
6     def op(self):
7         # YOUR CODE HERE
8         self.value = self.values[0]*self.values[1]
9
```

```
1 ### Tests for `Multiply`
2
3 e = V(2) * 3
4 assert_equal(e.compute(), 6)
5
6 e = (V(2) + 3) * V(4)
7 assert_equal(e.compute(), 20)
8
```

## ▼ Implementing autogradient

The next step consists in implementing autogradient. Consider an expression $e = E(x_0, \ldots, x_n)$, computed as function of its children expressions $x_0, \ldots, x_n$.

The goal of the autogradient computation is to accumulate, in each node of the expression, the gradient of the loss with respect to the node's value. For instance, if the gradient is $2$, we know that if we increase the value of the expression by $\Delta$, then the value of the loss is increased by $2\Delta$. We accumulate the gradient in the field `self.gradient` of the expression.

We say *accumulate* the gradient, because we don't really do:

```
self.gradient = ...
```

Rather, we have a method `e.zero_gradient()` that sets all gradients to 0, and we then *add* the gradient to this initial value of 0:

```
self.gradient += ...
```

We will explain later in detail why we do so; for the moment, just accept it.

### Computaton of the gradient

In the computation of the autogradient, the expression will receive as input the value $\partial L/\partial e$, where $L$ is the loss, and $e$ the value of the expression. The quantity $\partial L/\partial e$ is the gradient of the loss with respect to the expression value.

With this input, the method `compute_gradient` of Expr must do the following:

- It must *add $\partial L/\partial e$* to the gradient `self.gradient` of the expression.
- It must compute for each child $x_i$ the partial derivative $\partial e/\partial x_i$, via a call to the method `derivate`. The method `derivate` is implemented not for `Expr`, but for each specific operator, such as `Plus`, `Multiply`, etc: each operator knows how to compute the derivative with respect to its arguments.
- It must propagate to each child $x_i$ the gradient $\frac{\partial L}{\partial e} \cdot \frac{\partial e}{\partial x_i}$, by calling the method `compute_gradient` of the child with argument $\frac{\partial L}{\partial e} \cdot \frac{\partial e}{\partial x_i}$.

```
 1 def expr_derivate(self):
 2     """This method computes the derivative of the operator at the expression
 3     node.  It needs to be implemented in derived classes, such as Plus,
 4     Multiply, etc."""
 5     raise NotImplementedError()
 6
 7 Expr.derivate = expr_derivate
 8
 9 def expr_zero_gradient(self):
10     """Sets the gradient to 0, recursively for this expression
11     and all its children."""
12     self.gradient = 0
13     for e in self.children:
14         if isinstance(e, Expr):
15             e.zero_gradient()
16
17 Expr.zero_gradient = expr_zero_gradient
18
19 def expr_compute_gradient(self, de_loss_over_de_e=1):
20     """Computes the gradient.
21     de_loss_over_de_e is the gradient of the output.
22     de_loss_over_de_e will be added to the gradient, and then
23     we call for each child the method compute_gradient,
24     with argument de_loss_over_de_e * d expression / d child.
25     The value d expression / d child is computed by self.derivate. """
26     pass ### PLACEHOLDER FOR YOUR SOLUTION.
27
28 Expr.compute_gradient = expr_compute_gradient
```

Let us endow our operators `V`, `Plus`, `Multiply` with the `derivate` method, so you can see how it works in practice.

```
1 class V(Expr):
2     """This class represents a variable.  The derivative rule corresponds
3     to d/dx x = 1, but note that it will not be called, since the children
4     of a variable are just numbers."""
```

```
 5
 6    def assign(self, v):
 7        """Assigns a value to the variable.  Used to fit a model, so we
 8        can assign the various input values to the variable."""
 9        self.children = [v]
10
11    def op(self):
12        self.value = self.values[0]
13
14    def derivate(self):
15        return [1.] # This is not really used.
16
17
18 class Plus(Expr):
19    """An addition expression.  The derivative rule corresponds to
20    d/dx (x+y) = 1, d/dy (x+y) = 1"""
21
22    def op(self):
23        self.value = self.values[0] + self.values[1]
24
25    def derivate(self):
26        return [1., 1.]
27
28
29 class Multiply(Expr):
30    """A multiplication expression. The derivative rule corresponds to
31    d/dx (xy) = y, d/dy(xy) = x"""
32
33    def op(self):
34        self.value = self.values[0] * self.values[1]
35
36    def derivate(self):
37        return [self.values[1], self.values[0]]
```

Let us comment on some subtle points, before you get to work at implementing `compute_gradient`.

**`zero_gradient`** : First, notice how in the implementation of `zero_gradient`, when we loop over the children, we check whether each children is an `Expr` via isinstance(e, Expr). In general, we have to remember that children can be either `Expr`, or simply numbers, and of course numbers do not have methods such as `zero_gradient` or `compute_gradient` implemented for them.

**`derivate`** : Second, notice how `derivate` is not implemented in `Expr` directly, but rather, only in the derived classes such as `Plus`. The derivative of the expression with respect to its arguments depends on which function it is, obviously.

For `Plus`, we have $e = x_0 + x_1$, and so:

$$\frac{\partial e}{\partial x_0} = 1 \qquad \frac{\partial e}{\partial x_1} = 1 \ ,$$

because $d(x + y)/dx = 1$. Hence, the `derivate` method of `Plus` returns

$$\left[ \frac{\partial e}{\partial x_0}, \ \frac{\partial e}{\partial x_1} \right] \ = \ [1, 1] \ .$$

For `Multiply`, we have $e = x_0 \cdot x_1$, and so:

$$\frac{\partial e}{\partial x_0} = x_1 \qquad \frac{\partial e}{\partial x_1} = x_0 \ ,$$

because $d(xy)/dx = y$. Hence, the `derivate` method of `Plus` returns

$$\left[ \frac{\partial e}{\partial x_0}, \ \frac{\partial e}{\partial x_1} \right] \ = \ [x_1, x_0] \ .$$

**Calling `compute` before `compute_gradient`:** Lastly, a very important point: when calling `compute_gradient`, we will assume that `compute` has *already* been called. In this way, the value of the expression, and its children, are available for the computation of the gradient. Note how in

With these clarifications, we ask you to implement the `compute_gradient` method, which again must:

- *add $\partial L/\partial e$* to the gradient `self.gradient` of the expression;
- compute $\frac{\partial e}{\partial x_i}$ for each child $x_i$ by calling the method `derivate` of itself;
- propagate to each child $x_i$ the gradient $\frac{\partial L}{\partial e} \cdot \frac{\partial e}{\partial x_i}$, by calling the method `compute_gradient` of the child with argument $\frac{\partial L}{\partial e} \cdot \frac{\partial e}{\partial x_i}$.

```
1 ### Exercise: Implementation of `compute_gradient`
2
3 def expr_compute_gradient(self, de_loss_over_de_e=1):
4     """Computes the gradient.
5     de_loss_over_de_e is the gradient of the output.
6     de_loss_over_de_e will be added to the gradient, and then
7     we call for each child the method compute_gradient,
8     with argument de_loss_over_de_e * d expression / d child.
9     The value d expression / d child is computed by self.derivate. """
10     # YOUR CODE HERE
11     self.gradient+=de_loss_over_de_e
12     for i in range(len(self.children)):
13         if isinstance(self.children[i],Expr):
14             d = self.derivate()
15             self.children[i].compute_gradient(de_loss_over_de_e *d[i] )
16
17
18
19
20
```

```
21
22
23
24
25
26
27 Expr.compute_gradient = expr_compute_gradient
```

```
 1 ### Tests for `compute_gradient`
 2
 3 # First, the gradient of a sum.
 4 vx = V(3)
 5 vz = V(4)
 6 y = vx + vz
 7 assert_equal(y.compute(), 7)
 8 y.zero_gradient()
 9 y.compute_gradient()
10 assert_equal(vx.gradient, 1.)
11
12 # Second, the gradient of a product.
13 vx = V(3)
14 vz = V(4)
15 y = vx * vz
16 assert_equal(y.compute(), 12)
17 y.zero_gradient()
18 y.compute_gradient()
19 assert_equal(vx.gradient, 4)
20 assert_equal(vz.gradient, 3)
21
22 # Finally, the gradient of the product of sums.
23
24 vx = V(1)
25 vw = V(3)
26 vz = V(4)
27 y = (vx + vw) * (vz + 3)
28 assert_equal(y.compute(), 28)
29 y.zero_gradient()
30 y.compute_gradient()
31 assert_equal(vx.gradient, 7)
32 assert_equal(vz.gradient, 4)
33
```

## ▾ Why do we accumulate gradients?

We are now in the position of answering the question of why we accumulate gradients. There are
two reasons.

## Multiple variable occurrence

The most important reason why we need to *add* to the gradient of each node is that nodes, and in particular, variable nodes, can occur in multiple places in an expression tree. To compute the total influence of the variable on the expression, we need to *sum* the influence of each occurrence. Let's see this with a simple example. Consider the expression $y = x \cdot x$. We can code it as follows:

```
1 vx = V(2.) # Creates a variable vx and initializes it to 2.
2 y = vx * vx
```

For $y = x^2$, we have $dy/dx = 2x = 4$, given that $x = 2$. How is this reflected in our code?

Our code considers separately the left and right occurrences of `vx` in the expression; let us denote them with $vx_l$ and $vx_r$. The expression can be written as $y = vx_l \cdot vx_r$, and we have that $\partial y/\partial vx_l = vx_r = 2$, as $vx_r = 2$. Similarly, $\partial y/\partial vx_r = 2$. These two gradients are added to `vx.gradient` by the method `compute_gradient`, and we get that the total gradient is 4, as desired.

```
1 y.compute() # We have to call compute() before compute_gradient()
2 y.zero_gradient()
3 y.compute_gradient()
4 print("gradient of vx:", vx.gradient)
5 assert_equal(vx.gradient, 4)
```

```
gradient of vx: 4.0
```

## Multiple data to fit

The other reason why we need to tally up the gradient is that in general, we need to fit a function to more than one data point. Assume that we are given a set of inputs $x_1, x_2, \ldots, x_n$ and desired outputs $y_1, y_2, \ldots, y_n$. Our goal is to approximate the desired ouputs via an expression $e(x, \theta)$ of the input, according to some parameters $\theta$. Our goal is to choose the parameters $\theta$ to minimize the sum of the square errors for the points:

$$L_{tot}(\theta) = \sum_{i=1}^{n} L_i(\theta) \, ,$$

where $L_i(\theta) = (e(x_i, \theta) - y_i)^2$ is the loss for a single data point. The gradient $L_{tot}(\theta)$ with respect to $\theta$ can be computed by adding up the gradients for the individual points:

$$\frac{\partial}{\partial \theta} L_{tot}(\theta) = \sum_{i=1}^{n} \frac{\partial}{\partial \theta} L_i(\theta) \, .$$

To translate this into code, we will build an expression $e(x, \theta)$ involving the input $x$ and the parameters $\theta$, and an expression

$$L = (e(x, \theta) - y)^2$$

for the loss, involving $x$, $y$ and $\theta$. We will then zero all gradients via `zero_gradient`. Once this is done, we compute the loss $L$ for each point, and then the gradient $\partial L/\partial \theta$ via a call to `compute_gradient`. The gradients for all the points will be added, yielding the gradient for

## ▾ Rounding up the implementation

Now that we have implemented autogradient, as well as the operators `Plus` and `Multiply`, it is time to implement the remaining operators:

- `Minus`
- `Divide` (no need to worry about division by zero)
- `Power`, representing exponentiation (the `**` operator of Python)
- and the unary minus `Negative`.

```
1 ### Exercise: Implementation of `Minus`, `Divide`, `Power`, and `Negative`
2
3 import math
4
5 class Minus(Expr):
6     """Operator for x - y"""
7
8     def op(self):
9         # YOUR CODE HERE
10        self.value = self.values[0]-self.values[1]
11    def derivate(self):
12        # YOUR CODE HERE
13        return [1, -1]
14 class Divide(Expr):
15     """Operator for x / y"""
16
17     def op(self):
18        # YOUR CODE HERE
19        self.value = self.values[0]/self.values[1]
20    def derivate(self):
21        # YOUR CODE HERE
22        return [self.values[1]/(self.values[1]**2), -self.values[0]/(self.values[1]
23 class Power(Expr):
24     """Operator for x ** y"""
25
26     def op(self):
27        # YOUR CODE HERE
28        self.value = self.values[0]**self.values[1]
29    def derivate(self):
30        # YOUR CODE HERE
```

```
30          # YOUR CODE HERE
31          return [self.values[1]*(self.values[0]**(self.values[1]-1)) , math.log(self
32 class Negative(Expr):
33     """Operator for -x"""
34
35     def op(self):
36         # YOUR CODE HERE
37         self.value = -1* self.values[0]
38     def derivate(self):
39         # YOUR CODE HERE
40         return [-1.]
```

Here are some tests.

```
 1 ### Tests for `Minus`
 2
 3 # Minus.
 4 vx = V(3)
 5 vy = V(2)
 6 e = vx - vy
 7 assert_equal(e.compute(), 1.)
 8 e.zero_gradient()
 9 e.compute_gradient()
10 assert_equal(vx.gradient, 1)
11 assert_equal(vy.gradient, -1)
12
```

```
 1 ### Tests for `Divide`
 2
 3 from nose.tools import assert_almost_equal
 4
 5 # Divide.
 6 vx = V(6)
 7 vy = V(2)
 8 e = vx / vy
 9 assert_equal(e.compute(), 3.)
10 e.zero_gradient()
11 e.compute_gradient()
12 assert_equal(vx.gradient, 0.5)
13 assert_equal(vy.gradient, -1.5)
14
```

```
 1 ### Tests for `Power`
 2
 3 from nose.tools import assert_almost_equal
 4
 5 # Power.
 6 vx = V(2)
 7 vy = V(3)
```

```
 8 e = vx ** vy
 9 assert_equal(e.compute(), 8.)
10 e.zero_gradient()
11 e.compute_gradient()
12 assert_equal(vx.gradient, 12.)
13 assert_almost_equal(vy.gradient, math.log(2.) * 8., places=4)
14
```

```
 1 ### Tests for `Negative`
 2
 3 from nose.tools import assert_almost_equal
 4
 5 # Negative
 6 vx = V(6)
 7 e = - vx
 8 assert_equal(e.compute(), -6.)
 9 e.zero_gradient()
10 e.compute_gradient()
11 assert_equal(vx.gradient, -1.)
12
```

## ▾ Optimization

Let us use our ML framework to fit a parabola to a given set of points. Here is our set of points:

```
 1 points = [
 2     (-2, 2.7),
 3     (-1, 3),
 4     (0, 1.3),
 5     (1, 2.4),
 6     (3, 5.5),
 7     (4, 6.2),
 8     (5, 9.1),
 9 ]
```

Let us display these points.

```
 1 import matplotlib
 2 import matplotlib.pyplot as plt
 3
 4 matplotlib.rcParams['figure.figsize'] = (8.0, 3.)
 5 params = {'legend.fontsize': 'large',
 6           'axes.labelsize': 'large',
 7           'axes.titlesize':'large',
 8           'xtick.labelsize':'large',
 9           'ytick.labelsize':'large'}
10 matplotlib.rcParams.update(params)
```
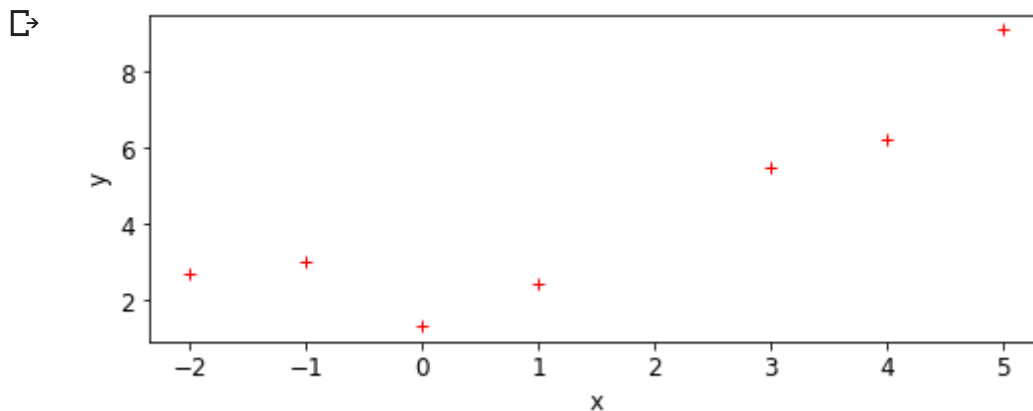
```
11
12 def plot_points(points):
13     fig, ax = plt.subplots()
14     xs, ys = zip(*points)
15     ax.plot(xs, ys, 'r+')
16     plt.xlabel('x')
17     plt.ylabel('y')
18     plt.show()
```

```
1 plot_points(points)
```



To fit a parabola to these points, we will build an `Expr` that represents the equation $\hat{y} = ax^2 + bx + c$, where $\hat{y}$ is the value of $y$ predicted by our parabola. If $\hat{y}$ is the predicted value, and $y$ is the observed value, to obtain a better prediction of the observations, we minimize the loss $L = (\hat{y} - y)^2$, that is, the square prediction error. Written out in detail, our loss is:

$$L = (y - \hat{y})^2 = (y - (ax^2 + bx + c))^2 .$$

Here, $a$, $b$, $c$ are parameters that we need to tune to minimize the loss, and obtain a good fit between the parabola and the points. This tuning, or training, is done by repeating the following process many times:

- Zero the gradient
- For each point:
  - Set the values of x, y to the value of the point.
  - Compute the expression giving the loss.
  - Backpropagate. This computes all gradients with respect to the loss, and in particular, the gradients of the coefficients $a$, $b$, $c$.
- Update the coefficients $a$, $b$, $c$ by taking a small step in the direction of the negative gradient (negative, so that the loss decreases).

```
1 va = V(0.)
2 vb = V(0.)
```

```
3 vc = V(0.)
4 vx = V(0.)
5 vy = V(0.)
6
7 oy = va * vx * vx + vb * vx + vc
8
9 loss = (vy - oy) * (vy - oy)
```

Below, implement the "for each point" part of the above informal description. Hint: this takes about 4-5 lines of code.

```
 1 def fit(loss, points, params, delta=0.0001, num_iterations=4000):
 2
 3     for iteration_idx in range(num_iterations):
 4         loss.zero_gradient()
 5         total_loss = 0.
 6         for x, y in points:
 7             ### You need to implement here the computaton of the
 8             ### loss gradient for the point (x, y).
 9             total_loss += loss.value
10         if (iteration_idx + 1) % 100 == 0:
11             print("Loss:", total_loss)
12         for vv in params:
13             vv.assign(vv.value - delta * vv.gradient)
14     return total_loss
```

```
 1 ### Exercise: Implementation of `fit`
 2
 3 def fit(loss, points, params, delta=0.0001, num_iterations=4000):
 4
 5     for iteration_idx in range(num_iterations):
 6         loss.zero_gradient()
 7         total_loss = 0.
 8         for x, y in points:
 9             # YOUR CODE HERE
10             vx.assign(x)
11             vy.assign(y)
12             loss.value = loss.compute()
13             loss.compute_gradient()
14             total_loss += loss.value
15         if (iteration_idx + 1) % 100 == 0:
16             print("Loss:", total_loss)
17         for vv in params:
18             vv.assign(vv.value - delta * vv.gradient)
19     return total_loss
```

Let's train the coefficients `va`, `vb`, `vc`:

```
1
```

```
1 from nose.tools import assert_less
2
3 lv = fit(loss, points, [va, vb, vc])
4 assert_less(lv, 2.5)
```

```
Loss: 15.691480263172831
Loss: 13.628854145973717
Loss: 11.95710461470721
Loss: 10.577288435238483
Loss: 9.421709915417068
Loss: 8.442894724380599
Loss: 7.606627717852742
Loss: 6.887533665385938
Loss: 6.266252079473012
Loss: 5.727613949776414
Loss: 5.259450432876111
Loss: 4.851802101057622
Loss: 4.4963837663728965
Loss: 4.18621382152429
Loss: 3.9153507183966223
Loss: 3.6787002636071753
Loss: 3.471870598609928
Loss: 3.2910600093895632
Loss: 3.13296792199636
Loss: 2.9947227347465053
Loss: 2.8738222326925422
Loss: 2.768083672161589
Loss: 2.6756014919475786
Loss: 2.594711177586694
Loss: 2.523958185206276
Loss: 2.462071090153853
Loss: 2.4079383059994615
Loss: 2.360587848696134
Loss: 2.3191697158813342
Loss: 2.28294052348259
Loss: 2.251250098021605
Loss: 2.2235297678961032
Loss: 2.199282133501638
Loss: 2.1780721263873506
Loss: 2.1595191931327227
Loss: 2.1432904612844434
Loss: 2.1290947632296184
Loss: 2.1166774098449848
Loss: 2.105815619569577
Loss: 2.096314520528721
```

Let's display the parameter values after the training:

```
1 print("a:", va.value, "b:", vb.value, "c:", vc.value)
```

```
a: 0.2676897030023671 b: 0.09446139849938508 c: 1.9725828646266883
```
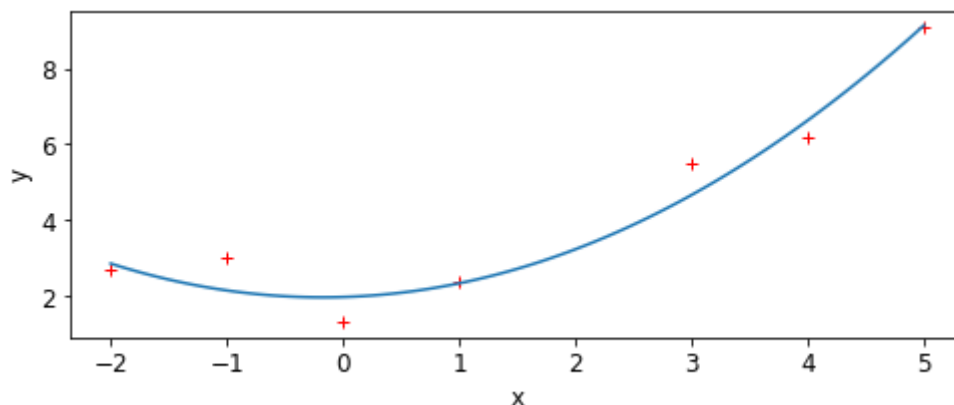
Let's display the points, along with the fitted parabola.

```
1 import numpy as np
2
3 def plot_points_and_y(points, vx, oy):
4     fig, ax = plt.subplots()
5     xs, ys = zip(*points)
6     ax.plot(xs, ys, 'r+')
7     x_min, x_max = np.min(xs), np.max(xs)
8     step = (x_max - x_min) / 100
9     x_list = list(np.arange(x_min, x_max + step, step))
10    y_list = []
11    for x in x_list:
12        vx.assign(x)
13        oy.compute()
14        y_list.append(oy.value)
15    ax.plot(x_list, y_list)
16    plt.xlabel('x')
17    plt.ylabel('y')
18    plt.show()
```

```
1 plot_points_and_y(points, vx, oy)
```



This looks like a good fit!

Note that if we chose too large a learning step, we would not converge to a solution. A large step causes the parameter values to zoom all over the place, possibly missing by large amounts the (local) minima where you want to converge. In the limit where the step size goes to 0, and the number of steps to infinity, you are guaranteed (if the function is differentiable, and some other hypotheses) converge to the minimum; the problem is that it would take infinitely long. You will learn in a more in-depth ML class how to tune the step size.

```
1 # Let us reinitialize the variables.
2 va.assign(0)
```

```
3 vb.assign(0)
4 vc.assign(0)
5 # ... and let's use a big step size.
6 fit(loss, points, [va, vb, vc], delta=0.01, num_iterations=1000)
```

```
Loss: 7.445536528760376e+257
Loss: inf
Loss: nan
Loss: nan
Loss: nan
Loss: nan
Loss: nan
Loss: nan
Loss: nan
Loss: nan
Loss: nan
nan
```

A step size of 0.01 was enough to take us to infinity and beyond.

Let us now show you how to fit a simple linear regression: $y = ax + b$, so $L = (y - (ax + b))^2$.

```
 1 # Parameters
 2 # Sometimes you have to be careful about initial values.
 3 va = V(1.)
 4 vb = V(1.)
 5
 6 # x and y
 7 vx = V(0.)
 8 vy = V(0.)
 9
10 # Predicted y
11 oy = va * vx + vb
12
13 # Loss
14 loss = (vy - oy) * (vy - oy)
```

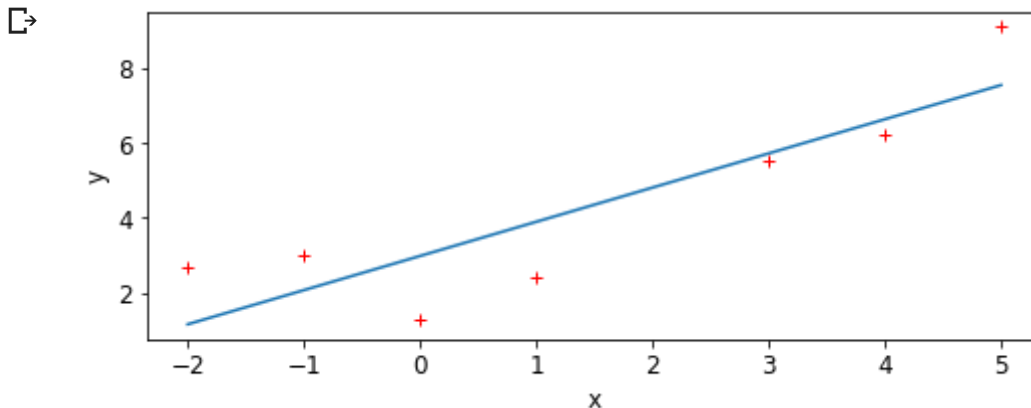```
1 fit(loss, points, [va, vb])
```

```
Loss: 28.04911931725314
Loss: 24.56807140584221
Loss: 22.04737547646302
Loss: 20.018461554483544
Loss: 18.363064879989466
Loss: 17.010221159926804
Loss: 15.904419410376471
Loss: 15.000526654509262
Loss: 14.261674191105309
Loss: 13.657727388420048
Loss: 13.164054060381849
Loss: 12.760519585409948
Loss: 12.43066568540981
Loss: 12.161039167477448
Loss: 11.940643231055738
Loss: 11.760488960015408
Loss: 11.613228706940784
Loss: 11.492856417067383
Loss: 11.394462669488892
Loss: 11.314034444913947
Loss: 11.248291453438622
Loss: 11.194552346912685
Loss: 11.150625359329684
Loss: 11.114718914975501
Loss: 11.085368558461209
Loss: 11.061377226459106
Loss: 11.041766425106331
Loss: 11.025736321831346
Loss: 11.012633123936265
```

```
1 plot_points_and_y(points, vx, oy)
```



## Exercises

Using the method illustrated above, fit the following equations to our set of points. Use `vx`, `xy` for $x$, $y$, and `va`, `vb`, `vc`, etc for the parameters. This is important, or the tests won't pass.

$$y = a^x + bx + c$$

```
1 ### Exercise: fit of y = a^x + bx + c
```

```
 2
 3 vx = V(0.)
 4 vy = V(0.)
 5 va = V(1.)
 6 vb = V(0.)
 7 vc = V(0.)
 8 # Define below what is oy and loss.
 9 oy = va**vx + vb*vx + vc
10 loss = (vy-oy)*(vy-oy)
11 # YOUR CODE HERE
12 fit(loss, points, [va, vb, vc])
```

⊓→  Loss: 13.22031673565597
    Loss: 11.141654640359056
    Loss: 9.524337885770317
    Loss: 8.258639654229485
    Loss: 7.261959756098997
    Loss: 6.472240243458916
    Loss: 5.842777193129943
    Loss: 5.338271943130829
    Loss: 4.93188598839005
    Loss: 4.603069161367653
    Loss: 4.335966482106359
    Loss: 4.118250607697591
    Loss: 3.940264273483409
    Loss: 3.794387517804188
    Loss: 3.6745678050846706
    Loss: 3.575968472470752
    Loss: 3.494703513553483
    Loss: 3.427635753770303
    Loss: 3.3722219198375187
    Loss: 3.326392689534276
    Loss: 3.288459066630091
    Loss: 3.2570387471581963
    Loss: 3.2309978039038927
    Loss: 3.2094042107284664
    Loss: 3.191490593868335
    Loss: 3.1766242293232927
    Loss: 3.164282770875983
    Loss: 3.1540345391718048
    Loss: 3.1455224617311255
    Loss: 3.1384509501803577
    Loss: 3.1325751510442505
    Loss: 3.127692122085746
    Loss: 3.1236335760432823
    Loss: 3.120259903985071
    Loss: 3.117455245995913
    Loss: 3.1151234209610212
    Loss: 3.113184562377402
    Loss: 3.111572335343445
    Loss: 3.1102316326277197
    Loss: 3.109116666132855
    3.109116666132855

```
 1 ### Tests for convergence of fit of y = a^x + bx + c
 2
```

Now, fit:

$$y = a \cdot 2^x + b \cdot 2^{-x} + cx^3 + dx^2 + ex + f$$

Use a small enough step size, and a sufficient number of iterations, to obtain a final loss of no more than 2.5.

Hint: write `vx * vx * vx`, not `vx ** 3`, etc, since as currently written, the `**` operator cannot handle a negative basis.

```
 1 ### Exercise: fit of y = a 2^x + b 2^{-x} + c x^3 + d x^2 + e x + f
 2 vx = V(0.)
 3 vy = V(0.)
 4 va = V(1.)
 5 vb = V(1.)
 6 vc = V(0.)
 7 vd = V(0.)
 8 ve = V(0.)
 9 vf = V(0.)
10 # Define here what is oy and what is the loss.
11 oy = (va*(2**vx))+(vb*(2**-vx))+(vc*vx*vx*vx)+(vd*vx*vx)+(ve*vx)+vf
12 loss =(vy-oy)*(vy-oy)
13 # YOUR CODE HERE
14 fit(loss, points, [va, vb, vc, vd, ve, vf], delta=0.0000037, num_iterations=200000)
```

```
Loss: 15.776692186667747
Loss: 15.45498304192241
Loss: 15.147787278963065
Loss: 14.853985226726893
Loss: 14.572565903477239
Loss: 14.302615829221551
Loss: 14.043309004157614
Loss: 13.793897931308285
Loss: 13.553705574242992
Loss: 13.32211815218934
Loss: 13.0985786850528
Loss: 12.882581210006673
Loss: 12.673665599505027
Loss: 12.47141291790349
Loss: 12.275441260439719
Loss: 12.085402024204846
Loss: 11.900976566002639
Loss: 11.721873206707228
Loss: 11.54782454595238
Loss: 11.378585054765374
Loss: 11.213928917143638
Loss: 11.053648094603274
Loss: 10.897550590443158
Loss: 10.745458892898395
Loss: 10.597208578533843
Loss: 10.452647059176279
Loss: 10.311632457430129
Loss: 10.174032597382995
Loss: 10.039724098507886
Loss: 9.908591562020703
Loss: 9.78052684007432
Loss: 9.655428379175635
Loss: 9.533200630110436
Loss: 9.413753517467885
Loss: 9.297001962576903
Loss: 9.182865454313387
Loss: 9.071267662815131
Loss: 8.962136091659813
Loss: 8.855401764524878
Loss: 8.750998942763857
Loss: 8.648864870705522
Loss: 8.548939545815221
Loss: 8.451165511156212
Loss: 8.355487667856059
Loss: 8.261853105521825
Loss: 8.170210948762726
Loss: 8.080512218170064
Loss: 7.992709704276649
Loss: 7.90675785317106
Loss: 7.822612662580752
Loss: 7.740231587360472
Loss: 7.659573453433536
Loss: 7.580598379332499
Loss: 7.503267704573782
Loss: 7.427543924180919
Loss: 7.353390628741772
Loss: 7.2807724494485715
```

```
Loss:  7.209655007627211
Loss:  7.140004868312745
Loss:  7.07178949747398
Loss:  7.004977222531232
Loss:  6.939537195847651
Loss:  6.8754393609075235
Loss:  6.81265442092474
Loss:  6.751153809650292
Loss:  6.690909664172089
Loss:  6.6318947995209365
Loss:  6.5740826849157425
Loss:  6.51744742149816
Loss:  6.461963721421729
Loss:  6.407606888174511
Loss:  6.3543527980263885
Loss:  6.302177882503095
Loss:  6.251059111798912
Loss:  6.2009739790486185
Loss:  6.151900485387674
Loss:  6.103817125735722
Loss:  6.056702875246133
Loss:  6.010537176368751
Loss:  5.965299926478887
Loss:  5.9209714660300845
Loss:  5.877532567191842
Loss:  5.834964422937747
Loss:  5.793248636552295
Loss:  5.752367211528142
Loss:  5.712302541827489
Loss:  5.673037402484516
Loss:  5.634554940527233
Loss:  5.596838666199466
Loss:  5.559872444465295
Loss:  5.523640486779782
Loss:  5.48812734311127
Loss:  5.45331789420183
Loss:  5.419197344053528
Loss:  5.385751212629192
Loss:  5.35296532875732
Loss:  5.320825823231576
Loss:  5.289319122096107
Loss:  5.258431940108489
Loss:  5.22815127437287
Loss:  5.198464398136308
Loss:  5.169358854741858
Loss:  5.140822451732372
Loss:  5.112843255099512
Loss:  5.085409583672536
Loss:  5.058510003642275
Loss:  5.0321333232153735
Loss:  5.006268587394808
Loss:  4.9809050728823605
Loss:  4.956032283099364
Loss:  4.931639943322195
Loss:  4.9077179959288255
Loss:  4.884256595753436
Loss:  4.861246105545932
Loss:  4.838677091533381
```