

# Crunching Numbers with NumPy and Cython

Brian Lester

Interactions

January 2, 2019

## Me

- Work on NLP at Interactions specializing in Deep Learning
- Published/Maintain several python packages
  - string-distance
  - text-rank
  - mead-baseline

- We use chat bots to “facilitate customer care interaction” aka to build better chat bots
- I mostly do DL and am a big Nerd so if you every want to talk about ML for NLP I’m game
- Python is my daily driver
- Have several packages on PyPI including ones written in cython (string-distance), the subject of today’s talk

## The Problem

Suppose I have a numpy array of  $N$  points in 2-space; what's the super-fast, pure-numpy way to get a symmetric,  $N \times N$  matrix of pairwise manhattan distances? Here's a simple example and a slow, naive solution:

Figure: Ann Arbor Machine Learning and Data Science Slack

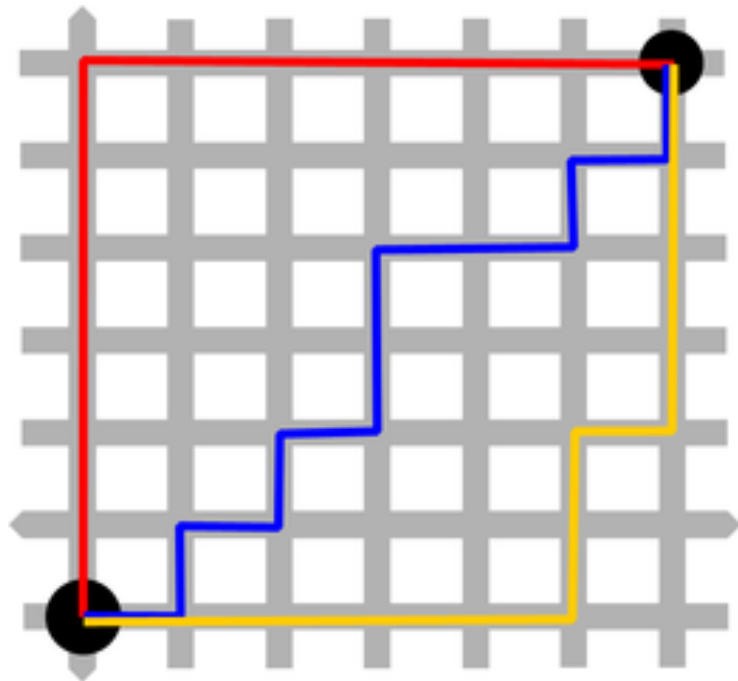
- a2mads is the Ann Arbor Machine Learning and Data Science Slack
- There aren't a lot of messages flying about but when someone asks a question there is normally a lot of discussion
- This user was asking that given a list of points what is the distance between a point and every other point in the list for each item in the list
- There was a long thread about how to optimize his solution, we will walk through some of them.

## Manhattan Distance

- The distance between points measured along axes at right angles

- Distance if you can only move along the axes
- Contrast with euclidean distance where we can move diagonally and go straight there.
- Like you have directions to somewhere in Manhattan, you can't smash through a building, you need to drive along the streets

- Here we can see example of moving between points as the Manhattan distance does



## Manhattan Distance

- The distance between points measured along axes at right angles
- Also called the L1, cityblock, or taxi distance

- Distance if you can only move along the axes
- Contrast with euclidean distance where we can move diagonally and go straight there.
- Like you have directions to somewhere in Manhattan, you can't smash through a building, you need to drive along the streets

## Manhattan Distance

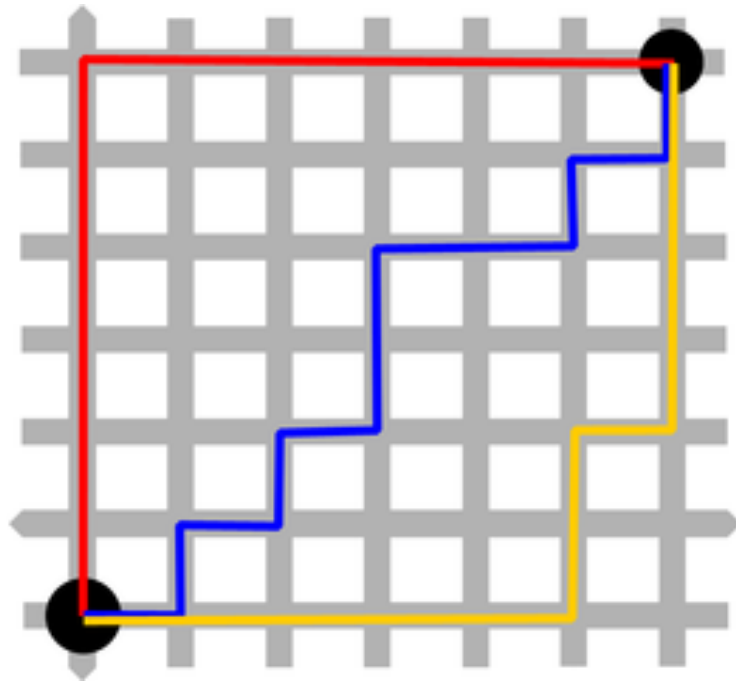
- While the Manhattan distance works for real numbers we will only be looking at Ints in code

A point with  $M$  dimensions is a vector of  $M$  numbers where the values at index  $i$  is the distance along the  $i^{th}$  axis.

$$X \in \mathbb{R}^M$$

$$Y \in \mathbb{R}^M$$

$$\text{ManhattanDistance}(X, Y) = \sum_i^M |X_i - Y_i|$$



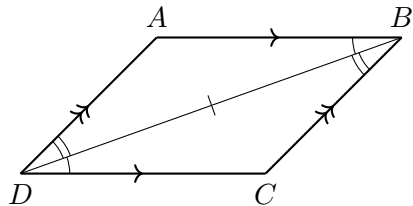
- As suggested by the calculation for Manhattan Distance all of these paths have the same length, the next few slides sketch out why these are all the same but that isn't very important to us
- Point out the absolute distance between point dims on the red line, 6 up 6 over for a distance of 12
- All paths are length 12



## Proof of Equal Lengths

Proof.

let  $\overline{AB} \parallel \overline{DC}$  and  $\overline{AD} \parallel \overline{BC}$  by  
the definition of a parallelogram.



$$\angle ABD \cong \angle CDB$$

$$\angle ADB \cong \angle CBD$$

$$\overline{DB} = \overline{DB}$$

$$\triangle ADB \cong \triangle CBD$$

$$\overline{AB} = \overline{DC}$$

$$\overline{AD} = \overline{BC}$$



- Quick overview of how on a regular grid you can move around a block in either direction and it will be the same length.
- DB is a transversal of both sets of parallel lines
- ABD and CDB are congruent because they are alternate interior angles
- ADB and DBC are congruent because they are alternate interior angles
- It is the same length as itself
- Triangles are congruent by angle side angle
- AB and CD are the same because they are the same sides of congruent triangles
- AD and CB are the same because they are the same sides of congruent triangles

## All Paths Equal

- Let  $R$  be a move 1 step *right* and  $U$  be a move 1 step *up*.
- Let the source point be  $(0,0)$  in our new coordinate system and the target be  $(p,q)$ .
- Any possible path from source to target will have  $p$   $R$  moves and  $q$   $U$  moves.
- There are  $p + q$  moves in total.
- There are  $\binom{p+q}{p}$  possible paths all of the same length
- If you make an extra  $R$  or  $U$  move you will not be at your target and will need to make extra steps so it will not be the shortest path.

- This is a quick overview on why all paths between two points have the same length
- You can scramble the order of moves and the path still has the same number of  $R$ s and  $U$ s but ends at the same spot
- It also sketches why the path length is always the length + height of the rectangle we make with the points as boxes
- This matches the formula where we take the absolute difference between points and sum across dimensions
- you can have a more rigorous proof via induction / contradiction that is based on the fact that if there was any shorter way to get to the next point then we could have gotten to this point through that but I didn't write it out here.
- The basic point is that because all paths are equal we can just look at the differences between the start and end, we don't need to care about the path take there.

# NumPy

- NumPy: A collection of linear algebra functions that operate on a ndarray

- ndarray has a data type
  - all the items in it need to be the same and they can be stored contiguously in memory.
  - So instead of following pointers for to the object for each element in a python list we can access the element directly
  - This makes access fast and because they are together in memory we get better cache utilization
- ndarray is a multi dimensional array
  - a thing is a scalar
  - a list of scalars is an array (or vector)
  - a list of arrays is a matrix
  - a list of matrices is a ndarray
  - You can index like `x[i, j, k]`
  - can have an arbitrary dimensions
- A C extension is code written in C that can be called from python. You can use the speed of C and let go of the GIL. Loses ease of Python

# NumPy

- **NumPy**: A collection of linear algebra functions that operate on a ndarray
- ndarray is a typed container that can represent multi-dimensional arrays

- ndarray has a data type
  - all the items in it need to be the same and they can be stored contiguously in memory.
  - So instead of following pointers for to the object for each element in a python list we can access the element directly
  - This makes access fast and because they are together in memory we get better cache utilization
- ndarray is a multi dimensional array
  - a thing is a scalar
  - a list of scalars is an array (or vector)
  - a list of arrays is a matrix
  - a list of matrices is a ndarray
  - You can index like `x[i, j, k]`
  - can have an arbitrary dimensions
- A C extension is code written in C that can be called from python. You can use the speed of C and let go of the GIL. Loses ease of Python

# NumPy

- **NumPy**: A collection of linear algebra functions that operate on a `ndarray`
- `ndarray` is a typed container that can represent multi-dimensional arrays
- The majority of functionality is implemented in C extensions for speed

- `ndarray` has a data type
  - all the items in it need to be the same and they can be stored contiguously in memory.
  - So instead of following pointers for to the object for each element in a python list we can access the element directly
  - This makes access fast and because they are together in memory we get better cache utilization
- `ndarray` is a multi dimensional array
  - a thing is a scalar
  - a list of scalars is an array (or vector)
  - a list of arrays is a matrix
  - a list of matrices is a `ndarray`
  - You can index like `x[i, j, k]`
  - can have an arbitrary dimensions
- A C extension is code written in C that can be called from python. You can use the speed of C and let go of the GIL. Loses ease of Python

## Cython

- **Cython** is a compiler for Python and Cython

- Cython is a superset of python, adds things like static types
- Compiling just python gives a slight speed boost
- You can also write code that releases the GIL so multithreading is possible
- We will see an example of why calling python code is nice

## Cython

- **Cython** is a compiler for Python and Cython
- Cython is also a language that looks more like python than C, but adds types

- Cython is a superset of python, adds things like static types
- Compiling just python gives a slight speed boost
- You can also write code that releases the GIL so multithreading is possible
- We will see an example of why calling python code is nice

## Cython

- **Cython** is a compiler for Python and Cython
- Cython is also a language that looks more like python than C, but adds types
- It creates C extensions

- Cython is a superset of python, adds things like static types
- Compiling just python gives a slight speed boost
- You can also write code that releases the GIL so multithreading is possible
- We will see an example of why calling python code is nice



## Cython

- **Cython** is a compiler for Python and Cython
- Cython is also a language that looks more like python than C, but adds types
- It creates C extensions
- It can compile pure python code

- Cython is a superset of python, adds things like static types
- Compiling just python gives a slight speed boost
- You can also write code that releases the GIL so multithreading is possible
- We will see an example of why calling python code is nice

## Cython

- **Cython** is a compiler for Python and Cython
- Cython is also a language that looks more like python than C, but adds types
- It creates C extensions
- It can compile pure python code
- You can call python code from inside the cython code (and therefore from inside the C extension)

- Cython is a superset of python, adds things like static types
- Compiling just python gives a slight speed boost
- You can also write code that releases the GIL so multithreading is possible
- We will see an example of why calling python code is nice

# Cython

- **Cython** is a compiler for Python and Cython
- Cython is also a language that looks more like python than C, but adds types
- It creates C extensions
- It can compile pure python code
- You can call python code from inside the cython code (and therefore from inside the C extension)
- Makes it easy to wrap C libraries

- Cython is a superset of python, adds things like static types
- Compiling just python gives a slight speed boost
- You can also write code that releases the GIL so multithreading is possible
- We will see an example of why calling python code is nice

## Cython

- **Cython** is a compiler for Python and Cython
- Cython is also a language that looks more like python than C, but adds types
- It creates C extensions
- It can compile pure python code
- You can call python code from inside the cython code (and therefore from inside the C extension)
- Makes it easy to wrap C libraries
- Some large libraries like **SpaCy** are all cython

- Cython is a superset of python, adds things like static types
- Compiling just python gives a slight speed boost
- You can also write code that releases the GIL so multithreading is possible
- We will see an example of why calling python code is nice

# Cython

- **Cython** is a compiler for Python and Cython
- Cython is also a language that looks more like python than C, but adds types
- It creates C extensions
- It can compile pure python code
- You can call python code from inside the cython code (and therefore from inside the C extension)
- Makes it easy to wrap C libraries
- Some large libraries like **SpaCy** are all cython
- `pip install cython`

- Cython is a superset of python, adds things like static types
- Compiling just python gives a slight speed boost
- You can also write code that releases the GIL so multithreading is possible
- We will see an example of why calling python code is nice

Let get to the code

## Python Brute

```
def pairwise_manhattan_python_v1(
    points: List[List[int]]
) -> List[List[int]]:
    results = []
    # For each pair of points
    for x in points:
        dists = []
        for y in points:
            # Manhattan distance
            dist.append(
                sum(abs(x_i - y_i) for x_i, y_i in zip(x, y))
            )
        results.append(dist)
    return results
```

- It is 2020, python2 is dead, lets all start using type hints
- Input is a list of  $N$  points where each point is a List of  $M$  integers
- This is the most brute force approach, for all combinations of points calculate the Manhattan distance

## Python Manhattan Distance

```
sum(abs(x_i - y_i) for x_i, y_i in zip(x, y))
```

- This is the calculation of manhattan distance directly translated from Math



## Python Brute

```
def pairwise_manhattan_python_v1(
    points: List[List[int]]
) -> List[List[int]]:
    results = []
    # For each pair of points
    for x in points:
        dists = []
        for y in points:
            # Manhattan distance
            dist.append(
                sum(abs(x_i - y_i) for x_i, y_i in zip(x, y))
            )
        results.append(dist)
    return results
```

- $O(N^2M)$  in runtime
- $O(N^2)$  in memory

- Big-O notation
- for  $x$  in points  $\rightarrow$  a loop over  $N$  points  $O(N)$
- for  $y$  in points  $\rightarrow$  a loop over  $N$  points for each point in the  $N$  points,  $O(N^2)$
- for  $x_i, y_i$  in  $\text{zip}(x, y) \rightarrow$  a loop over  $M$  dimensions,  $O(N^2M)$
- Memory used is  $O(N^2)$  for the result matrix, external memory datastructures are not used

## Python Brute

```
def pairwise_manhattan_python_v1(
    points: List[List[int]]
) -> List[List[int]]:
    results = []
    # For each pair of points
    for x in points:
        dists = []
        for y in points:
            # Manhattan distance
            dist.append(
                sum(abs(x_i - y_i) for x_i, y_i in zip(x, y))
            )
        results.append(dist)
    return results
```

- $O(N^2M)$  in runtime
- $O(N^2)$  in memory

- When optimizing you want to save all the time you can in the inner most loop
- every second wasted in the setup is only wasted once
- every second wasted in the inner loop is wasted over and over
- Python function calls have overhead so by manually inlining the Manhattan Distance calculation we get a speed up
- Without in-lining we pay for the call overhead  $O(N^2)$  times

## Python Cached

```
def pairwise_manhattan_python_v2(
    points: List[List[int]]
) -> List[List[int]]:
    # Preallocate the results
    results = [[None] * len(points) for _ in range(len(points))]
    # Look at every pair of you and the ones after you, distances
    # with points before you were calculated when they looked at you
    for i in range(len(points)):
        for j in range(i, len(points)):
            # Manhattan distance calculation
            dist = sum(
                abs(p1 - p2) for p1, p2 in zip(points[i], points[j])
            )
            # Manhattan distance is symmetric
            results[i][j] = dist
            results[j][i] = dist
    return results
```

- This is the only algorithmic improvement (in that we are doing less work) to our code
- The rest of our improvements are all computational and should all show similar speed curves
- The Manhattan distance is symmetric  
( $\text{ManhattanDistance}(X, Y) = \text{ManhattanDistance}(Y, X)$ )
- The distance from point  $i$  to point  $j$  is the same as the distance from point  $j$  to point  $i$

## Symmetric Manhattan Distance

$$\text{ManhattanDistance}(X, Y) = \sum_i^M |X_i - Y_i|$$

Proof.

$$\text{ManhattanDistance}(X, Y) = \text{ManhattanDistance}(Y, X)$$

Case (1)

let  $x - y > 0$

$$y - x < 0$$

$$|x - y| = x - y$$

$$|y - x| = -(y - x)$$

$$-(y - x) = x - y$$

- The Manhattan distance is a distance metric and it be one mathematically it needs to be symmetrical but we can see why here
- The absolute difference is the same when reversed so the sum is of the same values to the total sum is the same
- in the case where the difference is greater than one
- The switched must be negative
- Because we are positive value is just the difference
- The switched is the difference multiplied  $-1$  because it negative and need to be positive
- distribute the multiplication and rearrange

## Symetric Manhattan Distance

$$\text{ManhattanDistance}(X, Y) = \sum_i^M |X_i - Y_i|$$

Proof.

$$\text{ManhattanDistance}(X, Y) = \text{ManhattanDistance}(Y, X)$$

Case (2)

let  $x - y < 0$

$$y - x > 0$$

$$|y - x| = y - x$$

$$|x - y| = -(x - y)$$

$$-(x - y) = y - x$$

- in the case where the difference is less than one
- The switched must be positive
- Our result if the difference multiplied by  $-1$  because it negative and need to be positive
- Because the switched is positive its value is just the difference
- distribute the multiplication and rearrange

## Symmetric Manhattan Distance

$$\text{ManhattanDistance}(X, Y) = \sum_i^M |X_i - Y_i|$$

Proof.

$$\text{ManhattanDistance}(X, Y) = \text{ManhattanDistance}(Y, X)$$

Case (3)

let  $x - y = 0$

$$y - x = 0$$

$$|x - y| = 0$$

$$|y - x| = 0$$



- In the case we are zero the switch is zero so we have the same value.

## Python Cached

```
def pairwise_manhattan_python_v2(
    points: List[List[int]]
) -> List[List[int]]:
    # Preallocate the results
    results = [[None] * len(points) for _ in range(len(points))]
    # Look at every pair of you and the ones after you, distances
    # with points before you were calculated when they looked at you
    for i in range(len(points)):
        for j in range(i, len(points)):
            # Manhattan distance calculation
            dist = sum(
                abs(p1 - p2) for p1, p2 in zip(points[i], points[j])
            )
            # Manhattan distance is symmetric
            results[i][j] = dist
            results[j][i] = dist
    return results
```

- We can reuse the result from one direction and skip computing it for the other
- Big O throws away constants so the is still  $O(N^2M)$
- Theoretically this isn't much faster but in parctice because we are skipping so much work it is faster
- We could use a similar trick to reduce memory usage by just having user of the results always index with some rule where the tuple of points need to be sorted or something similar, not worth it so still  $O(N^2)$  in memory.
- You could optimize a little bit here by initalizing the results matrix with zeros and skipping computing your distance to your self.

## Numpy Brute

- This is a pure port of original python version including the duplicated computation
- This version introduces a lot of numpy concepts but we might as well look at the next version to start

```
def pairwise_manhattan_numpy_v1(  
    points: List[np.ndarray]  
) -> np.ndarray:  
    results = np.zeros((len(points), len(points)), dtype=np.int32)  
    for i in range(len(points)):  
        for j in range(len(points)):  
            results[i, j] = np.sum(np.abs(points[i] - points[j]))  
    return results
```



## Numpy Cached

```
def pairwise_manhattan_numpy_v2(
    points: List[np.ndarray]
) -> np.ndarray:
    results = np.zeros((len(points), len(points)), dtype=np.int32)
    for i in range(len(points)):
        for j in range(i, len(points)):
            dist = np.sum(np.abs(points[i] - points[j]))
            results[i, j] = dist
            results[j, i] = dist
    return results
```

- Preallocating the results as a large matrix
- Loop over all unique pairs of points
- Calculate the manhattan distance with numpy and save results into the matrix
- This was basically the first solution in the a2mads thread

## Numpy Manhattan Distance

```
np.sum(np.abs(points[i] - points[j]))
```

- There are no loops visible in this code!
- This introduces us to vectorization!

## Vectorization

- Code that can transparently operate on lists of numbers

```
def add(x, y):  
    if isinstance(x, list):  
        return [x_i + y_i for x, y in zip(x, y)]  
    return x + y
```

- Moves the loops out of slow python into fast C

- A large part of performant python is trying to remove loops by vectorizing code
- The name implies it might be parallel processing but that depends on the back end choices
- Vectorization generally does these element wise applications but there are also reductions

## Reductions

- Like vectorization operates on a list
- Reduces a list into a scalar by performing some operation

```
def mean(x):  
    return sum(a for a in x) / len(x)
```

- This reduction is from a list to a scalar but as we get to more dims we can decide how far we want to reduce things

## Numpy Manhattan Distance

```
np.sum(np.abs(points[i] - points[j]))
```

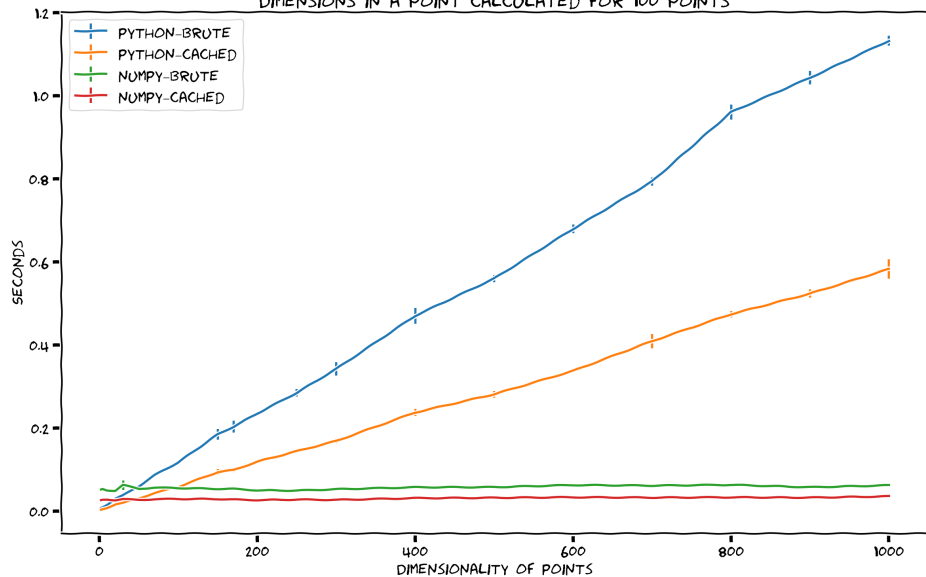
- points  $i$  and  $j$  are  $M$  dimensional points
- The subtraction and the `np.abs` is a vectorized function
- The `np.sum` is a reduction

## Numpy Cached

```
def pairwise_manhattan_numpy_v2(
    points: List[np.ndarray]
) -> np.ndarray:
    results = np.zeros((len(points), len(points)), dtype=np.int32)
    for i in range(len(points)):
        for j in range(i, len(points)):
            dist = np.sum(np.abs(points[i] - points[j]))
            results[i, j] = dist
            results[j, i] = dist
    return results
```

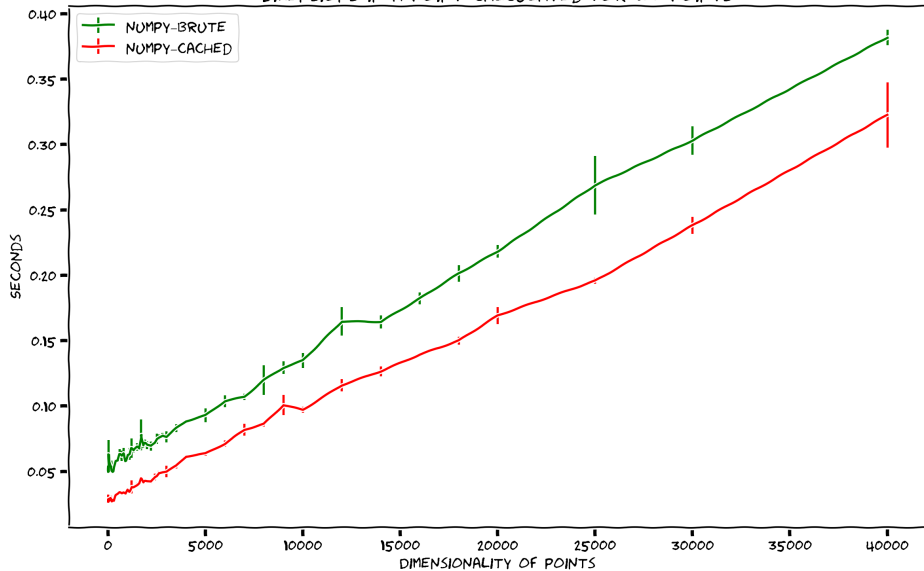
- Then only thing we are vectorizing here is the manhattan distance calculation
- We expect this to be about the same speed as the python v2 as we grow the number of points
- We expect this to show speed gains as we increase the number of dimensions

PAIRWISE MANHATTAN DISTANCE TIMING VS NUMBER OF DIMENSIONS IN A POINT CALCULATED FOR 100 POINTS



- This is a plot of the timing vs the number of dimensions in a point.
- Graphs created by running each implementation on randomly generated data (same data across implementations) 5 times. The point is the mean of the times and the error bar is the standard deviation of the times.
- For the `python_brute` and `python_cached` implementations we can see the predicted  $O(M)$  growth.
- At small input sizes we can see there is some overhead for converting lists into numpy that counteracts the advantage of vectorization

PAIRWISE MANHATTAN DISTANCE TIMING VS NUMBER OF DIMENSIONS IN A POINT CALCULATED FOR 100 POINTS



- Here we can see the same  $O(M)$  growth for numpy implemenetations where we have vectorized the manhattan distance
- This highlights how we are making computational improvments, not algorithmic ones
- I had to scale out to a lot more points to see this trend because when python code is in the graph it looks flat
- Note the very different scale on the y-axis



## Numpy Broadcast

- We can see that one of the loops is gone!
- If you guessed the answer is vectorization you are right
- This is specifically the idea of broadcasting which help you vectorize code

```
def pairwise_manhattan_numpy_v3(points: np.ndarray) -> np.ndarray:  
    results = []  
    for i in range(len(points)):  
        results.append(np.sum(np.abs(points[i] - points), axis=1))  
    return np.stack(results)
```

# Broadcasting

- Vectorization can speed up code

# Broadcasting

- Vectorization can speed up code
- But what do I do when the shapes don't match?

## Broadcasting

- Vectorization can speed up code
- But what do I do when the shapes don't match?
- If I want to add a scalar to a vector do I have to copy the scalar into a vector to vectorize it?

- Copying would be a waste of memory (need to have the scalar multiple times) and computation (time taken to make the copies of the scalar)

## Broadcasting

- Vectorization can speed up code
- But what do I do when the shapes don't match?
- If I want to add a scalar to a vector do I have to copy the scalar into a vector to vectorize it?
- Broadcasting to the rescue!

- Copying would be a waste of memory (need to have the scalar multiple times) and computation (time taken to make the copies of the scalar)

## Roll Your Own Broadcasting

```
def add_scalar_to_vector(scalar, vector):  
    return [v + scalar for v in vector]
```

- The scalar is added to each element in the vector as if the scalar was expanded to a vector and a normal vectorized add was done.

## Broadcasting Rules

- 1 If there is a mismatch in the number of dimensions the shape with fewer is padded with 1s on its leading side.
  - $2 \times 3 + 3 \rightarrow 2 \times 3 + 1 \times 3$
- 2 If the size of some dimension doesn't match and one of the sizes is 1 it is repeated to match the other size.
  - $5 \times 4 + 5 \times 1 \rightarrow 5 \times 4 + 5 \times 4$
- 3 If there is a mismatch and neither size is 1 an exception is thrown.
  - $3 \times 4 + 4 \times 3 \rightarrow \text{ERROR}$

## Broadcasting Examples

$$\begin{array}{ccc|c} X \in \mathbb{R}^{2 \times 3} & Y \in \mathbb{R}^{1 \times 3} & Y' \in \mathbb{R}^{2 \times 3} & X + Y \in \mathbb{R}^{2 \times 3} \\ \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} & \begin{bmatrix} 1 & 2 & 1 \end{bmatrix} & \rightarrow \begin{bmatrix} 1 & 2 & 1 \\ 1 & 2 & 1 \end{bmatrix} & \begin{bmatrix} 2 & 4 & 4 \\ 5 & 7 & 7 \end{bmatrix} \end{array}$$

$$\begin{array}{ccc|c} X \in \mathbb{R}^{2 \times 3} & Z \in \mathbb{R}^{2 \times 1} & Z' \in \mathbb{R}^{2 \times 3} & X + Z \in \mathbb{R}^{2 \times 3} \\ \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} & \begin{bmatrix} 2 \\ 3 \end{bmatrix} & \rightarrow \begin{bmatrix} 2 & 2 & 2 \\ 3 & 3 & 3 \end{bmatrix} & \begin{bmatrix} 3 & 4 & 5 \\ 7 & 8 & 9 \end{bmatrix} \end{array}$$

- Here we see examples of broadcasting
- The first one is broadcasting over the rows
- The second one is broadcasting over the columns
- Stress that these copies don't actually happen but computation acts like it did



## Numpy Broadcast

```
def pairwise_manhattan_numpy_v3(points: np.ndarray) -> np.ndarray:  
    results = []  
    for i in range(len(points)):  
        results.append(np.sum(np.abs(points[i] - points), axis=1))  
    return np.stack(results)
```

- points  $i$  is  $\in \mathbb{R}^M$
- points is  $\in \mathbb{R}^{N \times M}$
- The result of broadcasting vectorization for points is  $[N, M]$
- There is this new kwarg for `np.sum`, `axis`
- This is a way to apply reductions to multi dimensional inputs

## Multidimensional Reduction

```
y = np.sum(x)
```

$$x \in \mathbb{R}^{4 \times 3}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \\ 10 & 11 & 12 \end{bmatrix}$$

$$y \in \mathbb{R}$$

$$78$$

- We said before that a reduction turned a vector into a scalar, but what about a high dim?
- Without an axis to reduce the whole things, results in 78

## Reduction Along an Axis

- The axis argument allows for a reduction that only reduces the number of dims by one
- Reduction over the axis 1 is applied across the columns for matrices resulting in a vector

```
y = np.sum(x, axis=1)
```

$$x \in \mathbb{R}^{4 \times 3}$$
$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \\ 10 & 11 & 12 \end{bmatrix}$$

$$y \in \mathbb{R}^4$$
$$\begin{bmatrix} 6 \\ 15 \\ 24 \\ 33 \end{bmatrix}$$

## Reduction Along an Axis

- Reduction over axis 0 is applied across the rows for matrices
- There is also a `keepdims` argument to keep the single dimension of size 1 so instead of getting a vector of size 3 you get one of size [1, 3]

```
y = np.sum(x, axis=0)
```

$$x \in \mathbb{R}^{4 \times 3}$$
$$\downarrow \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \\ 10 & 11 & 12 \end{bmatrix}$$

$$y \in \mathbb{R}^3$$
$$[22 \quad 26 \quad 30]$$

## Numpy Broadcast

```
def pairwise_manhattan_numpy_v3(points: np.ndarray) -> np.ndarray:
    results = []
    for i in range(len(points)):
        results.append(np.sum(np.abs(points[i] - points), axis=1))
    return np.stack(results)
```

- Reducing  $[N, M]$  over axis 1 gives us  $N$
- Results is a list of vectors of size  $N$
- `np.stack` puts them on top of each other to form a single matrix
- Side Note: Type hints for array shapes has no support which sucks. To make broadcasting work you often need specific shapes but you can't communicate these constraints via type hints

## Numpy Broadcast

- One thing to notice here is that we are back to doing redundant computation by calculating all pairs of points!
- It is possible to remove this redundancy but it will get a little weird

```
def pairwise_manhattan_numpy_v3(points: np.ndarray) -> np.ndarray:  
    results = []  
    for i in range(len(points)):  
        results.append(np.sum(np.abs(points[i] - points), axis=1))  
    return np.stack(results)
```

## Numpy Broadcast Cached

```
def pairwise_manhattan_numpy_broadcast_cached(
    points: np.ndarray
) -> np.ndarray:
    results = np.zeros((len(points), len(points)), dtype=np.int32)
    dists = []
    for i in range(len(points)):
        dists.append(np.sum(np.abs(points[i] - points[i:]), axis=1))
    for i, dist in enumerate(dists):
        results[i, i:] = dist
    return results + results.T
```

- Here we are back to preallocating our results
- When we compute the distances we only compute before by slicing the points matrix
- We get back a list of tensors that are not the same shape (each one is on smaller than the previous)
- In python we copy these into the regular matrix so that it is a upper triangular matrix
  - Upper triangular matrix means it only has values above the diagonal
  - Only value when the column index  $j$  is  $\geq$  the row index  $i$
- We transpose (the rows become columns and columns rows) and add. Because we are upper triangle (lower values are all zero) the add just copies upper values in
- This would double the diag values but because we know Manhattan distance is 0 with ourselves the diag will stay 0

## Numpy Broadcast Cached

```
def pairwise_manhattan_numpy_broadcast_cached(
    points: np.ndarray
) -> np.ndarray:
    results = np.zeros((len(points), len(points)), dtype=np.int32)
    dists = []
    for i in range(len(points)):
        dists.append(np.sum(np.abs(points[i] - points[i:]), axis=1))
    for i, dist in enumerate(dists):
        results[i, i:] = dist
    return results + results.T
```

- Still  $O(N^2M)$  runtime.
- Still  $O(N^2)$  memory.

- We are adding a loop over the  $N$  points (copying the values into the results)  $O(N)$
- The addition is  $O(N^2)$
- Whole thing is  $O(N^2M + N^2 + N)$  but we can throw away extra for Big-O
- Even though we have slow python loops at large values we are saving so much computation it is worth it
- Sometimes a strange mix of vectorized and python is the best way.



## Numpy Double Broadcast

```
def pairwise_manhattan_numpy_v4(points: np.ndarray) -> np.ndarray:  
    exp_points = np.expand_dims(points, 1) # [N, 1, M]  
    return np.sum(np.abs(exp_points - points), axis=-1)
```


- All the loops are gone!
- We add a dimension of size 1 into the points
- When we do the subtraction the points  $[N, M]$  is expanded to  $[1, N, M]$
- $[N, 1, M] + [1, N, M] \rightarrow [N, N, M]$
- Reduction over axis  $-1$  (same as negative indexing in python, sum over the last dim)
- Reduction gives us  $[N, N]$  aka our answer

## Numpy Double Broadcast

```
def pairwise_manhattan_numpy_v4(points: np.ndarray) -> np.ndarray:  
    exp_points = np.expand_dims(points, 1) # [N, 1, M]  
    return np.sum(np.abs(exp_points - points), axis=-1)
```


- This idea of double broadcasting is very powerful and is the corner stone of why the Conditional Random Field (CRF) (which is a NLP model) I wrote for mead-baseline is so fast.
- Now that we are double broadcasting there is no way to express the fact that we can skip the repeated work of calculating the distance between symmetric pairs of points.

## Broadcasting Micro Benchmark

In [5]: 

```
1 %%timeit
2 pairwise_manhattan_numpy_v3(points)
```

41.3  $\mu$ s  $\pm$  411 ns per loop (mean  $\pm$  std. dev. of 7 runs, 10000 loops each)

In [6]: 

```
1 %%timeit
2 pairwise_manhattan_numpy_v4(points)
```

8.15  $\mu$ s  $\pm$  43.2 ns per loop (mean  $\pm$  std. dev. of 7 runs, 100000 loops each)

- This was a promising implementation in the thread
- Here we can see on a small micro bench mark (4 points of dim 2) that shows a lot of promise
- As I scaled it out speeds weren't meaningfully different/were worse from the single broadcast
- Lets look at why it was worse

## Numpy Double Broadcast

```
def pairwise_manhattan_numpy_v4(points: np.ndarray) -> np.ndarray:
    exp_points = np.expand_dims(points, 1) # [N, 1, M]
    return np.sum(np.abs(exp_points - points), axis=-1)
```

- Still  $O(N^2M)$  in runtime
- Now  $O(N^2M)$  in memory

- We are going to revisit Big-O for this implementation
- The runtime is still the same, remember the loops are still there, they are just in C not python
- We have this intermediate result of size  $[N, N, M]$  so our memory usage is higher.  $O(N^2M)$
- When you run out of memory you start swapping which involves disk and is therefore super slow
- Due to memory usage I could only get this version to work for limited points and only a few dimensions.

points	dims	Single Memory	Double Memory
10000	2	381.47MiB	762.94MiB
10000	10		3.73GiB
10000	150		55.88GiB
80000	2	23.84GiB	47.68GiB
80000	150		3.49PiB

- This is an example of why it is important to benchmark in a realistic scenario.

## Cython Implementation

```
@cython.wraparound(False)
@cython.boundscheck(False)
cpdef int[:, :] pairwise_manhattan(int[:, :] points):
    cdef int n = points.shape[0]
    cdef int m = points.shape[1]
    results = np.zeros((n, n), dtype=np.int32)
    cdef int[:, :] results_view = results
    cdef int i, j, k, dist
    for i in range(n):
        for j in range(i, n):
            dist = 0
            for k in range(m):
                dist += abs(points[i, k] - points[j, k])
            results_view[i, j] = dist
            results_view[j, i] = dist
    return results
```

- The cython code looks complex but if we zoom in we can see it looks a lot like our original python code
- There are a lot of annotations that help speed things up.
- we can just write loops, we don't have to think about vectorization

## Cython Decorators

```
@cython.wraparound(False)
@cython.boundscheck(False)
```

- These are directives that tell cython how to compile code
- `wraparound` turns off the ability to do negative indexing
- `boundscheck` turns off raising an error when you access past the end of an array
- Both of these ease of use things require the length of the array. C arrays don't carry their lengths so we need to ask the python object for the length. Accessing python objects are slow and require the GIL

## Cython Function Definition

```
cpdef int[:, :] pairwise_manhattan(int[:, :] points):
```

- cpdef is a utility to make this function callable by python code and other cython code. C function calls are faster and there is slight overhead from the cpdef. Cython in a tight inner loop (like a function of manhattan distance) should use cdef
- int[:, :] before the function name is the return type
- It is also a type of the input (points)
- We will get to what the types are in a second

## Cython Variable Declaration

- `cdef` is used to statically declare variables. This tells the compiler what the type is so it can compile it to C.

```
cdef int n = points.shape[0]
cdef int m = points.shape[1]
results = np.zeros((n, n), dtype=np.int32)
cdef int[:, :] results_view = results
cdef int i, j, k, dist
```



## Cython Variable Declaration

```
cdef int n = points.shape[0]
cdef int m = points.shape[1]
results = np.zeros((n, n), dtype=np.int32)
cdef int[:, :] results_view = results
cdef int i, j, k, dist
```

- `results` is a dynamically created numpy array, it is known to the python interpreter, the gil is needed to make it, the ref counter is incremented by one when we make it, and it is a python object that will be safe to use in python after we return it. Making dynamic python objects so easily is a strength of cython.

```

+11:     results = np.zeros((n, n), dtype=np.int32)
    __Pyx_GetModuleGlobalName(__pyx_t_1, __pyx_n_s_np); if (unlikely(!__pyx_t_1)) __PYX_ERR(0, 11, __pyx_L1_error)
    __Pyx_GOTREF(__pyx_t_1);
    __pyx_t_2 = __Pyx_PyObject_GetAttrStr(__pyx_t_1, __pyx_n_s_zeros); if (unlikely(!__pyx_t_2)) __PYX_ERR(0, 11, __pyx_L1_error)
    __Pyx_GOTREF(__pyx_t_2);
    __Pyx_DECREF(__pyx_t_1); __pyx_t_1 = 0;
    __pyx_t_1 = __Pyx_PyInt_From_int(__pyx_v_n); if (unlikely(!__pyx_t_1)) __PYX_ERR(0, 11, __pyx_L1_error)
    __Pyx_GOTREF(__pyx_t_1);
    __pyx_t_3 = __Pyx_PyInt_From_int(__pyx_v_n); if (unlikely(!__pyx_t_3)) __PYX_ERR(0, 11, __pyx_L1_error)
    __Pyx_GOTREF(__pyx_t_3);
    __pyx_t_4 = PyTuple_New(2); if (unlikely(!__pyx_t_4)) __PYX_ERR(0, 11, __pyx_L1_error)
    __Pyx_GOTREF(__pyx_t_4);
    __Pyx_GIVEREF(__pyx_t_1);
    PyTuple_SET_ITEM(__pyx_t_4, 0, __pyx_t_1);
    __Pyx_GIVEREF(__pyx_t_3);
    PyTuple_SET_ITEM(__pyx_t_4, 1, __pyx_t_3);
    __pyx_t_1 = 0;
    __pyx_t_3 = 0;
    __pyx_t_3 = PyTuple_New(1); if (unlikely(!__pyx_t_3)) __PYX_ERR(0, 11, __pyx_L1_error)
    __Pyx_GOTREF(__pyx_t_3);
    __Pyx_GIVEREF(__pyx_t_4);
    PyTuple_SET_ITEM(__pyx_t_3, 0, __pyx_t_4);
    __pyx_t_4 = 0;
    __pyx_t_4 = __Pyx_PyDict_NewPresized(1); if (unlikely(!__pyx_t_4)) __PYX_ERR(0, 11, __pyx_L1_error)
    __Pyx_GOTREF(__pyx_t_4);
    __Pyx_GetModuleGlobalName(__pyx_t_1, __pyx_n_s_np); if (unlikely(!__pyx_t_1)) __PYX_ERR(0, 11, __pyx_L1_error)
    __Pyx_GOTREF(__pyx_t_1);
    __pyx_t_5 = __Pyx_PyObject_GetAttrStr(__pyx_t_1, __pyx_n_s_int32); if (unlikely(!__pyx_t_5)) __PYX_ERR(0, 11, __pyx_L1_error)
    __Pyx_GOTREF(__pyx_t_5);
    __Pyx_DECREF(__pyx_t_1); __pyx_t_1 = 0;
    if (PyDict_SetItem(__pyx_t_4, __pyx_n_s_dtype, __pyx_t_5) < 0) __PYX_ERR(0, 11, __pyx_L1_error)
    __Pyx_DECREF(__pyx_t_5); __pyx_t_5 = 0;
    __pyx_t_5 = __Pyx_PyObject_Call(__pyx_t_2, __pyx_t_3, __pyx_t_4); if (unlikely(!__pyx_t_5)) __PYX_ERR(0, 11, __pyx_L1_error)
    __Pyx_GOTREF(__pyx_t_5);
    __Pyx_DECREF(__pyx_t_2); __pyx_t_2 = 0;

```

- Here we can see some of the code that is generated when creating the results python object. Having to interface with the CPython API like this would be difficult without cython

## Cython Variable Declaration

```
cdef int n = points.shape[0]
cdef int m = points.shape[1]
results = np.zeros((n, n), dtype=np.int32)
cdef int[:, :] results_view = results
cdef int i, j, k, dist
```

- `int[:, :]` is a typed memoryview
- a type memoryview is a view of the memory that backs an object that supports the buffer interface
- the buffer interface specific way that an object lays out its data in memory so other things can access the content. Things like NumPy arrays and the stdlib array class support this interface
- `results_view` is a typed memory view of the memory that holds the underlying data of results. This view lets us access the object from C, without gil

## Cython Manhattan Distance

```
from libc.stdlib cimport abs
cimport cython

@cython.wraparound(False)
@cython.boundscheck(False)
cdef inline int manhattan_distance(int[:] x, int[:] y):
    cdef int i
    cdef int m = x.shape[0]
    cdef int dist = 0
    for i in range(m):
        dist += abs(x[i] - y[i])
    return dist
```

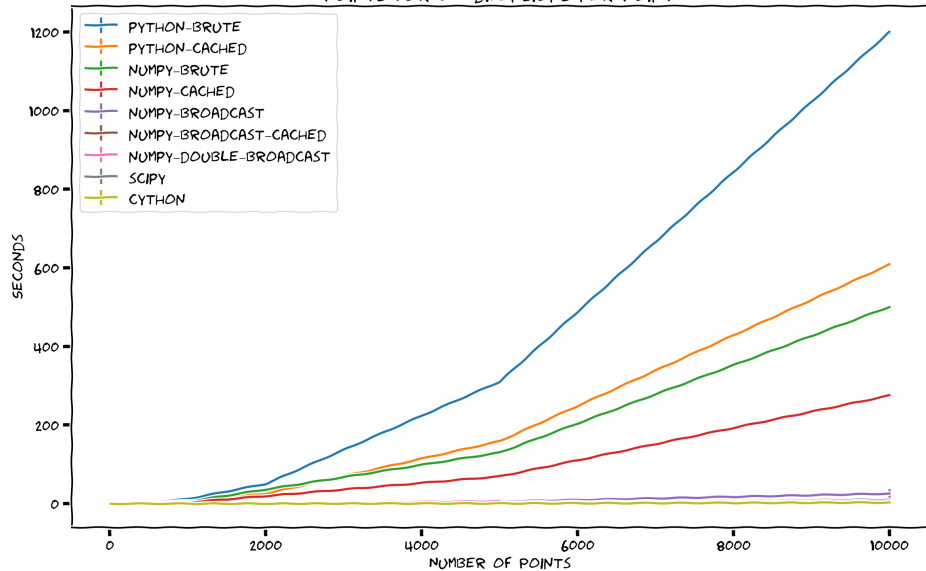
- I didn't do this in the code so it looks like the others but I would if this was real production code.
- cdef means this is only callable from cython but it has fast C function calls instead of slow python calls
- the inline keyword means when this is compiled this code is directly inserted, not called as a function
- This gives us the speed of having code inline but the reusability/nice SWE practice of having a separate function for it.
- We import abs from the C standard library so we always use the C version, if not cython will swap the version (C vs python) based on the datatypes and we wouldn't be able to release the GIL

## Cython Implementation

```
@cython.wraparound(False)
@cython.boundscheck(False)
cpdef int[:, :] pairwise_manhattan(int[:, :] points):
    cdef int n = points.shape[0]
    cdef int m = points.shape[1]
    results = np.zeros((n, n), dtype=np.int32)
    cdef int[:, :] results_view = results
    cdef int i, j, k, dist
    for i in range(n):
        for j in range(i, n):
            dist = 0
            for k in range(m):
                dist += abs(points[i, k] - points[j, k])
            results_view[i, j] = dist
            results_view[j, i] = dist
    return results
```

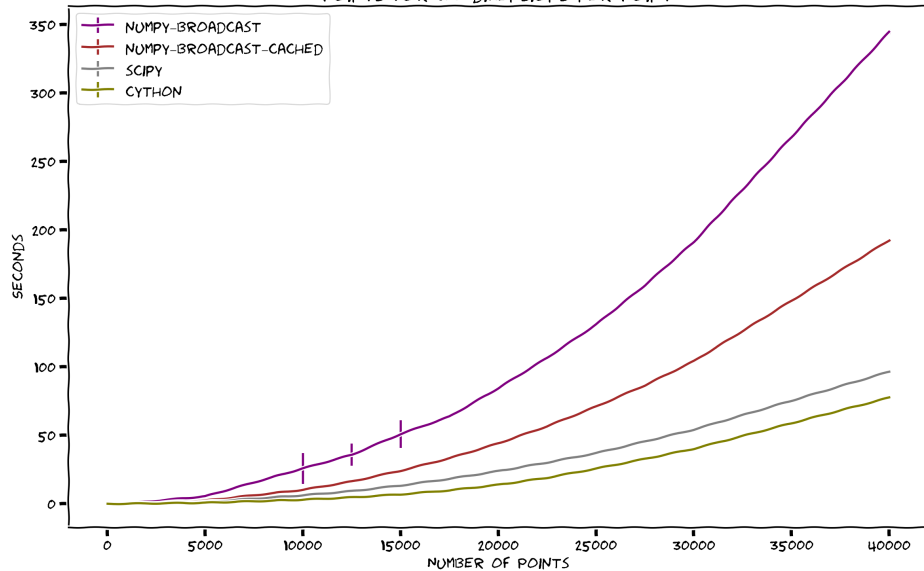
- There are a lot of optimizations left of the table, for example operating on batch points at a time in the inner loops for better cache coherence
- we could skip computing the distance to ourselves because we know it is zero
- These kind of optimizations are very hard (impossible) to express with NumPy
- We didn't use it but Cython also supports defining objects and structs which makes it amenable to speeding up code that isn't all numeric arrays
- NumPy is basically array only if you want speed.

PAIRWISE MANHATTAN DISTANCE TIMING VS NUMBER OF POINTS FOR 100 DIMENSIONS PER POINT



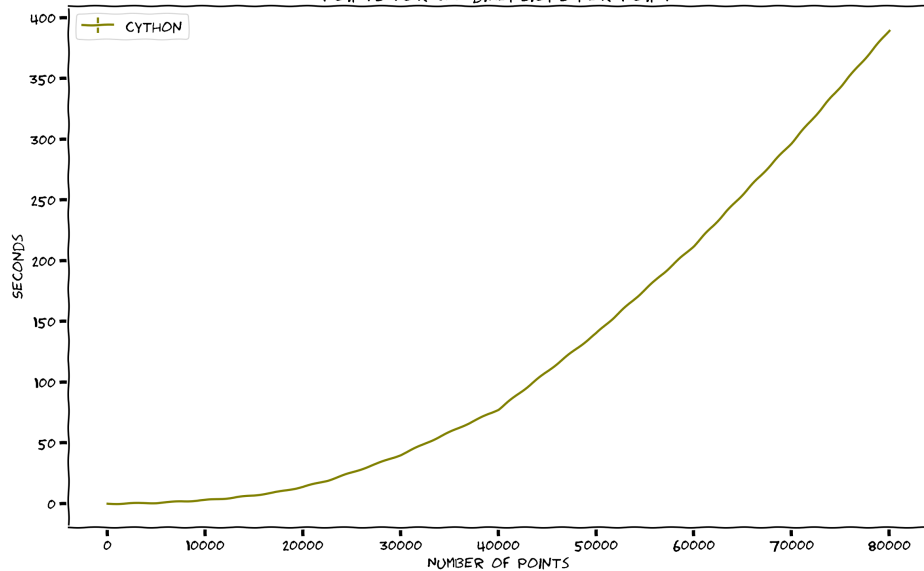
- Here we can see the timing curves of various implementations as we scale up the number of points.
- We can see the  $O(N^2M)$  growth for our algos
- The numpy broadcasting, scipy, and cython solutions are hard to see, we'll zoom in
- SciPy is a set of more specific/advanced scientific functions that includes a manhattan distance function

PAIRWISE MANHATTAN DISTANCE TIMING VS NUMBER OF  
POINTS FOR 100 DIMENSIONS PER POINT



- Here we are looking at only some implementations
  - numpy\_broadcast
  - numpy\_broadcast\_cached
  - scipy
  - cython
- As mentioned earlier the double broadcasting solution didn't show meaningful gains over the broadcasting and I actually couldn't get it to run with 100 dimensional points.
- We can start to see the  $O(N^2M)$  growth

PAIRWISE MANHATTAN DISTANCE TIMING VS NUMBER OF  
POINTS FOR 100 DIMENSIONS PER POINT



- Finally we are scaling out to a lot of points.
- cython only
- the  $O(N^2M)$  growth is obvious
- We are at 8 times as many points but only a third of the runtime as the python impl



## Interactive Cython

- Cython is a compiled language
- This makes rapid iteration harder than in python
- Jupyter notebooks have nice cython features to speed up exploratory work
- Notebook: [Jupyter Notebook](#)

- Cython C extensions are normally built via the `setup.py` file
- A change in cython code isn't reflected until you recompile it
- I am not a huge fan of jupyter notebooks but they can compile cython code for to allow for quick iteration
- Show the notebook
  - A python solution
  - A cython compiled version of the pure python
  - A cython version with types
  - shows how just compiling python code can give slight speed boosts
  - Vizualization where yellow lines are ones that interact with the python interperter and are therefore slow

In [1]: %load\_ext Cython

In [3]: %%cython -a

```
def pairwise_manhattan_cython_untyped(points):
    results = [[None] * len(points) for _ in range(len(points))]
    for i in range(len(points)):
        for j in range(i, len(points)):
            dist = sum(abs(p1 - p2) for p1, p2 in zip(points[i], points[j]))
            results[i][j] = dist
            results[j][i] = dist
    return results
```

Out[3]:

Generated by Cython 0.29.14

Yellow lines hint at Python interaction.

Click on a line that starts with a "+" to see the C code that Cython generated for it.

```
1:
+2: def pairwise_manhattan_cython_untyped(points):
+3:     results = [[None] * len(points) for _ in range(len(points))]
+4:     for i in range(len(points)):
+5:         for j in range(i, len(points)):
+6:             dist = sum(abs(p1 - p2) for p1, p2 in zip(points[i], points[j]))
+7:             results[i][j] = dist
+8:             results[j][i] = dist
+9:     return results
```

- We use the %load\_ext Cython magic so our notebook can compile Cython code
- Here we see the pure python implementation being compiled with Cython
- %%cython -a compiles the code in this cell and the -a adds annotation to the output
- The annotated view below tells us which lines (yellow) interact with the python interpreter
- General rule of thumb is yellow == slow

Out[4]:

Generated by Cython 0.29.14

Yellow lines hint at Python interaction.

Click on a line that starts with a "+" to see the C code that Cython generated for it.

```
+01: from libc.stdlib cimport abs
    02: cimport cython
    03: from cython.parallel cimport prange
+04: import numpy as np
    05:
    06: @cython.wraparound(False)
    07: @cython.boundscheck(False)
+08: cpdef int[:, :] pairwise_manhattan(int[:, :] points):
+09:     cdef int n = points.shape[0]
+10:     cdef int m = points.shape[1]
+11:     results = np.zeros((n, n), dtype=np.int32)
+12:     cdef int[:, :] results_view = results
    13:     cdef int i, j, k, dist
+14:     for i in range(n):
+15:         for j in range(i, n):
+16:             dist = 0
+17:             for k in range(m):
+18:                 dist += abs(points[i, k] - points[j, k])
+19:                 results_view[i, j] = dist
+20:                 results_view[j, i] = dist
+21:     return results
```

- This is the annotated version of the cython code
- We see some yellow, for example creating the results object but this is only done once
- The whole actual calculation is all white, there is no python, only C

- The compiled python version gives a slight boost
- The cython verions blows them out of the water

```
In [6]: %%timeit  
pairwise_manhattan_python(points)  
  
1.39 s ± 22 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)
```

```
In [7]: %%timeit  
pairwise_manhattan_cython_untyped(points)  
  
1.15 s ± 12 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)
```

```
In [8]: %%timeit  
pairwise_manhattan(points)  
  
1.24 ms ± 5.32 µs per loop (mean ± std. dev. of 7 runs, 1000 loops each)
```

## Conclusion

- Vectorization and Broadcasting in NumPy lets us push loops from python into C
- Cython lets us write fast loops and is useful for speeding up difficult to vectorize code
- If you have a custom use case cython can often beat the more general NumPy solution

- As we saw in the numpy broadcasting it was a lot of extra work/complexity to remove the redundant computation which is easy in cython.
- As we saw in the numpy double broadcasting we can't remove the redundant computation.

## Links

- Notebook:  
<https://nbviewer.jupyter.org/github/blester125/MIPy-Talk-Jan-2-2020/blob/master/scripts/cython-jupyter-example.ipynb>
- Slides: <https://github.com/blester125/MIPy-Talk-Jan-2-2020/blob/master/slides/slides.pdf>
- Twitter: <https://twitter.com/BrianLester125>
- GitHub: <https://github.com/blester125>

- The link to the slides is in a GitHub repo that has all these implementations, code to reproduce the graphs, and a docker container that has optimized versions of NumPy installed.