

---

# **Quantitative Big Imaging - Advanced segmentation**

**Anders Kaestner**

**Mar 27, 2024**



# CONTENTS

0.1	Advanced segmentation . . . . .	1
0.2	Where segmentation fails . . . . .	3
0.3	Apply thresholds . . . . .	8
0.4	More Complicated Images . . . . .	18
0.5	Adding Information . . . . .	21
0.6	Clustering / Classification (Unsupervised) . . . . .	25
0.7	Quad trees . . . . .	34
0.8	Superpixels . . . . .	36
0.9	Probabilistic Models of Segmentation . . . . .	39
0.10	Summary . . . . .	39
0.11	Supervised Segmentation Approaches . . . . .	40
0.12	Basic Methods Overview . . . . .	41
0.13	Classification . . . . .	41
0.14	Linear Regression . . . . .	47
0.15	Decision trees . . . . .	50
0.16	Pipelines . . . . .	57
0.17	Classification using a pipeline . . . . .	61
0.18	Regression using a pipeline . . . . .	66
0.19	Assessment of multi-category data . . . . .	68
0.20	Segmentation (Pixel Classification) . . . . .	70
0.21	Deep learning with a U-Net . . . . .	88
0.22	Summary . . . . .	100



This is the lecture notes for the 5th lecture of the Quantitative big imaging class given during the spring semester 2021 at ETH Zurich, Switzerland.

## 0.1 Advanced segmentation

### 0.1.1 Today's Outline

- Motivation
  - Many Samples
  - Difficult Samples
  - Training / Learning
  - Thresholding
  - Automated Methods
  - Hysteresis Method
- 

- Feature Vectors
- K-Means Clustering
- Superpixels
- Working with Segmented Images
- Contouring

### 0.1.2 Literature / Useful References

- Superpixels

### 0.1.3 Load some needed modules

```
from skimage.io import imread
import matplotlib.pyplot as plt

import numpy as np
from skimage.morphology import dilation, opening, disk
from collections import OrderedDict
from skimage.data import page
import skimage.filters as flt
import pandas as pd
from skimage.filters import gaussian, median, threshold_triangle
from sklearn.cluster import KMeans
from sklearn.datasets import make_blobs
from matplotlib.colors import ListedColormap
import matplotlib.colors as colors
%matplotlib inline

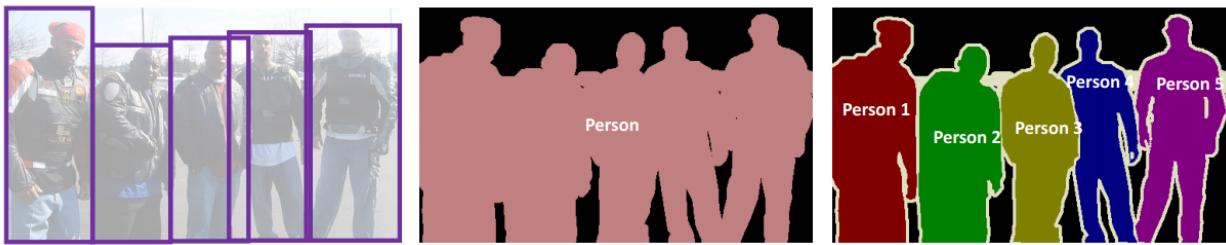
colors = plt.rcParams['axes.prop_cycle'].by_key()['color']
```

### 0.1.4 Previously on QBI

- Image acquisition and representations
  - Enhancement and noise reduction
  - Understanding models and interpreting histograms
  - Ground Truth and ROC Curves
- 
- Choosing a threshold
  - Examining more complicated, multivariate data sets
  - Improving segmentation with morphological operations
  - Filling holes
  - Connecting objects
  - Removing Noise
  - Partial Volume Effect

### 0.1.5 Different types of segmentation

When we talk about image segmentation there are different meanings to the word. In general, segmentation is an operation that marks up the image based on pixels or pixel regions. This is a task that has been performed since beginning of image processing as it is a natural step in the workflow to analyze images - we must know which regions we want to analyze and this is a tedious and error prone task to perform manually. Looking at the figure below we two type of segmentation.



Object Detection

Semantic Segmentation

Instance Segmentation

- Object detection - identifies regions containing an object. The exact boundary does not matter so much here. We are only interested in a bounding box.
- Semantic segmentation - classifies all the pixels of an image into meaningful classes of objects. These classes are “semantically interpretable” and correspond to real-world categories. For instance, you could isolate all the pixels associated with a cat and color them green. This is also known as dense prediction because it predicts the meaning of each pixel.
- Instance segmentation - Identifies each instance of each object in an image. It differs from semantic segmentation in that it doesn't categorize every pixel. If there are three cars in an image, semantic segmentation classifies all the cars as one instance, while instance segmentation identifies each individual car.

## 0.2 Where segmentation fails

Segmentation using a single threshold is sensitive to different conditions...

In fact, segmentation is rarely an obvious task. What you want to find is often obscured by other image features and unwanted artefacts from the experiment technique. If you take a glance at the painting by Bev Doolittle, you quickly spot the red fox in the middle. Looking a little closer at the painting, you'll recognize two Indians on their spotted ponies. This example illustrates the problems you will encounter when you start to work with image segmentation.



Fig. 1: Cases making the segmentation task harder than just applying a single threshold.

*Woodland Encounter* Bev Doolittle

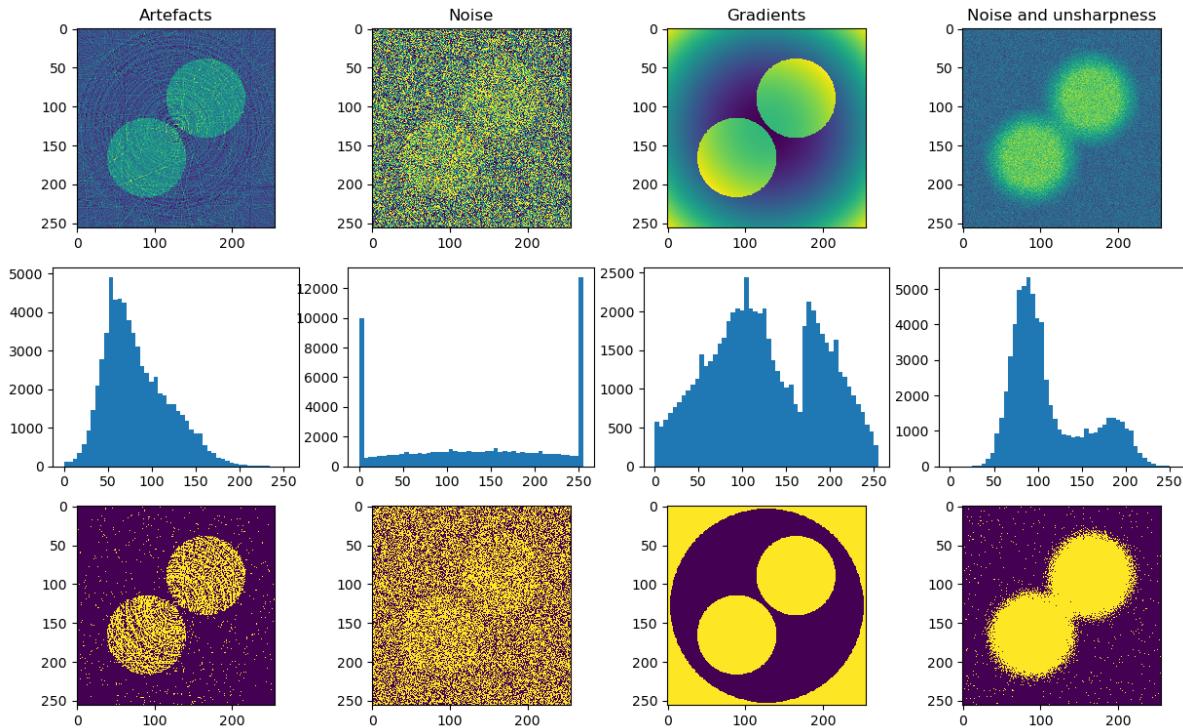
### 0.2.1 Typical image features that makes life harder

The basic segmentation shown in the previous lecture can only be used under good conditions when the classes are well separated. Images from many experiments are unfortunately not well-behaved in many cases. The figure below shows four different cases when an advanced technique is needed to segment the image. Neutron images often show a combination of all cases.

The impact on the segmentation performance of all these cases can be reduced by proper experiment planning. Still, sometimes these improvements are not fully implemented in the experiment to be able to fulfill other experiment criteria.

```
files = ['class_art_in.png', 'class_wgn_in.png', 'class_cupped_in.png', 'class_wgn_
↪smooth_in.png']
titles = ['Artefacts', 'Noise', 'Gradients', 'Noise and unsharpness']

plt.figure(figsize=[15, 9])
for i in range(4) :
    img = imread('data/' + files[i]).mean(axis=2)
    plt.subplot(3, 4, i+1), plt.imshow(img, interpolation='none'), plt.title(titles[i])
    plt.subplot(3, 4, i+5), plt.hist(img.ravel(), bins=50)
    plt.subplot(3, 4, i+9), plt.imshow(120 < img, interpolation='none')
```



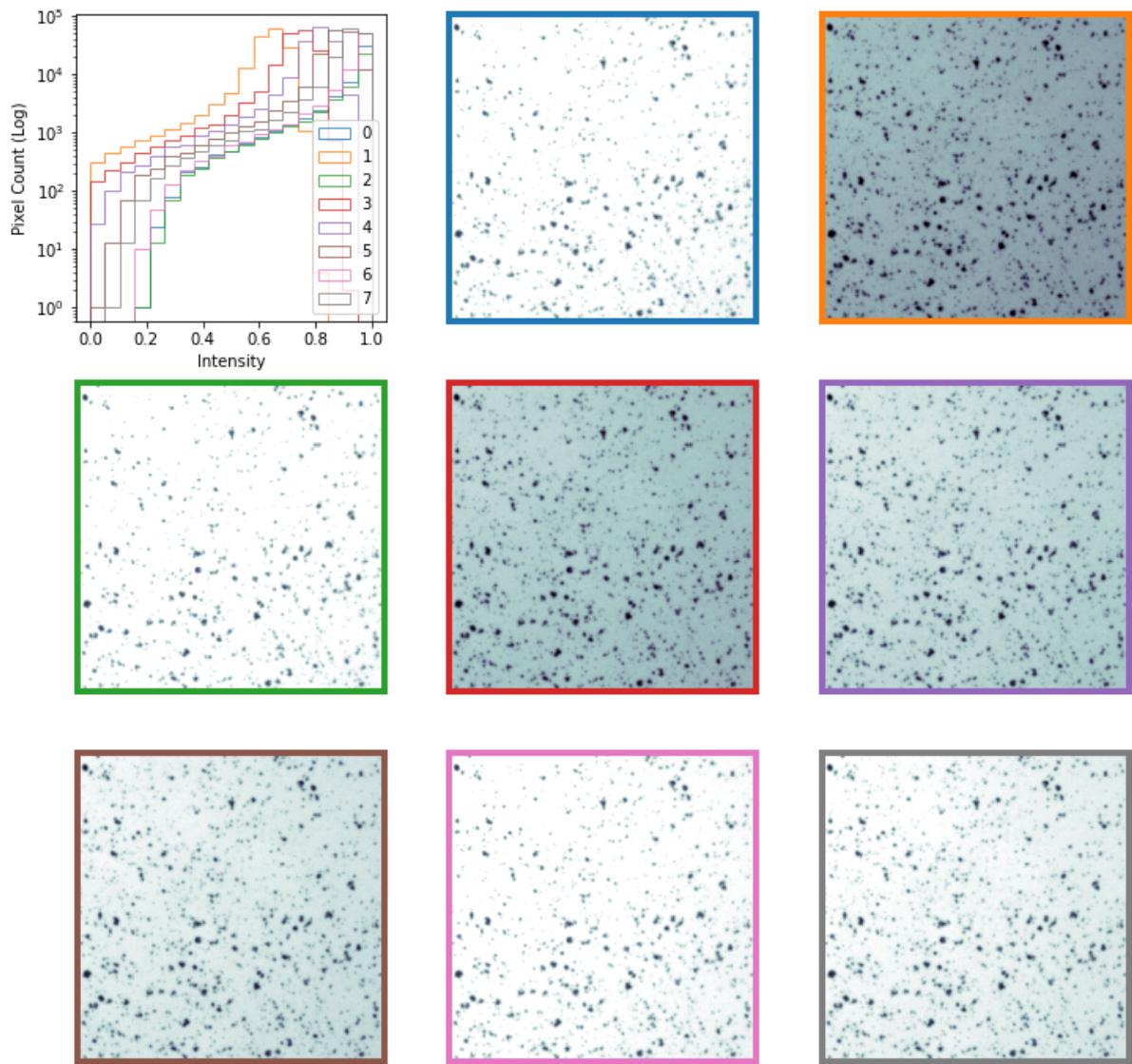
### 0.2.2 Inconsistent illumination on real data

Under realistic experiment conditions it may happen that the illumination fluctuates. This results in inconsistently illuminated images.

In this particular example of a cell colony, there is a random bias added to the images.

```
cell_img = imread("figures/Cell_Colony.jpg")/255.0
np.random.seed(2018)
m_cell_imgs = [cell_img+k for k in np.random.uniform(-0.25, 0.25, size = 8)]
fig, m_axs = plt.subplots(3, 3, figsize = (12, 12), dpi = 72)
ax1 = m_axs[0,0]
for i, (c_ax, c_img) in enumerate(zip(m_axs.flatten()[1:], m_cell_imgs)):
    ax1.hist(c_img.ravel(), np.linspace(0, 1, 20),
             label = '{}'.format(i), alpha = 0.75, histtype = 'step')
    c_ax.imshow(c_img, cmap = 'bone', vmin = 0, vmax = 1)
    #   c_ax.axis('off')
    c_ax.set_xticks([]);
    c_ax.set_yticks([]);
    for axis in ['top', 'bottom', 'left', 'right']:
        c_ax.spines[axis].set_linewidth(4)
        c_ax.spines[axis].set_color(colors[i])

    ax1.set_yscale('log', nonpositive = 'clip')
    ax1.set_ylabel('Pixel Count (Log)')
    ax1.set_xlabel('Intensity')
    ax1.legend();
```



### 0.2.3 Apply a threshold

When we apply a constant threshold of 0.6 to the images with different illuminations we just looked at you will see that there are great differences in how the cells are represented.

```
fig, m_axs = plt.subplots(3, len(m_cell_imgs),
                        figsize = (4*len(m_cell_imgs), 5*3),
                        dpi = 72)
THRESH_VAL = 0.6

for i, ((ax0, ax1, ax2), c_img) in enumerate(zip(m_axs.transpose(), m_cell_imgs)):
    ax0.imshow(c_img, cmap = 'bone', vmin = 0, vmax = 1)
    ax0.axis('off')
    ax1.hist(c_img.ravel()[(c_img.ravel()>THRESH_VAL)],
             np.linspace(0, 1, 100),
             label = 'Above Threshold',
```

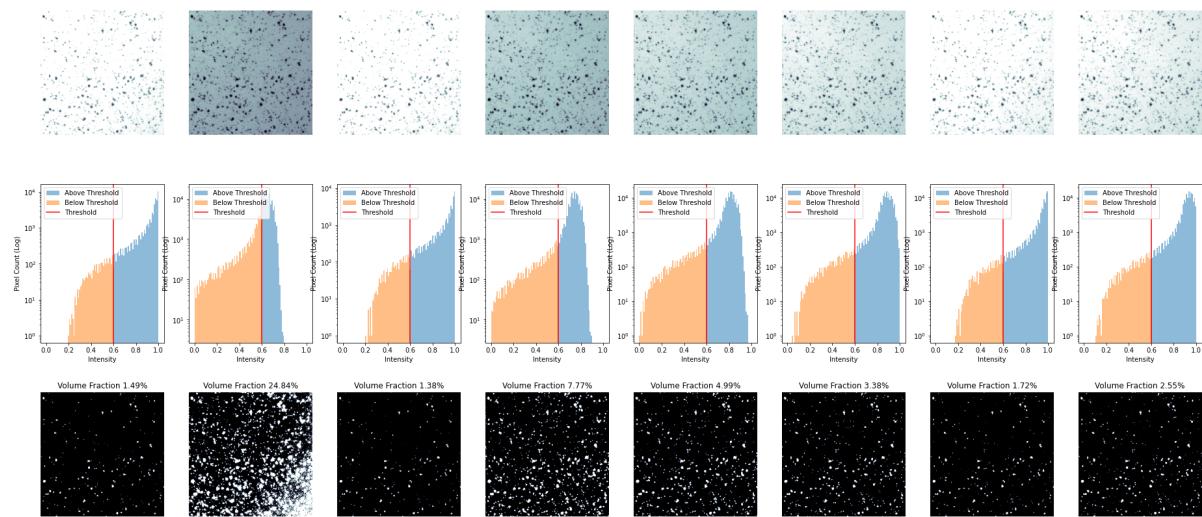
(continues on next page)

(continued from previous page)

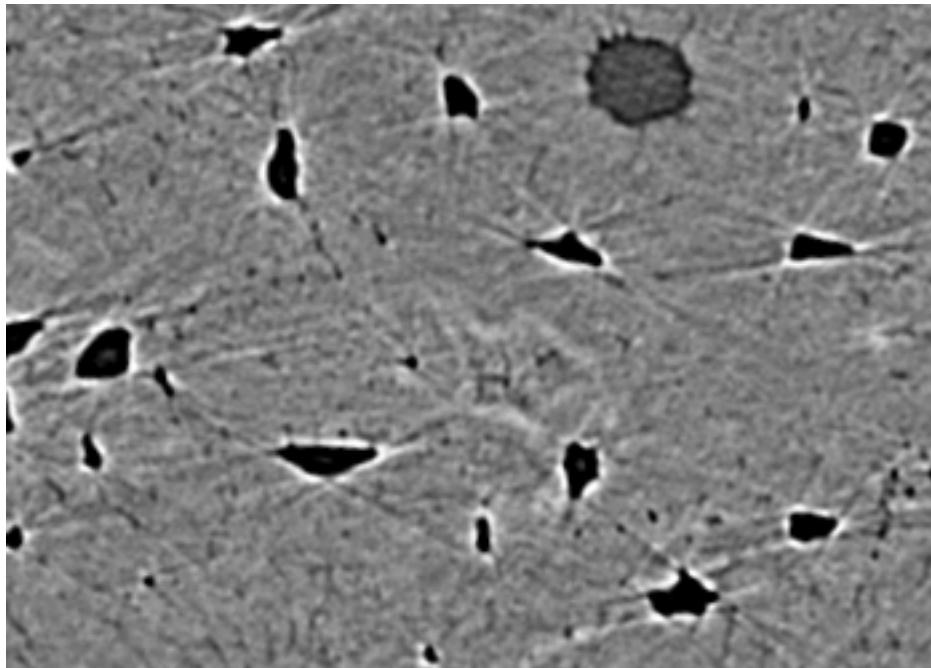
```

alpha = 0.5)
ax1.hist(c_img.ravel()[c_img.ravel()<THRESH_VAL],
         np.linspace(0, 1, 100),
         label = 'Below Threshold',
         alpha = 0.5)
ax1.axvline(THRESH_VAL,color='r',label='Threshold');
ax1.set_yscale('log', nonpositive = 'clip'); ax1.set_ylabel('Pixel Count (Log)');
ax1.set_xlabel('Intensity'); ax1.legend()
ax2.imshow(c_img<THRESH_VAL, cmap = 'bone', vmin = 0, vmax = 1)
ax2.set_title('Volume Fraction {0:2.2f}%.format(100*np.mean(c_img.ravel()<THRESH_
VAL))), ax2.axis('off'));

```



### 0.2.4 Where segmentation fails: Canaliculi



### Here is a bone slice

1. Find the larger cellular structures (osteocyte lacunae)
  2. Find the small channels which connect them together
- 

The first task is easy using a threshold and size criteria (we know how big the cells should be)

The second is much more difficult because the small channels having radii on the same order of the pixel size are obscured by

- partial volume effects
- and noise.

### 0.2.5 Where segmentation fails: Brain Cortex

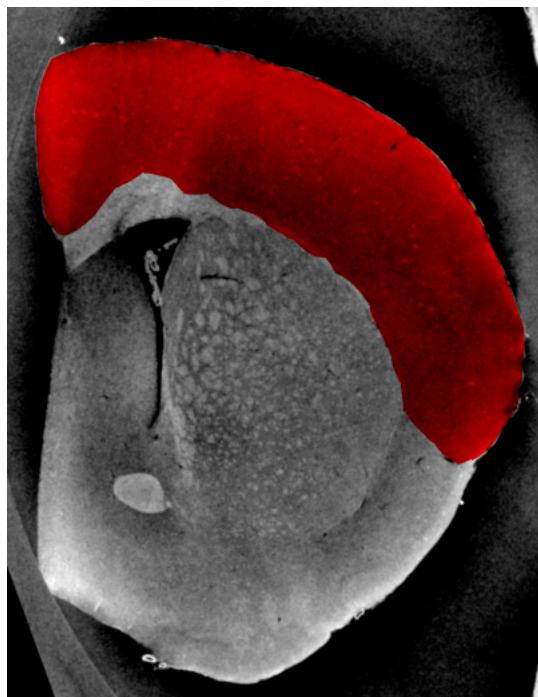


Fig. 2: Brain image with masked cortex.

- The cortex is barely visible to the human eye
- Tiny structures hint at where cortex is located

### Our problem

- A simple threshold is insufficient to finding the cortical structures
- Other filtering techniques are unlikely to magically fix this problem

## 0.3 Apply thresholds

```
%matplotlib inline
from skimage.io import imread
import matplotlib.pyplot as plt
import numpy as np

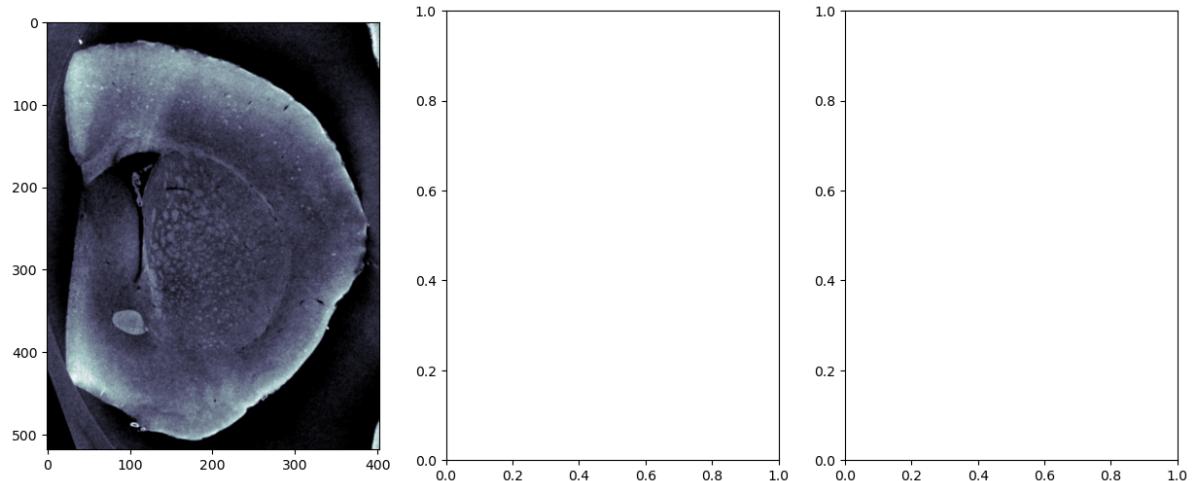
cortex_img = imread("figures/cortex.png")
np.random.seed(2018)
fig, (ax1, ax2, ax3) = plt.subplots(1, 3,
                                    figsize = (15, 6))
ax1.imshow(cortex_img, cmap = 'bone')

clst = np.array([[64,64,64],
                 [73,162,199],
                 [234,80,162],
                 [241,254,83]])/255.0
cmap = colors.ListedColormap(clst)

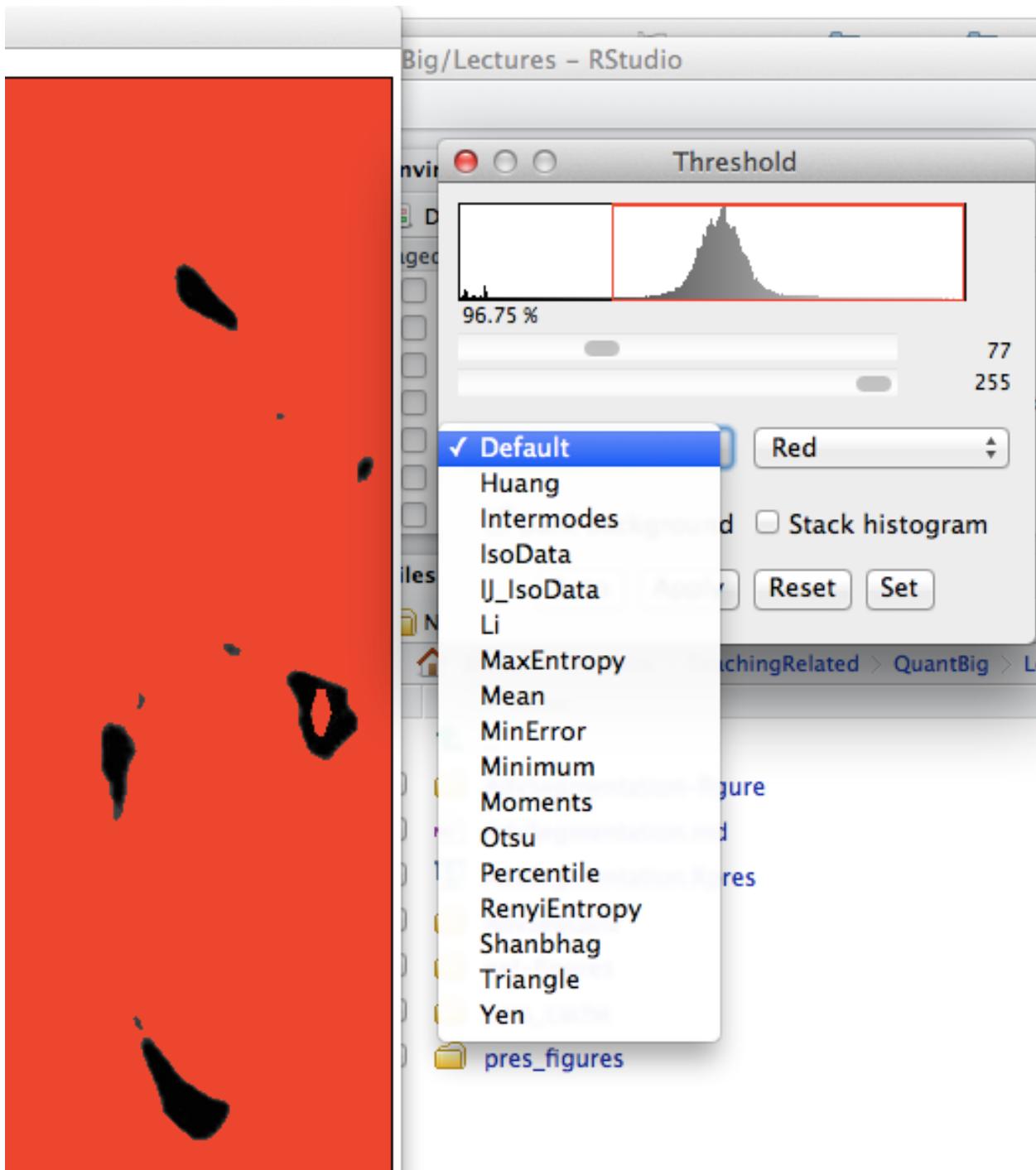
thresh_vals = np.linspace(cortex_img.min(), cortex_img.max(), 4+2)[:-1]
out_img = np.zeros_like(cortex_img)
for i, (t_start, t_end) in enumerate(zip(thresh_vals, thresh_vals[1:])):
    thresh_reg = (cortex_img>t_start) & (cortex_img<t_end)
    ax2.hist(cortex_img.ravel()[thresh_reg.ravel()],color=clst[i])
    out_img[thresh_reg] = i
ax3.imshow(out_img, cmap = cmap,interpolation='none');
```

```
-----
AttributeError                                                 Traceback (most recent call last)
Cell In[6], line 11
      5     ax1.imshow(cortex_img, cmap = 'bone')
      6     clst = np.array([[64,64,64],
      7                     [73,162,199],
      8                     [234,80,162],
      9                     [241,254,83]])/255.0
--> 11     cmap = colors.ListedColormap(clst)
     12     thresh_vals = np.linspace(cortex_img.min(), cortex_img.max(), 4+2)[:-1]
     13     out_img = np.zeros_like(cortex_img)

AttributeError: 'list' object has no attribute 'ListedColormap'
```



### 0.3.1 Automated Threshold Selection



Given that applying a threshold is such a common and significant step, there have been many tools developed to automatically (unsupervised) perform it. A particularly important step in setups where images are rarely consistent such as outdoor imaging which has varying lighting (sun, clouds). The methods are based on several basic principles.

### 0.3.2 Automated Methods

#### Histogram-based methods

Just like we visually inspect a histogram an algorithm can examine the histogram and find local minimums between two peaks, maximum / minimum entropy and other factors

- Otsu, Isodata, Intermodes, etc

#### Image based Methods

These look at the statistics of the threshold image themselves (like entropy) to estimate the threshold

#### Results-based Methods

These search for a threshold which delivers the desired results in the final objects. For example if you know you have an image of cells and each cell is between 200-10000 pixels the algorithm runs thresholds until the objects are of the desired size

- More specific requirements need to be implemented manually

### 0.3.3 Histogram-based Methods

Taking a typical image of a bone slice, we can examine the variations in calcification density in the image

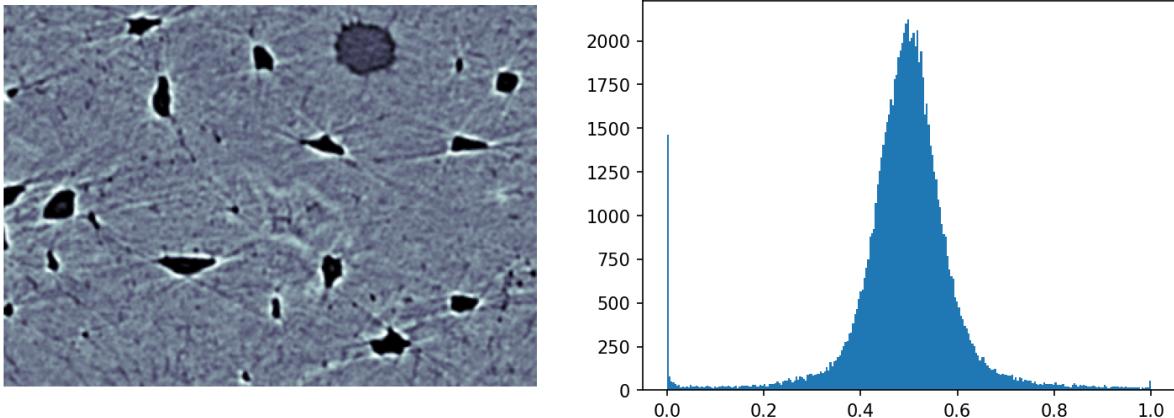
---

We can see in the histogram that there are two peaks, one at 0 (no absorption / air) and one at 0.5 (stronger absorption / bone)

```
%matplotlib inline
from skimage.io import imread
import matplotlib.pyplot as plt
import numpy as np

bone_img = imread("figures/bonegfilslice.png")/255.0
np.random.seed(2018)

fig, (ax1, ax2) = plt.subplots(1, 2, figsize = (12, 4), dpi = 150)
ax1.imshow(bone_img, cmap = 'bone')
ax1.axis('off')
ax2.hist(bone_img.ravel(), bins=256);
```



### 0.3.4 Histogram-Methods

---

#### Intermodes

- Take the point between the two modes (peaks) in the histogram

#### Otsu

Search and minimize intra-class (within) variance  $\sigma_w^2(t) = \omega_{bg}(t)\sigma_{bg}^2(t) + \omega_{fg}(t)\sigma_{fg}^2(t)$

#### Isodata

- Initialize with:  $\text{thresh} = \frac{\max(img) + \min(img)}{2}$
- *while* the thresh is changing
- $bg = img < \text{thresh}$ ,  $obj = img > \text{thresh}$
- Update  $\text{thresh} = (\text{avg}(bg) + \text{avg}(obj))/2$

### 0.3.5 Try All Thresholds

- opencv2
- scikit-image

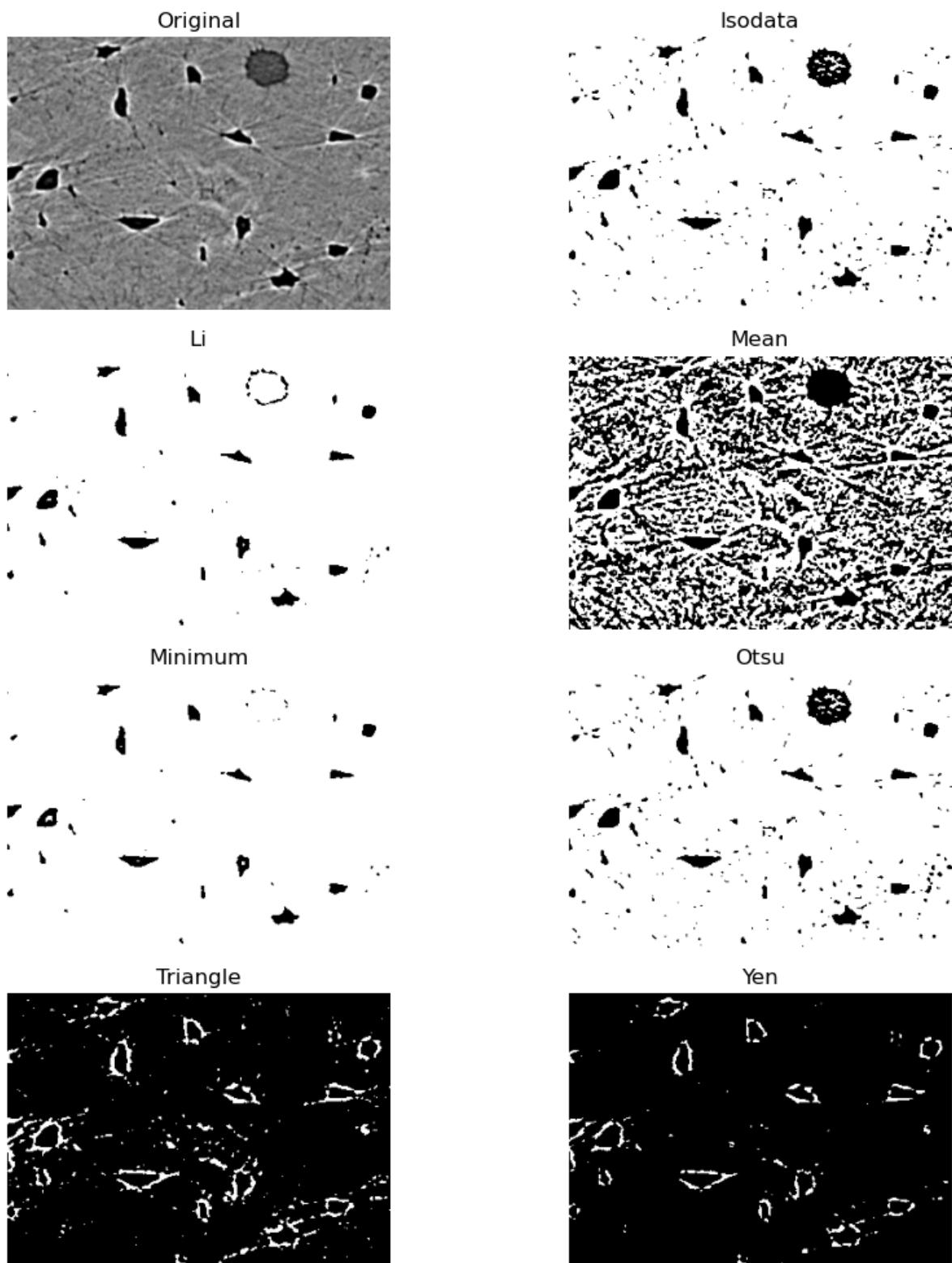
There are many methods and they can be complicated to implement yourself.

Both Fiji and scikit-image offers many of them as built in functions so you can automatically try all of them on your image.

```
bone_img = imread("figures/bonegfilslice.png")/255.0

from skimage.filters import try_all_threshold

fig, ax = try_all_threshold(bone_img, figsize=(10, 10), verbose=False);
```



### 0.3.6 Pitfalls of different thresholding methods

While an incredibly useful tool, there are many potential pitfalls to these automated techniques.

#### Histogram-based

These methods are very sensitive to the distribution of pixels in your image and may work really well on images with equal amounts of each phase but work horribly on images which have very high amounts of one phase compared to the others

#### Image-based

These methods are sensitive to noise and a large noise content in the image can change statistics like entropy significantly.

#### Results-based

These methods are inherently biased by the expectations you have.

- If you want to find objects between 200 and 1000 pixels - you will!
- they just might not be anything meaningful.

### 0.3.7 Realistic Approaches for Dealing with these Shortcomings

Imaging science rarely represents the ideal world and will never be 100% perfect. At some point we need to write our master's thesis, defend, or publish a paper. These are approaches for more qualitative assessment we will later cover how to do this a bit more robustly with quantitative approaches

#### Model-based

One approach is to

- try and simulate everything (including noise) as well as possible
- and to apply these techniques to many realizations of the same image
- and qualitatively keep track of how many of the results accurately identify your phase or not.

**Hint:** >95% seems to convince most biologists

#### Sample-based

- Apply the methods to each sample and keep track of which threshold was used for each one.
- Go back and apply each threshold to each sample in the image
- and keep track of how many of them are correct enough to be used for further study.

## Worst-case Scenario

Come up with the worst-case scenario (noise, misalignment, etc) and assess how unacceptable the results are. Then try to estimate the quartiles range (75% - 25% of images).

### 0.3.8 Hysteresis Thresholding

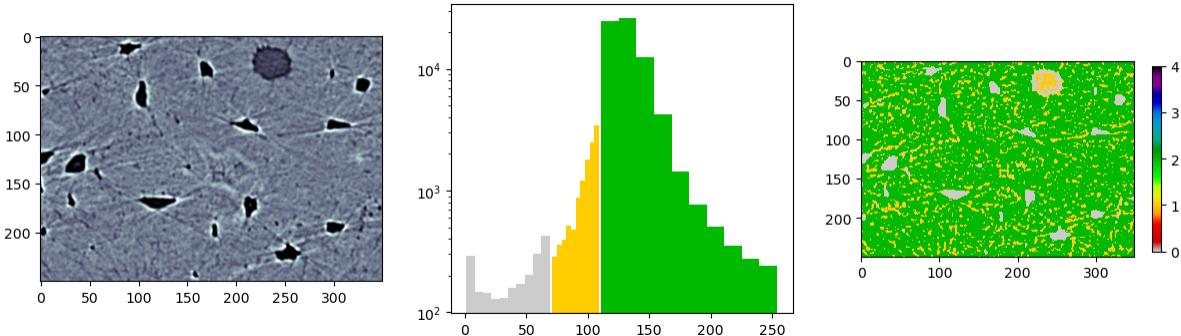
For some images a single threshold does not work

- large structures are very clearly defined
- smaller structures are difficult to differentiate (see partial volume effect)

ImageJ Source

```
%matplotlib inline
from skimage.io import imread
import matplotlib.pyplot as plt
import numpy as np

bone_img = imread("figures/bonegfilslice.png")
np.random.seed(2018)
fig, (ax1, ax2, ax3) = plt.subplots(1, 3,
                                    figsize = (15, 4))
cmap = plt.cm.nipy_spectral_r
ax1.imshow(bone_img, cmap = 'bone')
thresh_vals = [0, 70, 110, 255]
out_img = np.zeros_like(bone_img)
for i, (t_start, t_end) in enumerate(zip(thresh_vals, thresh_vals[1:])):
    thresh_reg = (bone_img>t_start) & (bone_img<t_end)
    ax2.hist(bone_img.ravel()[thresh_reg.ravel()], color = cmap(i/len(thresh_vals)))
    out_img[thresh_reg] = i
ax2.set_yscale("log", nonpositive='clip')
th_ax = ax3.imshow(out_img, cmap = cmap, vmin = 0, vmax = len(thresh_vals),
                    interpolation='none')
plt.colorbar(th_ax, shrink=0.6);
```



### 0.3.9 Goldilocks Situation

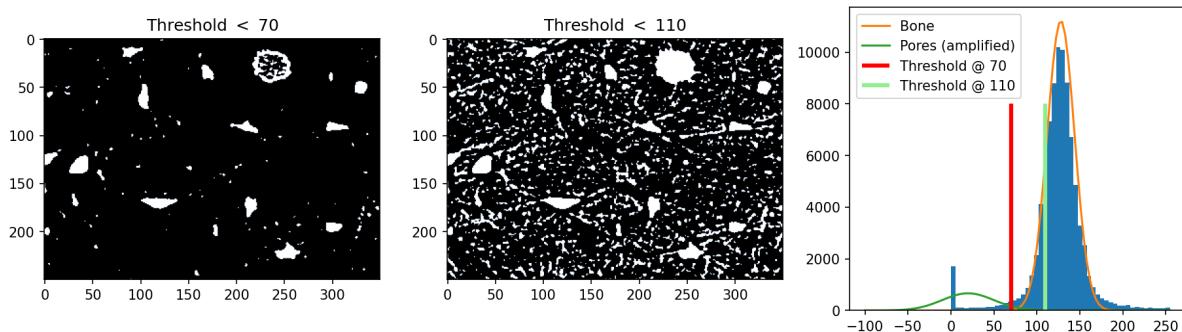
Here we end up with a goldilocks situation

- Mama bear and Papa Bear
- one is too low and the other is too high

The Goldilocks principle

```
fig, (ax1, ax2, ax3) = plt.subplots(1, 3, figsize = (15, 4), dpi=150)
ax1.imshow(bone_img

```



### Baby Bear

We can thus follow a process for ending up with a happy medium of the two

### Hysteresis Thresholding: Reducing Pixels

Now we apply the following steps.

1. Take the first threshold image with the highest (more strict) threshold
2. Remove the objects which are not cells (too small) using an opening operation.
3. Take a second threshold image with the higher value
4. Combine both threshold images
5. Keep the *between* pixels which are connected (by looking again at a neighborhood  $\mathcal{N}$ ) to the *air* voxels and ignore the other ones. This goes back to our original supposition that the smaller structures are connected to the larger structures

```
%matplotlib inline
from skimage.morphology import dilation, opening, disk
from collections import OrderedDict
```

### Thresholding with hysteresis

**Step 1** Apply several thresholds to the image

```
thresh_vals = [0, 70, 110, 255]
step_list = OrderedDict()
step_list['Strict Threshold'] = bone_img

```

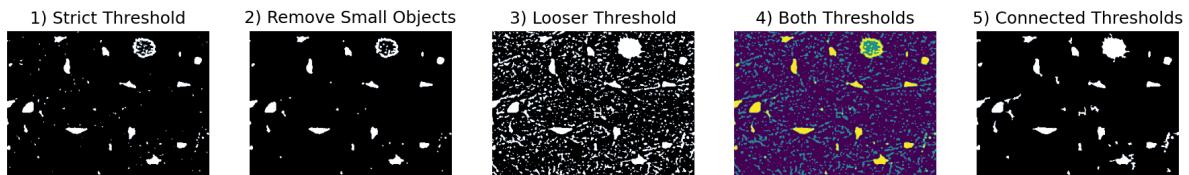
**Step 2** Combine the different images  $I_{ConnectedThresholds} = \delta(I_{NoSmallObjects}) \cap I_{LooseThresh}$

```
# the tricky part keeping the between images
step_list['Connected Thresholds'] = step_list['Remove Small Objects']

for i in range(10):
    step_list['Connected Thresholds'] = dilation(step_list['Connected Thresholds'] ,
                                                disk(1.8)) & step_list['Looser_
    Threshold']

fig, ax_steps = plt.subplots(1, len(step_list), figsize = (15, 5), dpi = 150)

for i, (c_ax, (c_title, c_img)) in enumerate(zip(ax_steps.flatten(), step_list.
    items()), 1):
    c_ax.imshow(c_img, cmap = 'bone' if c_img.max()<=1 else 'viridis')
    c_ax.set_title('%d %s' % (i, c_title)); c_ax.axis('off');
```



### 0.3.10 Multiple thresholds

It is not uncommon to have multiple classes in an image

Setting thresholds at 0.6 and 2.2 makes sense...

We see some misclassification that resembles a skin on the object!

### Investigation the virtual skin effect

Let's look at three different edge profiles in the image

The virtual skin appears because the great transition from class 1 to 3 passes over several pixels and some of these pixels have values in the interval of class 2. These will be misclassified as class 2 and appear as the virtual skin in the segmented image.

### How to avoid segmentation problems at the edges

Often the edge segmentation only involves a few pixels that can be handled in different ways. E.g.

- Hysteresis
- Guarded edges and region growing

The guarded edge method uses an edge enhancing filter like a Laplacian filter to create an edge image. It is not necessary to handle all edge but only those with high amplitudes in the Laplacian image. These pixels will be excluded from the thresholding initially. Once the obvious pixels have been assigned we can handle the edge using region growing to fill up the missing edge information. we could also use a Laplacian or Gaussian filter to increase the robustness for low SNR.

## 0.4 More Complicated Images

As we briefly covered last time, many measurement techniques produce quite rich data.

- Digital cameras produce 3 channels of color for each pixel (rather than just one intensity)
- MRI produces dozens of pieces of information for every voxel which are used when examining different *contrasts* in the system.
- Raman-shift imaging produces an entire spectrum for each pixel
- Coherent diffraction techniques produce 2- (sometimes 3) diffraction patterns for each point.  $I(x, y) = \hat{f}(x, y)$

### 0.4.1 Feature Vectors

A pairing between spatial information (position) and some other kind of information (value).  $\vec{x} \rightarrow \vec{f}$

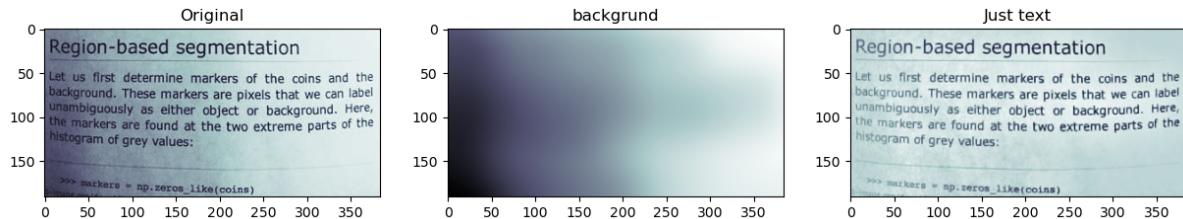
We are used to seeing images in a grid format where the position indicates the row and column in the grid and the intensity (absorption, reflection, tip deflection, etc) is shown as a different color. We take an example here of text on a page.

This is an alternative to the top-hat method last week.

```
from skimage.data import page
import skimage.morphology as morph

page_image = page()
background = 255*gaussian(page_image, 20.0)
#background = morph.closing(page_image, morph.disk(51))
just_text = median(page_image, morph.disk(1))-background
```

```
fig,ax = plt.subplots(1,3,figsize=[15,5])
ax[0].imshow(page_image, cmap = 'bone'); ax[0].set_title('Original');
ax[1].imshow(background, cmap = 'bone'); ax[1].set_title('background');
ax[2].imshow(just_text, cmap = 'bone'); ax[2].set_title('Just text');
```



The gradient in the original is removed using a combination of filters to reconstruct the illumination

$$\text{JustText} = \underbrace{\text{median}_{\text{disk}}(\text{img})}_{\text{Noise}} - \underbrace{G_{\sigma=20} * \text{img}}_{\text{Illumination}}$$

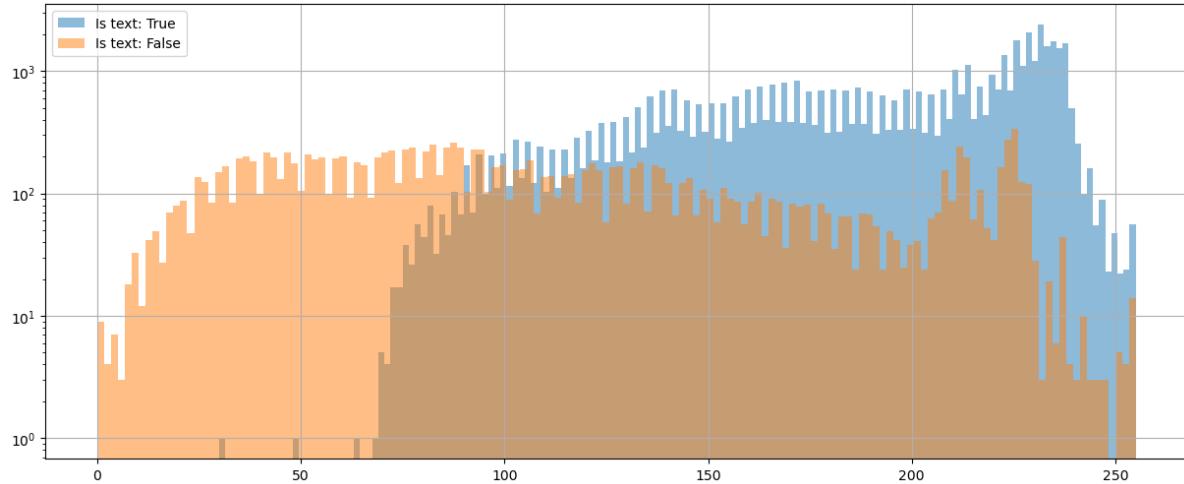
## 0.4.2 Let's create a feature table

```
xx, yy = np.meshgrid(np.arange(page_image.shape[1]),
                     np.arange(page_image.shape[0]))
page_table = pd.DataFrame(dict(x = xx.ravel(),
                                 y = yy.ravel(),
                                 intensity = page_image.ravel(),
                                 is_text = just_text.ravel()>0))
page_table.sample(10)
```

	x	y	intensity	is_text
38169	153	99	174	True
17741	77	46	161	True
38819	35	101	107	True
24130	322	62	234	True
15369	9	40	120	True
24015	207	62	120	False
6168	24	16	118	False
53301	309	138	227	True
8348	284	21	204	False
30094	142	78	181	True

## 0.4.3 Inspect the features - Intensity vs. IsText

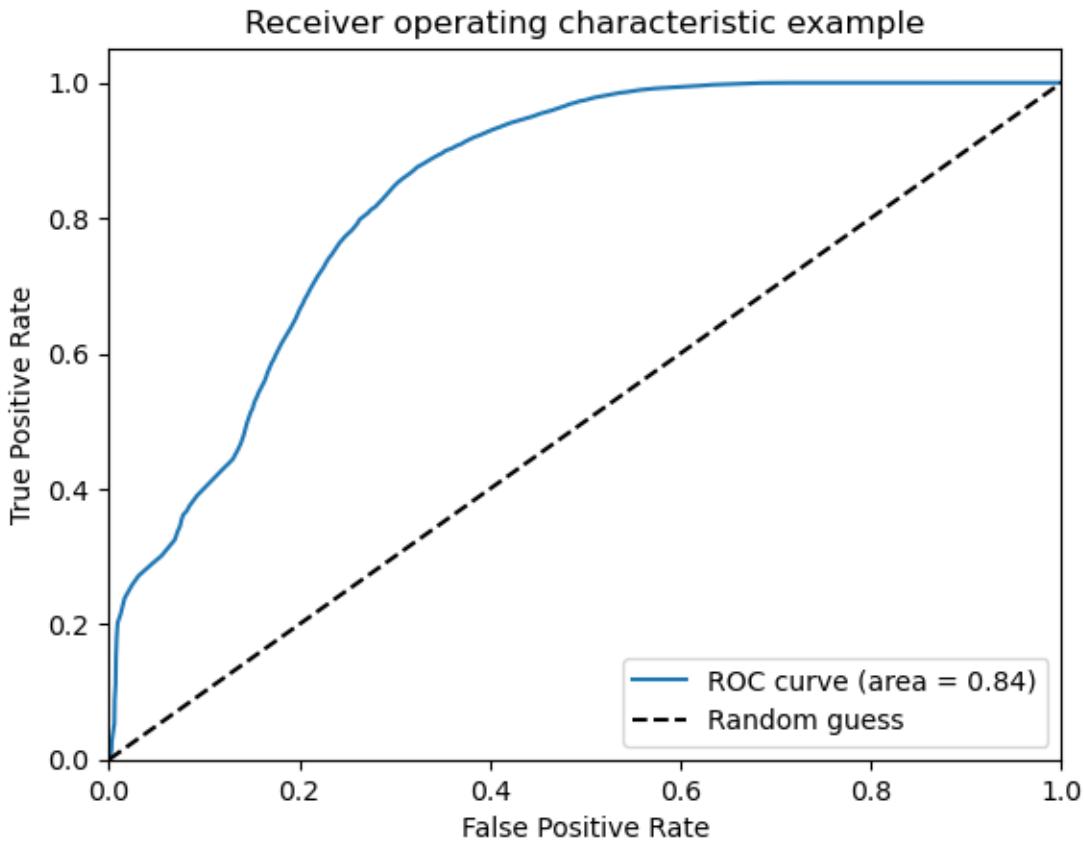
```
fig, ax1 = plt.subplots(1,1,figsize=(15,6))
for c_cat in [True,False] :
    page_table[page_table['is_text']==c_cat]['intensity'].hist(ax=ax1,
                                                               bins=150,
                                                               alpha=0.5,
                                                               label='Is text: {}'.
                                                               format(c_cat))
ax1.set_yscale("log", nonpositive='clip')
ax1.legend();
```



### 0.4.4 What does the ROC curve look like?

```
from sklearn.metrics import roc_curve, roc_auc_score
fpr, tpr, _ = roc_curve(page_table['is_text'], page_table['intensity'])
roc_auc      = roc_auc_score(page_table['is_text'], page_table['intensity'])
```

```
fig, ax = plt.subplots(1,1)
ax.plot(fpr, tpr, label='ROC curve (area = {:.2f})'.format(roc_auc))
ax.plot([0, 1], [0, 1], 'k--', label='Random guess')
ax.set_xlim([0.0, 1.0])
ax.set_ylim([0.0, 1.05])
ax.set_xlabel('False Positive Rate');
ax.set_ylabel('True Positive Rate')
ax.set_title('Receiver operating characteristic example');
ax.legend(loc="lower right");
```



## 0.5 Adding Information

Here we can improve the results by adding information.

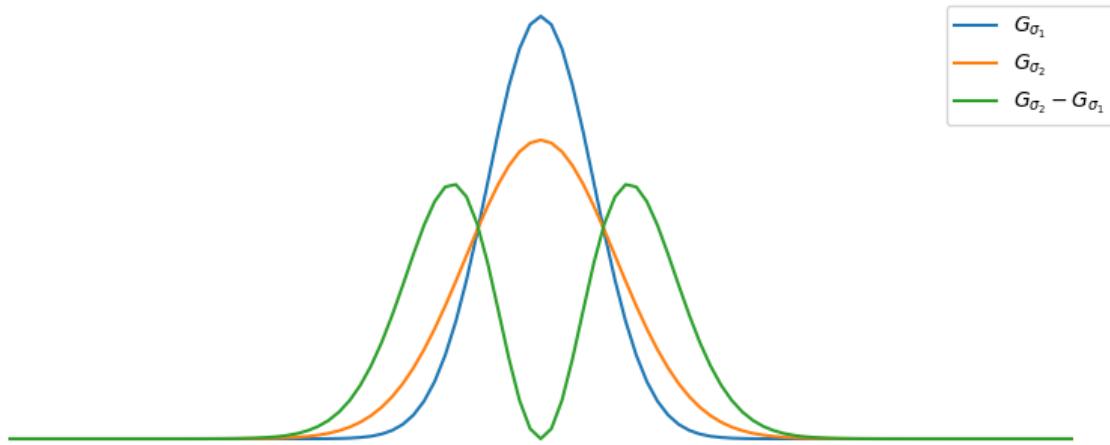
As we discussed in the third lecture (enhancement), edge-enhancing filters can be very useful for classifying images.

How about:

$$f_{DoG} = G_{\sigma_1} * f - G_{\sigma_2} * f = (G_{\sigma_1} - G_{\sigma_2}) * f, \quad \sigma_1 < \sigma_2$$

```
fig,ax=plt.subplots(1,figsize=(10,4))
x=np.linspace(-10,10,101)
g1=np.exp(-x**2/2)
g2=np.exp(-x**2/4)

ax.plot(x,g1/g1.sum(), label="$G_{\sigma_1}$")
ax.plot(x,g2/g2.sum(), label="$G_{\sigma_2}$")
ax.plot(x,(g2-g1)/(g2.sum()-g1.sum()),label="$G_{\sigma_2}-G_{\sigma_1}$")
ax.legend()
ax.axis("off");
```



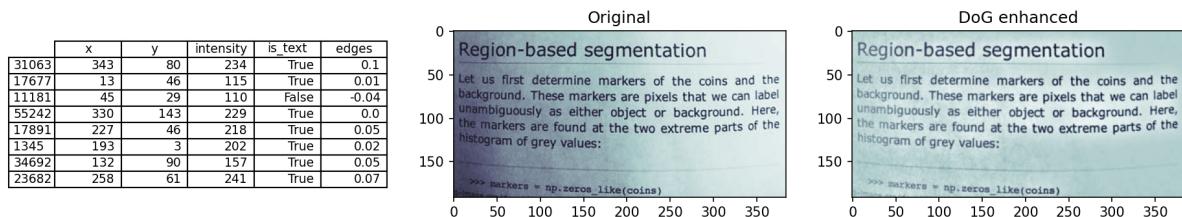
### 0.5.1 Testing with enhanced edges

```
def dog_filter(in_img, sig_1, sig_2):
    return gaussian(in_img, sig_1) - gaussian(in_img, sig_2)

page_edges = dog_filter(page_image, 0.5, 10)
page_table['edges'] = page_edges.ravel()
```

```
fig, ax = plt.subplots(1,3, figsize = (15, 4), dpi=150)
ax[1].imshow(page_image, cmap = 'bone'), ax[1].set_title('Original')
ax[2].imshow(page_edges, cmap = 'bone'), ax[2].set_title('DoG enhanced')

pd.plotting.table(data=page_table.sample(8).round(decimals=2), ax=ax[0], loc='center')
ax[0].axis('off');
```

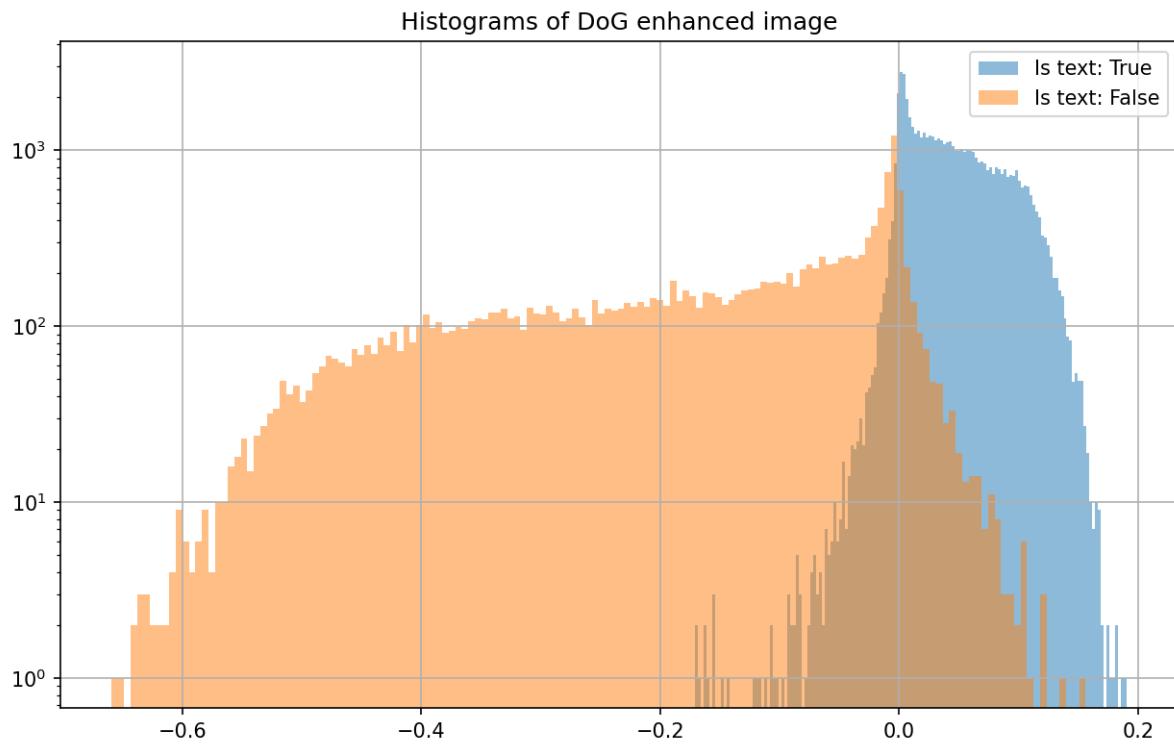


## 0.5.2 Histogram of enhanced image

Let's look at the gray level distribution after applying the enhancement filter.

```
fig, ax = plt.subplots(1, figsize=(10, 6), dpi=150)

for c_cat in [True, False] :
    page_table[page_table['is_text']==c_cat]['edges'].hist(ax=ax,
                                                          bins=150,
                                                          alpha=0.5,
                                                          label='Is text: {}'.format(c_cat))
ax.set_yscale("log", nonpositive='clip')
ax.legend();
ax.set_title('Histograms of DoG enhanced image'); ax.set_yscale("log", nonpositive=
    'clip'); ax.legend();
```



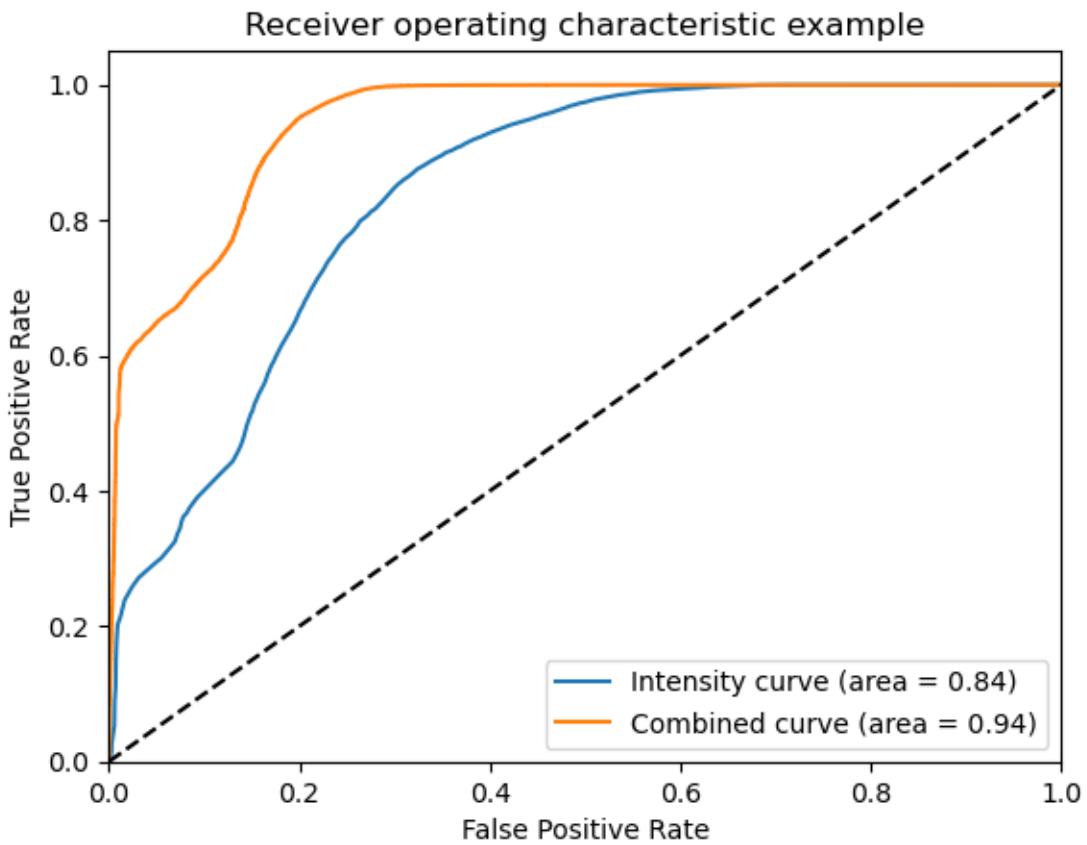
We see here that the text pixels have been compressed into a narrow interval with some tail toward the lower intensities.

## 0.5.3 Checking the ROC performance

Now, we can compute the ROC curve to see if the enhanced image performs better than the original.

```
from sklearn.metrics import roc_curve, roc_auc_score
fpr2, tpr2, _ = roc_curve(page_table['is_text'],
                           page_table['intensity']/1000.0+page_table['edges'])
roc_auc2      = roc_auc_score(page_table['is_text'],
                           page_table['intensity']/1000.0+page_table['edges'])
```

```
# Visualization
fig, ax = plt.subplots(1,1)
ax.plot(fpr, tpr, label='Intensity curve (area = %0.2f)' % roc_auc)
ax.plot(fpr2, tpr2, label='Combined curve (area = %0.2f)' % roc_auc2)
ax.plot([0, 1], [0, 1], 'k--')
ax.set_xlim([0.0, 1.0])
ax.set_ylim([0.0, 1.05])
ax.set_xlabel('False Positive Rate')
ax.set_ylabel('True Positive Rate')
ax.set_title('Receiver operating characteristic example')
ax.legend(loc="lower right");
```



We can see that it is doing better already by looking at the curves. Also the AUC confirms this which is great. We have made an improvement thanks to the preprocessing of the data.

### 0.5.4 Why does the second filter perform better?

We can see the reason for the performance improvement when we compare the histograms.

```
fig, (ax1, ax2) = plt.subplots(1,2, figsize = (15, 5))

for c_cat in [True, False] :
    page_table[page_table['is_text']==c_cat]['intensity'].hist(ax=ax1,
                                                               bins=150,
                                                               alpha=0.5,
```

(continues on next page)

(continued from previous page)

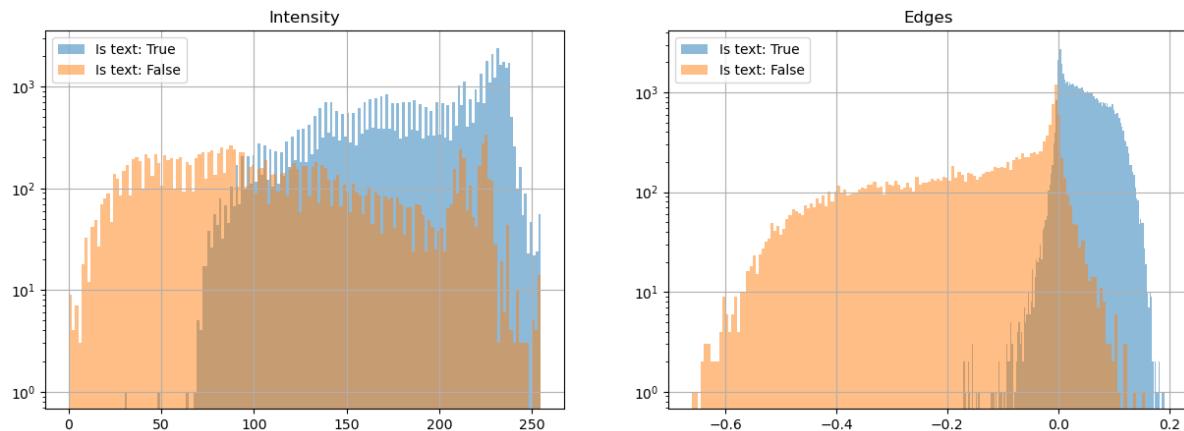
```

    ↵format(c_cat))
ax1.set_yscale("log", nonpositive='clip')
ax1.legend();
ax1.set_title('Intensity')

for c_cat in [True, False] :
    page_table[page_table['is_text']==c_cat]['edges'].hist(ax=ax2,
        bins=150,
        alpha=0.5,
        label='Is text: {0}'.format(c_cat))

ax2.set_yscale("log", nonpositive='clip')
ax2.legend();
ax2.set_title('Edges');

```



The two classes are more or less overlapping each other in the original image. After filtering the gray levels are more compact for each class with some slight overlap in the tail distributions.

## 0.6 Clustering / Classification (Unsupervised)

Unsupervised segmentation method tries to make sense of the that without any prior knowledge. They may need an initialization parameter telling how many classes are expected in the data, but beyond that you don't have to provide much more information.

- Automatic clustering of multidimensional data into groups based on a distance metric
- Fast and scalable to petabytes of data (Google, Facebook, Twitter, etc. use it regularly to classify customers, advertisements, queries)
- **Input** = feature vectors, distance metric, number of groups
- **Output** = a classification for each feature vector to a group

### 0.6.1 Example data

With clustering methods you aim to group data points together into a limited number of clusters. Here, we start to look at an example where each data point has two values.

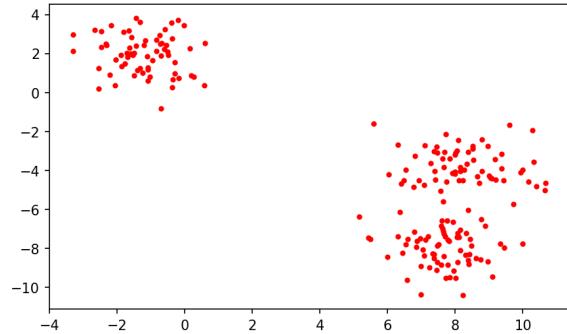
The test data is generated using the `make_blobs` function.

```
test_pts = pd.DataFrame(make_blobs(n_samples=200, random_state=2018) [
    0], columns=['x', 'y'])

fig, ax = plt.subplots(1,2, figsize = (15, 4), dpi=150)

ax[1].plot(test_pts.x, test_pts.y, 'r.')
pd.plotting.table(data=test_pts.sample(8).round(decimals=2), ax=ax[0], loc='center')
ax[0].axis('off');
```

	x	y
190	7.5	-7.8
96	8.43	-7.5
135	8.88	-6.84
73	5.16	-6.36
81	10.32	-3.54
22	8.75	-8.56
12	-2.22	0.91
119	-1.87	1.36



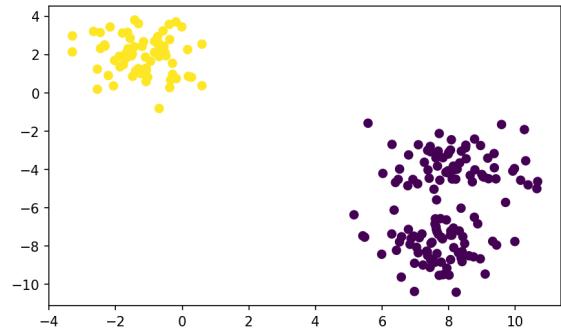
### First clustering attempt

The generated data set has two obvious clusters, but if look closer it is even possible to identify three clusters. We will now use this data to try the k-means algorithm.

```
km = KMeans(n_clusters=2, random_state=2018, n_init='auto')
n_grp = km.fit_predict(test_pts)
grp_pts = test_pts.copy()
grp_pts['group'] = n_grp
grp_pts.groupby(['group']).apply(lambda x: x.sample(5))
```

```
# Plotting
fig, ax = plt.subplots(1,2, figsize = (15, 4), dpi=150)
ax[1].scatter(test_pts.x, test_pts.y, c=n_grp)
pd.plotting.table(data=grp_pts.groupby(['group']).apply(lambda x: x.sample(5)), ax=ax[0], loc='center')
ax[0].axis('off');
```

	x	y	group
(0, 137)	7.936151356509987	-6.630712020263951	0.0
(0, 172)	8.062238824854576	-8.70500225647179	0.0
(0, 125)	7.370485145512017	-8.489127187418553	0.0
(0, 105)	6.295207189761015	-7.366814493077998	0.0
(0, 195)	7.720958193819376	-2.1207408297319548	0.0
(1, 72)	-1.0861012238651098	1.297087486519242	1.0
(1, 129)	-2.1666153280332363	3.444121311992933	1.0
(1, 197)	-2.0652721320971406	0.3682746189614694	1.0
(1, 118)	-1.07745532002573654	1.1922372983022362	1.0
(1, 198)	-0.012441330856636457	3.4512483449860842	1.0



## 0.6.2 K-Means Algorithm

We give as an initial parameter

- the number of groups we want to find
- and possibly a criteria for removing groups that are too similar

It is an iterative method that starts with a label image where each pixel has a random label assignment. Each iteration involves the following steps:

1. Compute current centroids based on the o current labels
2. Compute the value distance for each pixel to the each centroid. Select the class which is closest.
3. Reassign the class value in the label image.
4. Repeat until no pixels are updated.

The distance from pixel  $i$  to centroid  $j$  is usually computed as  $\|p_i - c_j\|_2$ .

### Increasing the number of clusters

In this example we will use the blob data we previously generated and to see how k-means behave when we select different numbers of clusters.

Let's try k-means with N=2,3,4 clusters:

```
clusters=[2, 3, 4, 5]
fig, axes = plt.subplots(1, len(clusters), figsize=(15, 4))

for ax,N in zip(axes,clusters) :
    km      = KMeans(n_clusters=N, random_state=2018,n_init='auto');
    n_grp = km.fit_predict(test_pts)

    ax.scatter(test_pts.x, test_pts.y, c=n_grp)
    ax.set_title('{0} groups'.format(N))
```

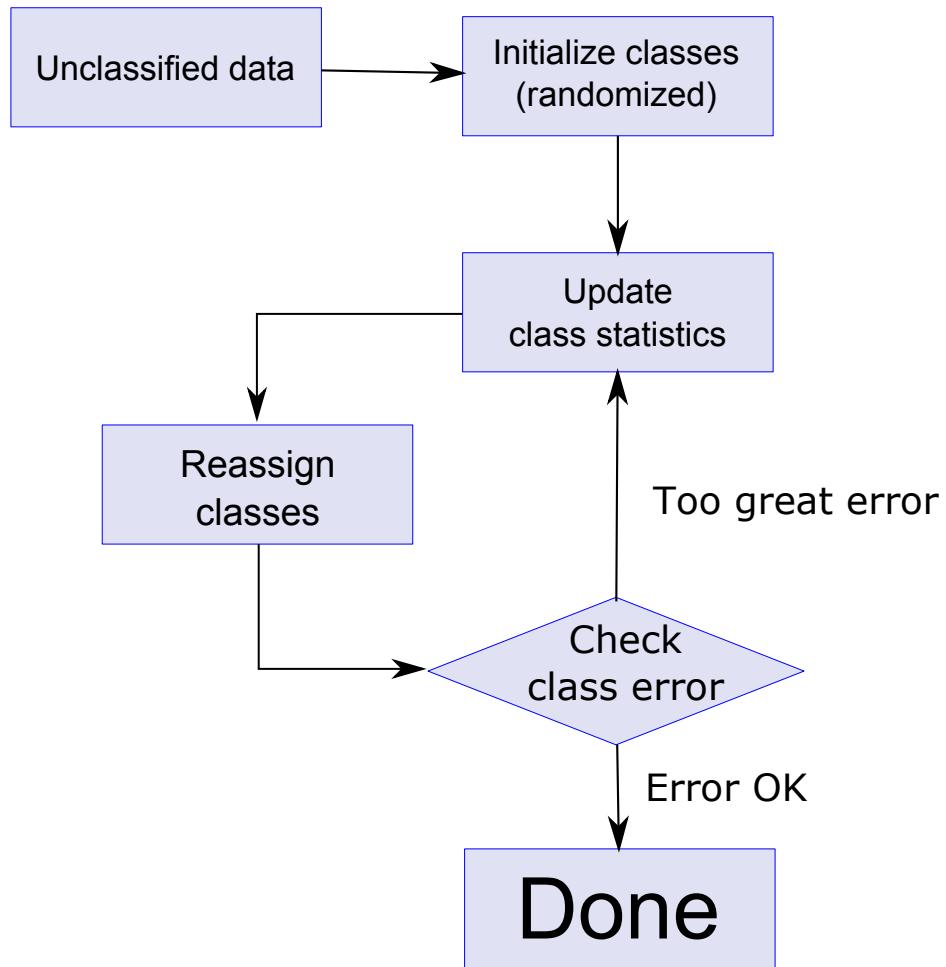
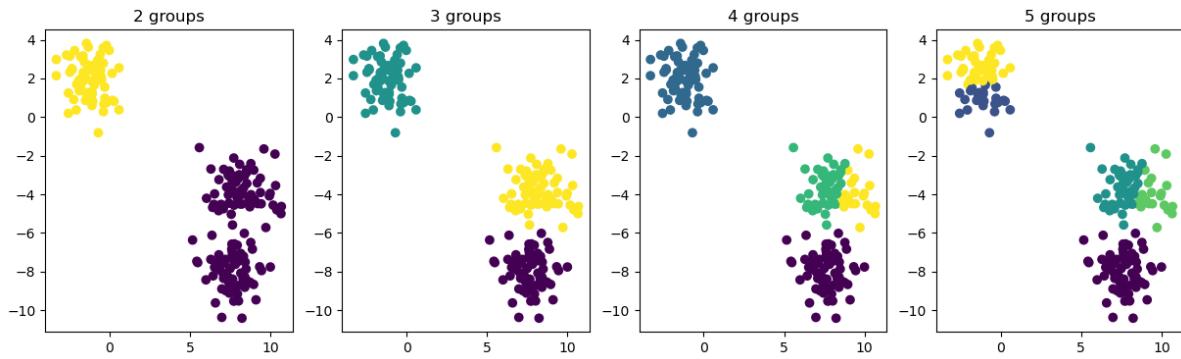


Fig. 1: Flow chart for the k-means clustering iterations.



**Note:** If you look for N groups you will always find N groups with K-Means, whether or not they make any sense

When we select two clusters there is a natural separation between the two clusters we easily spotted by just looking at the data. When the number of clusters is increased to three, we again see a cluster separation that makes sense. Now, when the number of clusters is increased yet another time we see that one of the clusters is split once more. This time it is however questionable if the number of clusters makes sense. From this example, we see that it is important to be aware of problems related to over segmentation.

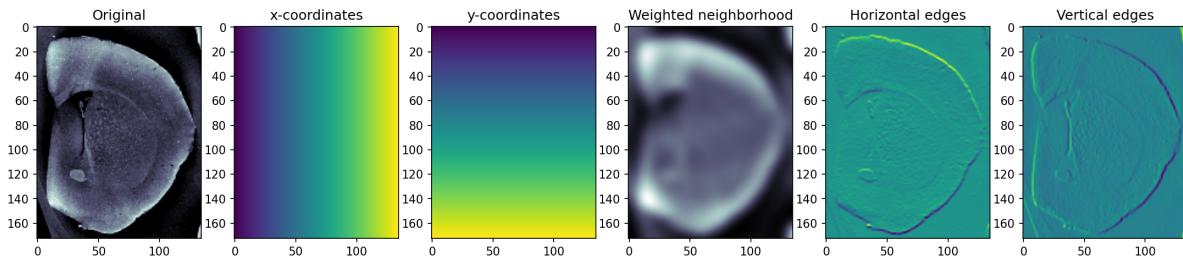
### 0.6.3 What vector space do we have?

- Sometimes represent physical locations (classify swiss people into cities)
- Can include intensity or color (K-means can be used as a thresholding technique when you give it image intensity as the vector and tell it to find two or more groups)
- Can also include orientation, shape, or in extreme cases full spectra (chemically sensitive imaging)

### 0.6.4 Add spatial information to k-means

It is important to note that k-means by definition is not position sensitive. Clustering is by definition not position sensitive, mainly measures distances between values and distributions. The position can however be included as additional components in the data vectors. You can also add neighborhood information using filtered images as additional components of the feature vectors.

```
cortex_img = imread("figures/cortex.png")[:, :, ::3]/1000.0
fig, ax = plt.subplots(1, 6, figsize=(18, 5), dpi=150);
xx, yy = np.meshgrid(np.arange(cortex_img.shape[1]),
                      np.arange(cortex_img.shape[0]))
ax[0].imshow(cortex_img, cmap='bone'), ax[0].set_title('Original')
ax[1].imshow(xx), ax[1].set_title('x-coordinates')
ax[2].imshow(yy), ax[2].set_title('y-coordinates')
ax[3].imshow(flt.gaussian(cortex_img, sigma=5), cmap='bone'), ax[3].set_title(
    'Weighted neighborhood')
ax[4].imshow(flt.sobel_h(cortex_img)), ax[4].set_title('Horizontal edges')
ax[5].imshow(flt.sobel_v(cortex_img)), ax[5].set_title('Vertical edges');
```



### K-Means Applied to Cortex Image

In this example we use position and intensity as feature vectors.

```

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from skimage.io import imread
%matplotlib inline

cortex_img = imread("figures/cortex.png")[:, :, ::3]/1000.0
np.random.seed(2018)
xx, yy = np.meshgrid(np.arange(cortex_img.shape[1]), np.arange(cortex_img.shape[0]))

cortex_df = pd.DataFrame(dict(x=xx.ravel(),
                               y=yy.ravel(),
                               intensity=cortex_img.ravel()))

```

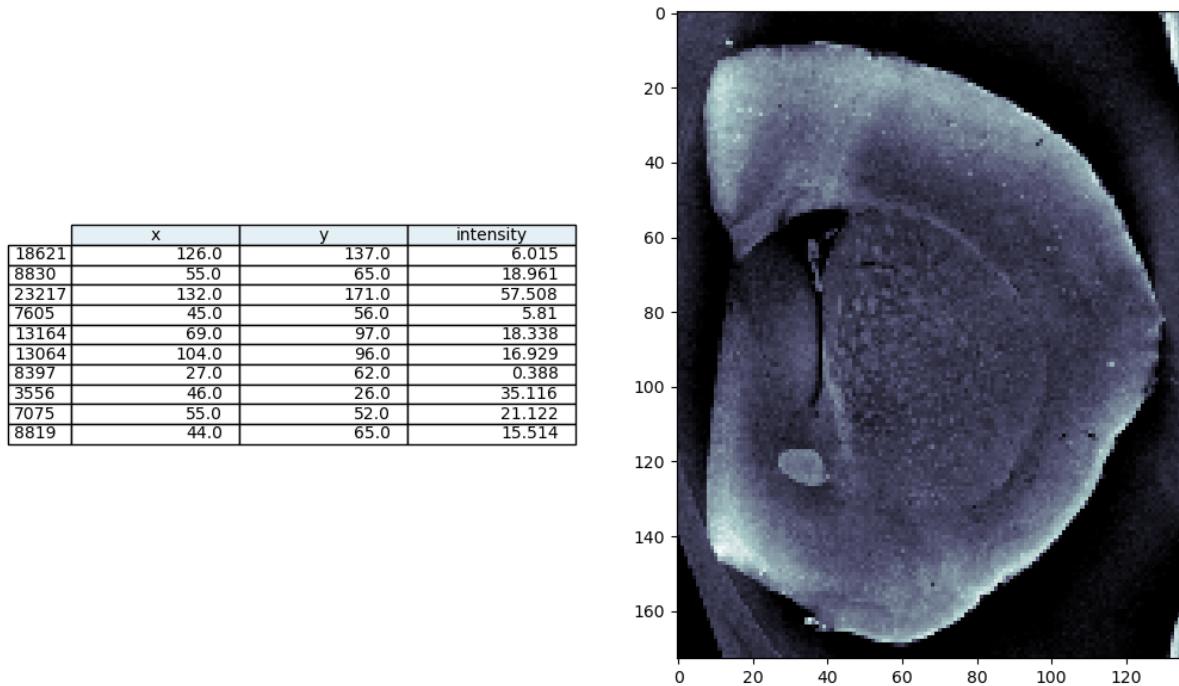
```

fig, (ax0, ax1) = plt.subplots(1, 2, figsize=(12, 8))
ax1.imshow(cortex_img, cmap='bone')

ccolors = plt.cm.BuPu(np.full(3, 0.1))

pd.plotting.table(data=cortex_df.sample(10), ax=ax0, loc='center', colColours=ccolors,
                   fontsize=20)
ax0.axis('off');

```



### First segmentation attempt with k-means on brain image

We use  $N_{clusters} = 4$

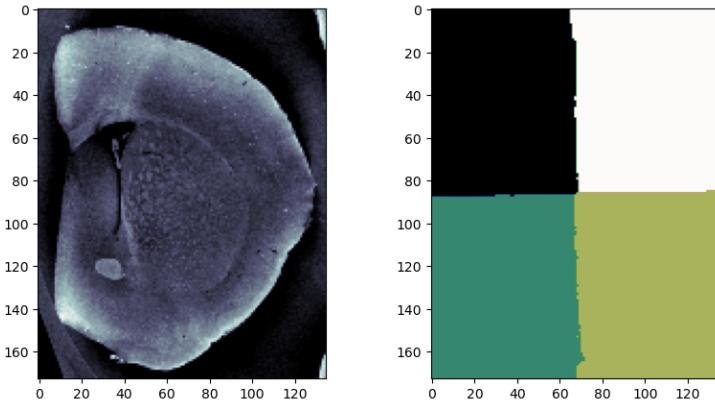
```
from sklearn.cluster import KMeans
km = KMeans(n_clusters=4, random_state=2018, n_init='auto')
cortex_df['group'] = km.fit_predict(cortex_df[['x', 'y', 'intensity']].values)
cortex_seg = cortex_df['group'].values.reshape(cortex_img.shape)

fig, (ax0, ax1, ax2) = plt.subplots(1, 3, figsize=(15, 5))
ax1.imshow(cortex_img, cmap='bone')
ax2.imshow(cortex_seg, cmap='gist_earth')

ccolors = plt.cm.BuPu(np.full(4, 0.1))

pd.plotting.table(data=cortex_df.groupby(['group']).apply(lambda x: x.sample(3)), ↴
    ax=ax0, loc='center', colColours=ccolors, fontsize=20)
ax0.axis('off');
```

	x	y	intensity	group
(0, 10455)	60.0	77.0	17.748	0.0
(0, 5977)	37.0	44.0	33.461	0.0
(0, 11652)	42.0	86.0	23.335	0.0
(1, 12448)	28.0	92.0	19.582	1.0
(1, 16085)	20.0	119.0	17.667	1.0
(1, 12975)	15.0	96.0	23.156	1.0
(2, 19830)	120.0	146.0	3.784	2.0
(2, 11742)	132.0	86.0	13.773	2.0
(2, 12222)	72.0	90.0	16.579	2.0
(3, 5513)	113.0	40.0	0.112	3.0
(3, 11157)	87.0	82.0	17.388	3.0
(3, 5922)	117.0	43.0	0.282	3.0



### Why is the image segmented like this?

#### Rescaling components

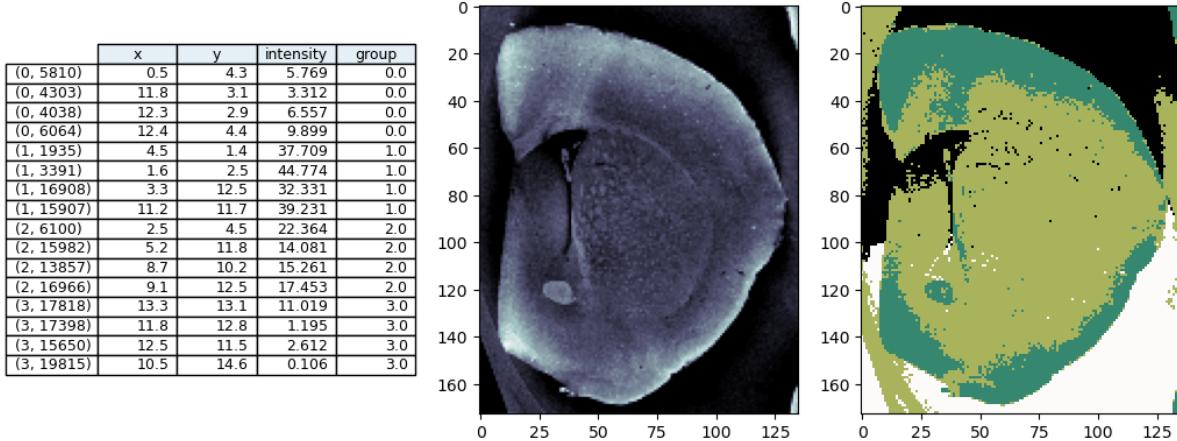
Since the distance is currently calculated by  $\|\vec{v}_i - \vec{v}_j\|$  and the values for the position is much larger than the values for the Intensity , Sobel or Gaussian they need to be rescaled so they all fit on the same axis  $\vec{v} = \{\frac{x}{10}, \frac{y}{10}, \text{Intensity}\}$

$N_{clusters}=4$

```
km = KMeans(n_clusters=4, random_state=2018, n_init='auto')
scale_cortex_df = cortex_df.copy()
scale_cortex_df.x = scale_cortex_df.x/10
scale_cortex_df.y = scale_cortex_df.y/10
scale_cortex_df['group'] = km.fit_predict(scale_cortex_df[['x', 'y', 'intensity']].values)
cortex_seg=scale_cortex_df['group'].values.reshape(cortex_img.shape)
```

```
fig, (ax0,ax1, ax2) = plt.subplots(1, 3,
                                    figsize=(12, 5))
ax1.imshow(cortex_img, cmap='bone')
ax2.imshow(cortex_seg,
           cmap='gist_earth', interpolation='none')

ccolors = plt.cm.BuPu(np.full(4, 0.1))
pd.plotting.table(data=scale_cortex_df.groupby(['group']).apply(lambda x: x.sample(4)), ax=ax0, loc='center', colColours=ccolors, fontsize=20)
ax0.axis('off');
```



### Let's try a different position scaling

$$\vec{v} = \left\{ \frac{x}{5}, \frac{y}{5}, \text{Intensity} \right\}$$

$N_{clusters}=5$

```

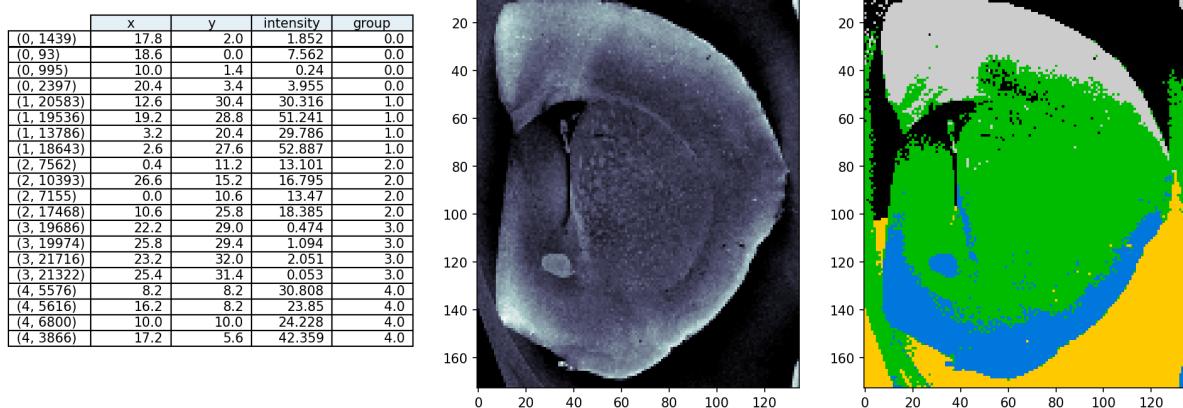
km = KMeans(n_clusters=5, random_state=2019)
scale_cortex_df = cortex_df.copy()
scale_cortex_df.x = scale_cortex_df.x/5
scale_cortex_df.y = scale_cortex_df.y/5
scale_cortex_df['group'] = km.fit_predict(
    scale_cortex_df[['x', 'y', 'intensity']].values)

fig, (ax0, ax1, ax2) = plt.subplots(1, 3,
                                    figsize=(15, 8), dpi=150)
ax1.imshow(cortex_img, cmap='bone')
ax2.imshow(scale_cortex_df['group'].values.reshape(cortex_img.shape),
           cmap='nipy_spectral')
scale_cortex_df.groupby(['group']).apply(lambda x: x.sample(3))

ccolors = plt.cm.BuPu(np.full(4, 0.1))
pd.plotting.table(data=scale_cortex_df.groupby(['group']).apply(lambda x: x.
    sample(4)), ax=ax0, loc='center', colColours=ccolors, fontsize=20)
ax0.axis('off');

/Users/kaestner/anaconda3/lib/python3.11/site-packages/sklearn/cluster/_kmeans.
    ↪py:870: FutureWarning: The default value of `n_init` will change from 10 to 'auto
    ↪' in 1.4. Set the value of `n_init` explicitly to suppress the warning
    warnings.warn(

```



### 0.6.5 When can clustering be used on images?

Clustering and in particular k-means are often used for imaging applications.

- Single images (Cortex example)

It does not make much sense to segment using only the image intensity with k-means. This would result in thresholding similar to the one provided by Otsu's method. If you, however, add spatial information like edges and positions it starts to be interesting to use k-means for a single image.

- Bimodal data

In recent years, several neutron imaging instruments have installed an X-ray source to provide complementary information to the neutron images. This is a great mix of information for k-means based segmentation. Another type of bimodal data is when a grating interferometry setup is used. This setup provides three images revealing different aspects of the samples. Again, here it is a great combination as these data are

- Hyper-spectral data (materials science example)

Many materials have a characteristic response to the neutron wavelength. This is used in many materials science experiments, in particular experiments performed at pulsed neutron sources. The data from such experiments result in a spectrum response for each pixel. This means each pixel can be a vector of >1000 elements.

## 0.7 Quad trees

Split the image in subregions until a criterion is fulfilled:

### 0.7.1 Principle

A quad tree is created by dividing image into four sections. Next you iterate the following steps until no more splits are made:

1. Compute metric (e.g. max-min or standard deviation) for each new sub-region
2. If the metric is greater than a threshold slit the region into four new regions
3. Otherwise leave it as is, the region is “constant” according to the criterion.

The figure below illustrates how an image is decomposed into a quad tree. Here you can see that the regions near edges are much smaller than in constant valued regions.

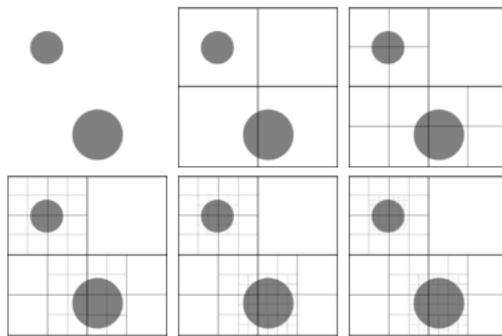


Fig. 1: A quad tree decomposition.

### 0.7.2 Quad tree example

Next we try to decompose an natural image as a quad tree.

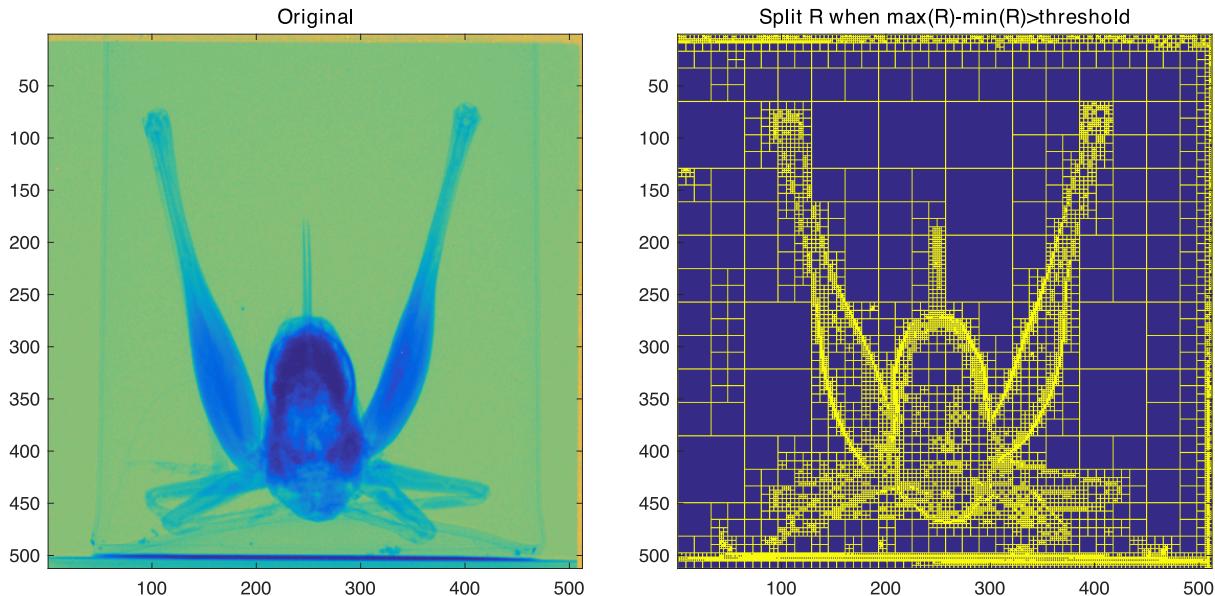


Fig. 2: A neutron radiograph of a grasshopper decomposed using a quad tree.

## 0.8 Superpixels

An approach for simplifying images by performing a clustering and forming super-pixels from groups of similar pixels.  
<https://ivrl.epfl.ch/research/superpixels>

DOI

### 0.8.1 Why use superpixels

Super pixels

- Drastically reduced data size,
  - Serves as an initial segmentation showing spatially meaningful groups
- 

### A super-pixel example

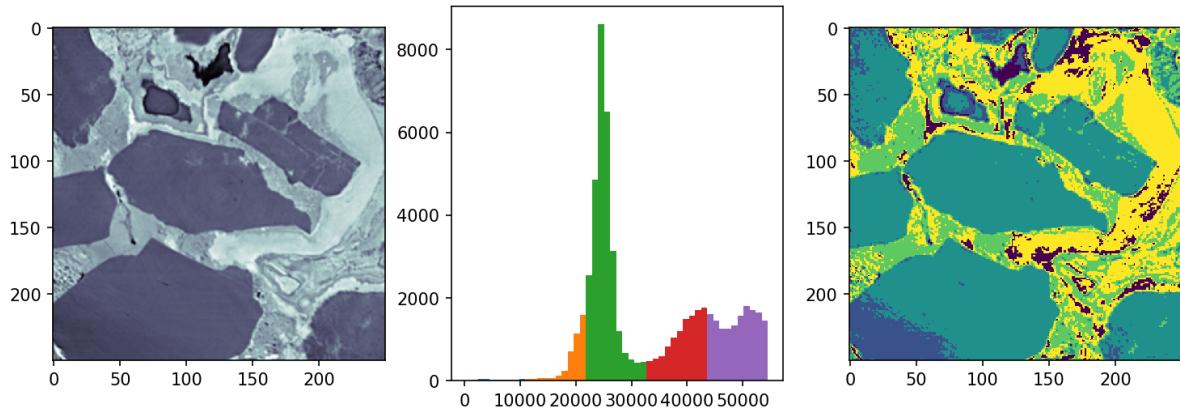
We start with an example of shale with multiple phases

- rock
- clay
- pore

### Basic thresholds

We start the analysis with using plain threshold based on the histogram. You can clearly see in the histogram that the chosen thresholds will produce many misclassified pixels. This approach will give a hint on the regions but is not very precise.

```
shale_img = imread("figures/shale-slice.tiff")
np.random.seed(2018)
fig, (ax1, ax2, ax3) = plt.subplots(1, 3,
                                    figsize=(12, 4), dpi=150)
ax1.imshow(shale_img, cmap='bone')
thresh_vals = np.linspace(shale_img.min(), shale_img.max(), 5+2)[:-1]
out_img = np.zeros_like(shale_img)
for i, (t_start, t_end) in enumerate(zip(thresh_vals, thresh_vals[1:])):
    thresh_reg = (shale_img > t_start) & (shale_img < t_end)
    ax2.hist(shale_img.ravel()[thresh_reg.ravel()])
    out_img[thresh_reg] = i
# ax3.imshow(out_img, cmap='gist_earth', interpolation='none');
ax3.imshow(out_img, cmap='viridis', interpolation='none');
```



## 0.8.2 Super pixels

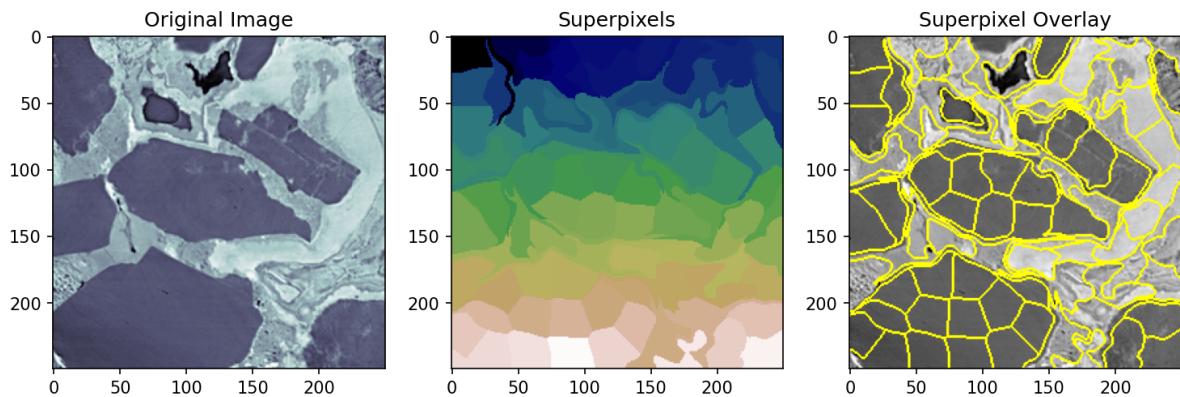
Super-pixels in as sense an evolution of the quad tree. They are not bound to the strict splitting scheme but can generate regions with limited variations of arbitrary shape.

Using the SLIC algorithm

```
from skimage.segmentation import slic, mark_boundaries

shale_segs = slic(shale_img,
                  n_segments=100,
                  compactness=5e-2,
                  sigma=3.0,
                  start_label=1,
                  channel_axis=None)
```

```
fig, (ax1, ax2, ax3) = plt.subplots(1, 3, figsize=(12, 4), dpi=150)
ax1.imshow(shale_img, cmap='bone')
ax1.set_title('Original Image')
ax2.imshow(shale_segs, cmap='gist_earth')
ax2.set_title('Superpixels')
ax3.imshow(mark_boundaries(shale_img, shale_segs))
ax3.set_title('Superpixel Overlay');
```

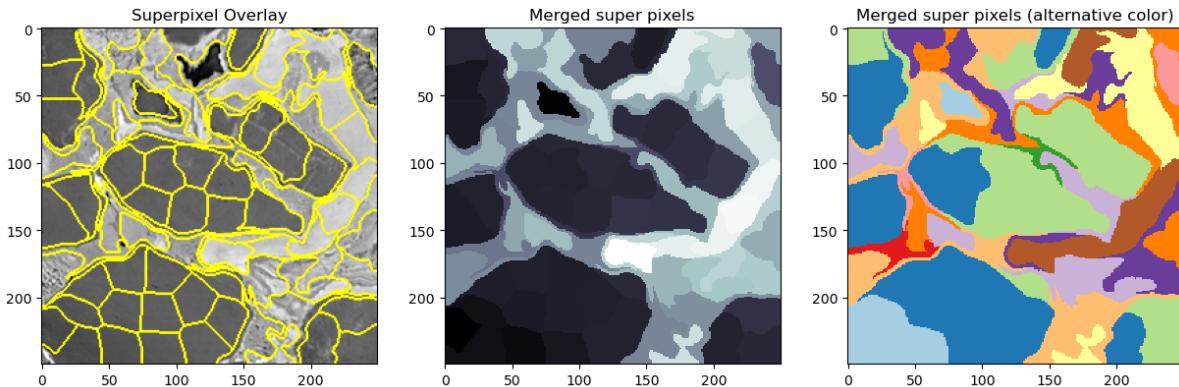


### 0.8.3 Merging super pixels

Usually, too many super pixels are identified. They are also not really representing the feature shapes in the image. The super-pixels currently also only have pixel indexes as values. In the next step we assign the average values of each super-pixel. This produces a patchy image that ideally is easier to interpret.

```
flat_shale_img = shale_img.copy()
for s_idx in np.unique(shale_segs.ravel()):
    flat_shale_img[shale_segs == s_idx] = np.mean(
        flat_shale_img[shale_segs == s_idx])
```

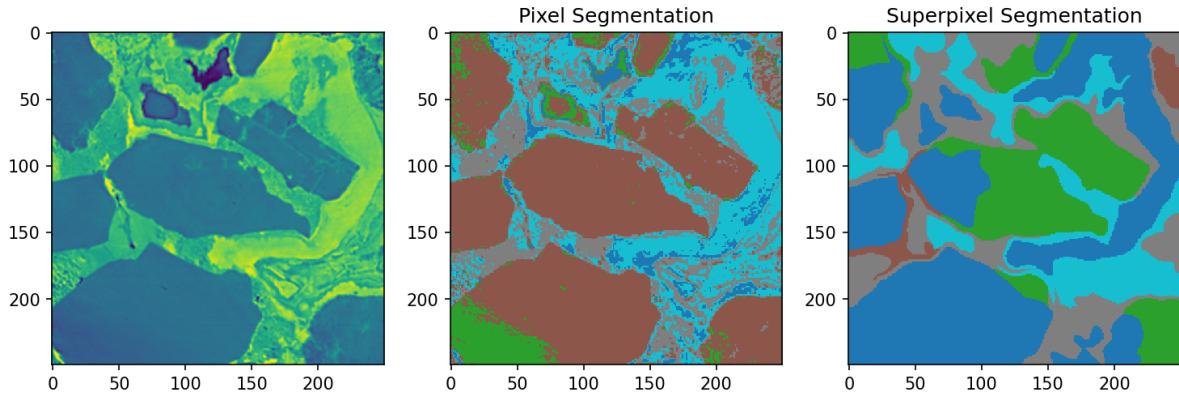
```
fig, (ax1, ax2, ax3) = plt.subplots(1, 3, figsize=(15, 5))
ax1.imshow(mark_boundaries(shale_img, shale_segs))
ax1.set_title('Superpixel Overlay')
ax2.imshow(flat_shale_img, cmap='bone', interpolation='none');
ax2.set_title('Merged super pixels')
ax3.imshow(flat_shale_img, cmap='Paired', interpolation='none');
ax3.set_title('Merged super pixels (alternative color)');
```



### 0.8.4 Segmentation using super pixels

```
thresh_vals = np.linspace(flat_shale_img.min(), flat_shale_img.max(), 5+2)[:-1]
sp_out_img = np.zeros_like(flat_shale_img)
for i, (t_start, t_end) in enumerate(zip(thresh_vals, thresh_vals[1:])):
    thresh_reg = (flat_shale_img > t_start) & (flat_shale_img < t_end)
    sp_out_img[thresh_reg] = i
```

```
fig, (ax1, ax2, ax3) = plt.subplots(1, 3, figsize=(12, 4), dpi=150)
ax1.imshow(shale_img, cmap='viridis')
ax2.imshow(out_img, cmap='tab10', interpolation='none')
ax2.set_title('Pixel Segmentation')
ax3.imshow(sp_out_img, cmap='tab10', interpolation='none')
ax3.set_title('Superpixel Segmentation');
```



## 0.9 Probabilistic Models of Segmentation

A more general approach is to use a probabilistic model to segmentation. We start with our image  $I(\vec{x}) \forall \vec{x} \in \mathbb{R}^N$  and we classify it into two phases  $\alpha$  and  $\beta$

$$P(\{\vec{x}, I(\vec{x})\}|\alpha) \propto P(\alpha) + P(I(\vec{x})|\alpha) + P\left(\sum_{x' \in \mathcal{N}} I(\vec{x}')|\alpha\right)$$

- $P(\{\vec{x}, f(\vec{x})\}|\alpha)$  the probability a given pixel is in phase  $\alpha$  given we know its position and value (what we are trying to estimate)
- $P(\alpha)$  probability of any pixel in an image being part of the phase (expected volume fraction of that phase)
- $P(I(\vec{x})|\alpha)$  probability adjustment based on knowing the value of  $I$  at the given point (standard threshold)
- $P(f(\vec{x}')|\alpha)$  are the collective probability adjustments based on knowing the value of a pixels neighbors (very simple version of [Markov Random Field](#) approaches)

## 0.10 Summary

### 0.10.1 Histogram based thresholding

- Otsu and others
- Hysteresis thresholding

### 0.10.2 Clustering - K-means

- Add position information

### 0.10.3 Similar region segmentations

- Quad trees
- Super pixels

## 0.11 Supervised Segmentation Approaches

### 0.11.1 Overview

1. Methods
2. Pipelines
3. Classification
4. Regression
5. Segmentation

### 0.11.2 Reading Material

- Introduction to Machine Learning: ETH Course
- Decision Forests for Computer Vision and Medical Image Analysis
- U-Net: Convolutional Networks for Biomedical Image Segmentation
- U-Net Website

#### Load some modules for the notebook

```
import seaborn as sns
import matplotlib.pyplot as plt
import pandas as pd
import numpy as np
from skimage.io import imread
from sklearn.datasets import make_blobs
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.tree import export_graphviz
import graphviz
from sklearn.tree import DecisionTreeClassifier
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor
#from pipe_utils import px_flatten_step, show_pipe, fit_img_pipe
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import RobustScaler
from sklearn.preprocessing import FunctionTransformer
from sklearn.cluster import KMeans
import plotsupport as ps
%matplotlib inline

plt.rcParams['font.family'] = ['sans-serif']
```

(continues on next page)

(continued from previous page)

```
plt.rcParams['font.sans-serif'] = ['DejaVu Sans']
plt.style.use('ggplot')
sns.set_style("whitegrid", {'axes.grid': False})
```

## 0.12 Basic Methods Overview

There are a number of supervised methods we can use for

- classification,
- regression
- and both.

There are a number of methods we can use for classification, regression and both. For the simplification of the material we will not make a massive distinction between classification and regression but there are many situations where this is not appropriate. Here we cover a few basic methods, since these are important to understand as a starting point for solving difficult problems. The list is not complete and importantly Support Vector Machines are completely missing which can be a very useful tool in supervised analysis. A core idea to supervised models is they have a training phase and a predicting phase.

### 0.12.1 Training

The training phase is when the parameters of the model are *learned* and involve putting inputs into the model and updating the parameters so they better match the outputs. This is a sort-of curve fitting (with linear regression it is exactly curve fitting).

- The training phase is when the parameters of the model are *learned*
- Used training data with ground truth.

### 0.12.2 Predicting

The predicting phase is once the parameters have been set applying the model to new datasets. At this point the parameters are no longer adjusted or updated and the model is frozen. Generally it is not possible to tweak a model any more using new data but some approaches (most notably neural networks) are able to handle this.

- Provides responses on inputs using the trained model.
- Uses new unseen data.

## 0.13 Classification

### 0.13.1 Lets create some data...

Here we create some bivariate data ‘blobs’ with Gaussian distribution. This time the blobs have classes assigned to them. The table

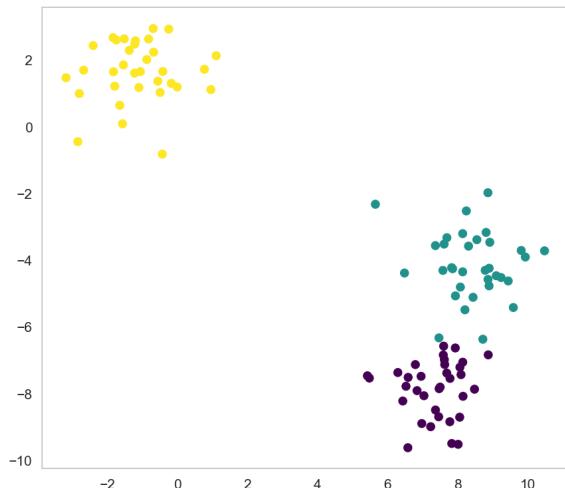
```

blob_data, blob_labels = make_blobs(n_samples=100,
                                     random_state=2018)
test_pts = pd.DataFrame(blob_data, columns=['x', 'y'])
test_pts['group_id'] = blob_labels

fig,ax = plt.subplots(1,2,figsize=(15,6),dpi=150)
ax[1].scatter(test_pts.x, test_pts.y,
               c=test_pts.group_id,
               cmap='viridis')
ccolors = plt.cm.BuPu(np.full(3, 0.1))
pd.plotting.table(data=test_pts.sample(10).round(decimals=2), ax=ax[0], loc='center',
                   colColours=ccolors)
ax[0].axis('off');

```

	x	y	group_id
21	8.06	-8.71	0.0
87	9.45	-4.62	1.0
11	-1.2	2.57	2.0
85	8.44	-5.11	1.0
47	7.23	-8.99	0.0
82	7.6	-6.58	0.0
83	8.31	-3.58	1.0
13	6.96	-7.48	0.0
32	7.04	-8.06	0.0
31	8.88	-6.84	0.0



### 0.13.2 Nearest Neighbor (or K Nearest Neighbors)

The technique is as basic as it sounds, it basically finds the nearest point to what you have put in.

The *k nearest neighbors* algorithm makes the inference based on a point cloud of training point. When the model is presented with a new point it computes the distance to the closest points in the model. The *k* in the algorithm name indicates how many neighbors should be considered. E.g. *k*=3 means that the major class of the three nearest neighbors is assigned the tested point.

Looking at the example below we would say that using three neighbors the

- Green hiker would claim he is in a spruce forest.
- Orange hiker would claim he is in a mixed forest.
- Blue hiker would claim he is in a birch forest.

```

plt.figure(figsize=[6,6],dpi=100)
plt.scatter(test_pts.x, test_pts.y,
            c=test_pts.group_id,
            cmap='viridis')
plt.plot(2,-2,'X',color='cornflowerblue',markersize=10,label='Point A')
plt.plot(7.7,-6.1,'o',color='darkorange',markersize=10,label='Point B')
plt.plot(8.5,-4,'P',color='deeppink',markersize=10,label='Point C')
plt.legend();

```

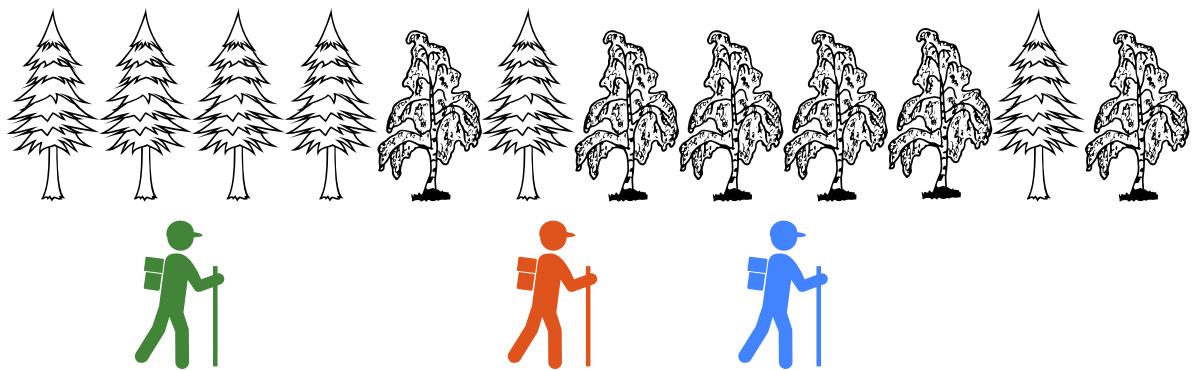
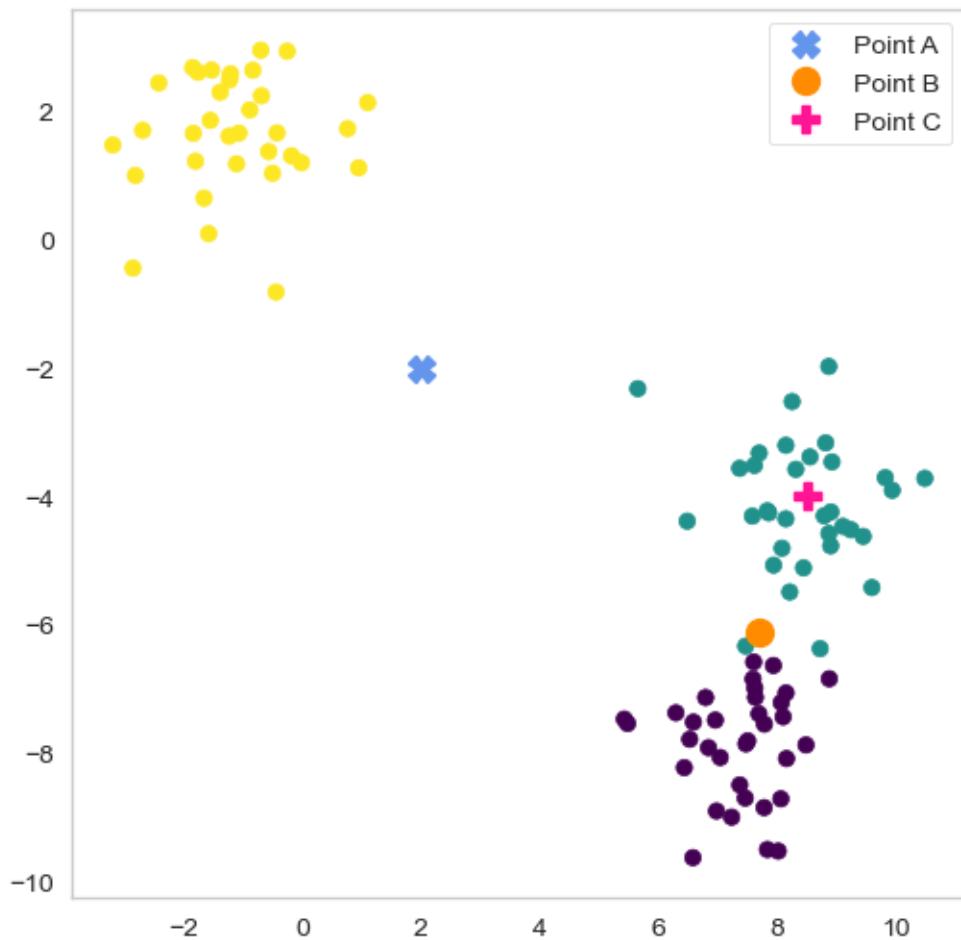


Fig. 1: Depending on the where the hiker is standing he makes the the conclusion that is either in a birch or spruce forest.



### 0.13.3 Text example

#### Training data

Value	1	2	3	4
Class	I	am	a	dog

#### Start the training

```
from sklearn.neighbors import KNeighborsClassifier
import numpy as np
k_class = KNeighborsClassifier(1)
k_class.fit(X=np.reshape([0, 1, 2, 3], (-1, 1)),
             y=['I', 'am', 'a', 'dog'])
```

```
KNeighborsClassifier(n_neighbors=1)
```

### 0.13.4 Nearest neighbor predictions

#### Basic test

Same values as training

```
print(k_class.predict(np.reshape([0, 1, 2, 3],
                                 (-1, 1))))
```

```
['I' 'am' 'a' 'dog']
```

#### Testing with different values

```
print("Input : ", k_class.predict(np.reshape([1.5], (1, 1))))
print("Input 100 : ", k_class.predict(np.reshape([100], (1, 1))))
```

```
Input : ['am']
Input 100 : ['dog']
```

### 0.13.5 Let's come back to the blob data

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import make_blobs
```

```
blob_data, blob_labels = make_blobs(n_samples=100,
                                      cluster_std=2.0,
                                      random_state=2018)
```

(continues on next page)

(continued from previous page)

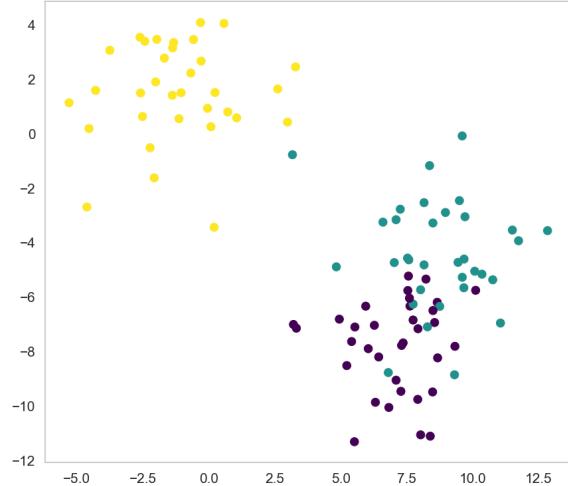
```

test_pts = pd.DataFrame(blob_data, columns=['x', 'y'])
test_pts['group_id'] = blob_labels

fig, ax = plt.subplots(1,2, figsize=(15,6), dpi=150)
ax[1].scatter(test_pts.x, test_pts.y, c=test_pts.group_id, cmap='viridis')
ccolors = plt.cm.BuPu(np.full(3, 0.1))
pd.plotting.table(data=test_pts.sample(10).round(decimals=2), ax=ax[0], loc='center',
                   colColours=ccolors)
ax[0].axis('off');

```

	x	y	group_id
21	8.48	-9.5	0.0
96	7.74	-6.27	1.0
10	-4.29	1.59	2.0
20	3.2	-7.02	0.0
73	2.59	1.5	2.0
72	-2.07	-1.63	2.0
53	9.49	-2.46	1.0
95	7.03	-4.74	1.0
6	-0.05	0.93	2.0
66	8.39	-11.12	0.0



### 0.13.6 Training the model using one neighbor

```

# Define classifier
k_class = KNeighborsClassifier(1)

# Train the classifier model with data
k_class.fit(test_pts[['x', 'y']], test_pts['group_id'])

```

```
KNeighborsClassifier(n_neighbors=1)
```

### Resulting prediction map for a single neighbor

```

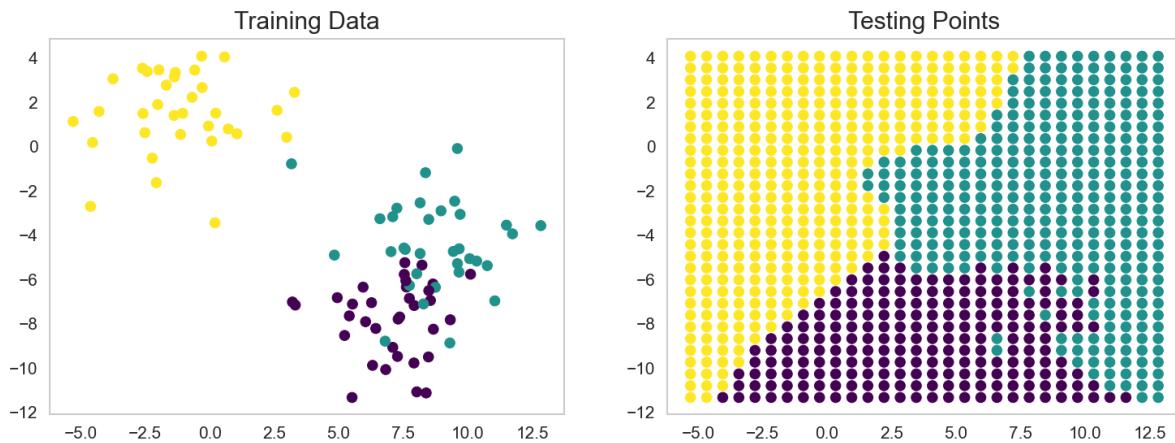
xx, yy = np.meshgrid(np.linspace(test_pts.x.min(), test_pts.x.max(), 30), np.
                     linspace(test_pts.y.min(), test_pts.y.max(), 30), indexing='ij');
grid_pts = pd.DataFrame(dict(x=xx.ravel(), y=yy.ravel()))
grid_pts['predicted_id'] = k_class.predict(grid_pts[['x', 'y']])

```

```

fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 4), dpi=150)
ax1.scatter(test_pts.x, test_pts.y, c=test_pts.group_id, cmap='viridis');
ax1.set_title('Training Data');
ax2.scatter(grid_pts.x, grid_pts.y, c=grid_pts.predicted_id, cmap='viridis');
ax2.set_title('Testing Points');

```



In this example we trained the classifier with some data points appearing in clouds. The next step is to see how the classifier performs on unknown data. This is done by creating a mesh of equidistant points, each point is then fed to the classifier and we can see where the decision boundaries are for the trained classifier. You can see that there are some irregularities along the boundaries between the classes. One reason is that the training data has overlapping classes and the other is the trained model.

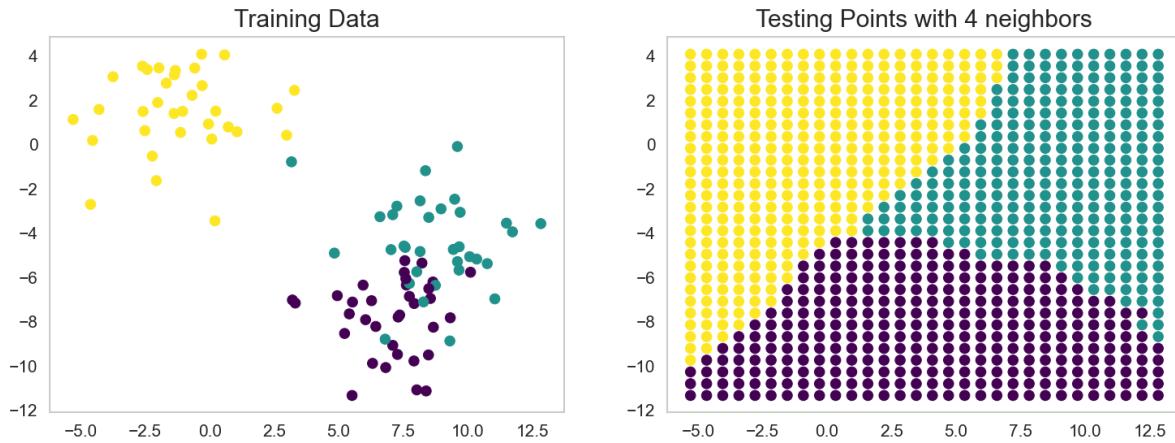
Here, we used a data frame for the presentation of the point tables.

### 0.13.7 Stabilizing Results - Increase number of neighbors

- We can see here that the result is thrown off by single points
- Prediction can be improved by using more than the nearest neighbor

```
k_class = KNeighborsClassifier(n_neighbors=4)
k_class = k_class.fit(test_pts[['x', 'y']], test_pts['group_id'])
xx, yy = np.meshgrid(np.linspace(test_pts.x.min(), test_pts.x.max(), 30),
                      np.linspace(test_pts.y.min(), test_pts.y.max(), 30),
                      indexing='ij')
grid_pts = pd.DataFrame(dict(x=xx.ravel(), y=yy.ravel()))
grid_pts['predicted_id'] = k_class.predict(grid_pts[['x', 'y']])
```

```
fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 4), dpi=150)
ax1.scatter(test_pts.x, test_pts.y, c=test_pts.group_id, cmap='viridis');
ax1.set_title('Training Data');
ax2.scatter(grid_pts.x, grid_pts.y, c=grid_pts.predicted_id, cmap='viridis');
ax2.set_title('Testing Points with 4 neighbors');
```



## 0.14 Linear Regression

- Linear regression is a fancy-name for linear curve fitting
- Fitting a line through points (sometimes in more than one dimension).
- It is a very basic method,
  - is easy to understand,
  - interpret
  - and fast to compute

### 0.14.1 We need data to fit...

```
from sklearn.linear_model import LinearRegression

x = np.array([0, 1, 2, 3])
y = np.array([10, 20, 30, 40] )

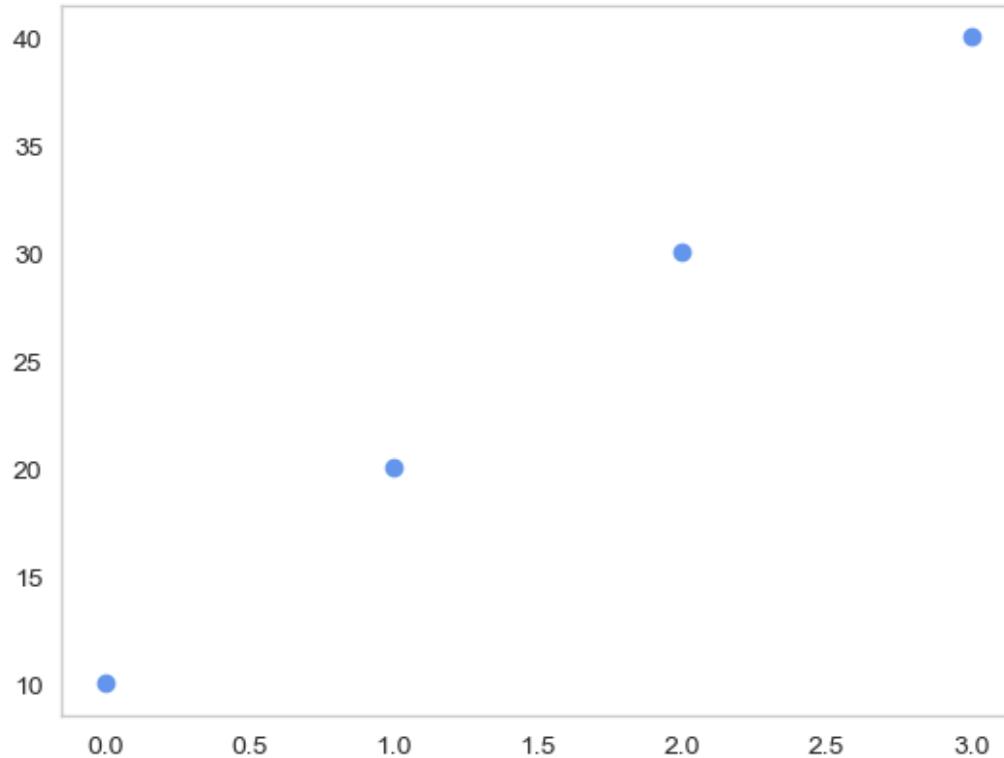
plt.plot(x,y,'o',color='cornflowerblue'); plt.title('Training data')

l_reg = LinearRegression()
l_reg.fit(X=np.reshape(x, (-1, 1)),y=y)

print("slope: {0:0.4}, intercept: {1:0.4}".format(l_reg.coef_[0], l_reg.intercept_))
```

slope: 10.0, intercept: 10.0

Training data



Let's try the model on some data points

```
print('An array of values:', l_reg.predict(np.reshape([0, 1, 2, 3], (-1, 1))))
print('x: -100, =>', l_reg.predict(np.reshape([-100], (1, 1))))
print('x: 500, =>', l_reg.predict(np.reshape([500], (1, 1))))
```

```
An array of values: [10. 20. 30. 40.]
x: -100 => [-990.]
x: 500 => [5010.]
```

### 0.14.2 Regression on blob data

```
from sklearn.datasets import make_blobs
import matplotlib.pyplot as plt

import numpy as np
import pandas as pd
```

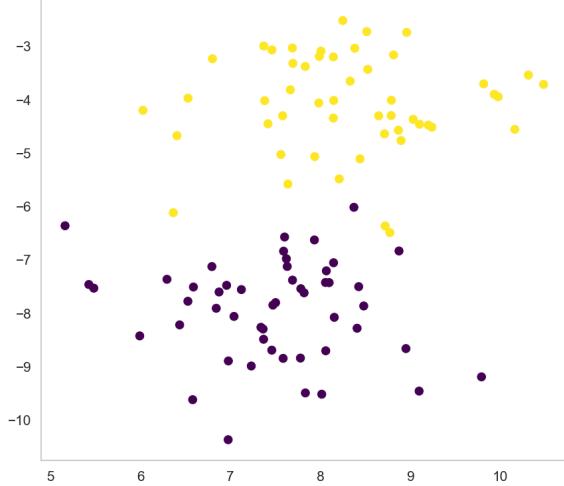
```
blob_data, blob_labels = make_blobs(centers=2, n_samples=100,
                                    random_state=2018)
test_pts = pd.DataFrame(blob_data, columns=['x', 'y'])
test_pts['group_id'] = blob_labels
```

```

fig, ax = plt.subplots(1,2, figsize=(15,6), dpi=150)
ax[1].scatter(test_pts.x, test_pts.y, c=test_pts.group_id, cmap='viridis')

ccolors = plt.cm.BuPu(np.full(3, 0.1))
pd.plotting.table(data=test_pts.sample(10).round(decimals=2), ax=ax[0], loc='center',
                   colColours=ccolors)
ax[0].axis('off');
    
```

	x	y	group_id
21	7.59	-6.84	0.0
63	7.47	-7.85	0.0
11	6.37	-6.12	1.0
53	8.53	-3.44	1.0
52	6.96	-7.48	0.0
96	10.32	-3.54	1.0
14	8.25	-2.52	1.0
89	6.98	-8.9	0.0
92	7.7	-3.33	1.0
32	8.07	-7.21	0.0



### 0.14.3 Train the regression model

```

l_reg = LinearRegression()
l_reg.fit(test_pts[['x', 'y']], test_pts['group_id'])
print('Slope', l_reg.coef_)
print('Offset', l_reg.intercept_)
    
```

```

Slope [0.04813137 0.20391973]
Offset 1.346929708594462
    
```

### Evaluate the regression model

```

xx, yy = np.meshgrid(np.linspace(test_pts.x.min(), test_pts.x.max(), 20),
                      np.linspace(test_pts.y.min(), test_pts.y.max(), 20),
                      indexing='ij')
grid_pts = pd.DataFrame(dict(x=xx.ravel(), y=yy.ravel()))
grid_pts['predicted_id'] = l_reg.predict(grid_pts[['x', 'y']])
    
```

```

fig, (ax1, ax2, ax3, ax4) = plt.subplots(1, 4, figsize=(15, 4), dpi=150)
ax1.scatter(test_pts.x, test_pts.y, c=test_pts.group_id, cmap='viridis');
ax1.set_title('Training Data')
ax2.scatter(grid_pts.x, grid_pts.y, c=grid_pts.predicted_id, cmap='viridis');
ax2.set_title('Testing Points')
ax3.imshow(grid_pts.predicted_id.values.reshape(
    
```

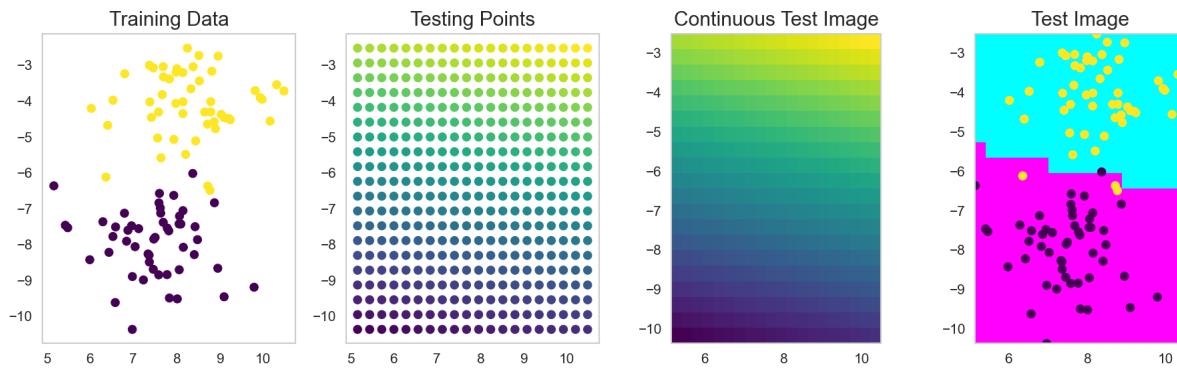
(continues on next page)

(continued from previous page)

```

xx.shape).T[::-1], cmap='viridis', extent=[test_pts.x.min(), test_pts.x.max(), -
→test_pts.y.min(), test_pts.y.max()])
ax3.set_title('Continuous Test Image');
ax4.imshow(grid_pts.predicted_id.values.reshape(
    xx.shape).T[::-1]<0.5, cmap='cool', extent=[test_pts.x.min(), test_pts.x.max(), -
→test_pts.y.min(), test_pts.y.max()])
ax4.scatter(test_pts.x, test_pts.y, c=test_pts.group_id, cmap='viridis');
ax4.set_title('Test Image');

```



## 0.15 Decision trees

- SciKit Learn documentation on trees
- SciKit Learn trees explained
- Hastie et al., Elements Of Statistical Learning, 2009 Section 9.2 Trees.

```

from sklearn.tree import export_graphviz
import sklearn.tree as tree
import graphviz
from sklearn.tree import DecisionTreeClassifier
import numpy as np
from IPython.display import SVG

def show_tree(in_tree):
    return graphviz.Source(export_graphviz(in_tree, out_file=None))

```

### 0.15.1 Create a decision tree classifier

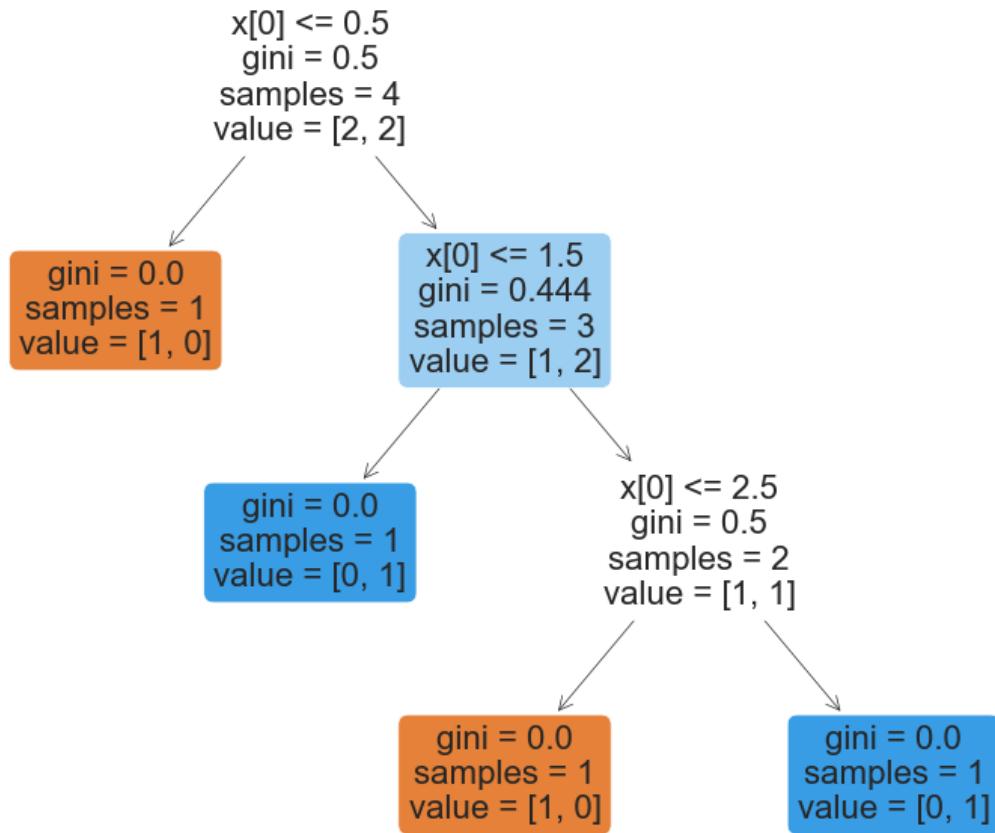
We want to identify odd numbers in the sequence [0, 1, 2, 3]

```

d_tree = DecisionTreeClassifier()
d_tree.fit(X=np.reshape([0, 1, 2, 3], (-1, 1)),
           y=[0, 1, 0, 1])
fig, axes = plt.subplots(nrows = 1, ncols = 1, figsize = (10,8))

tree.plot_tree(d_tree, filled=True, rounded=True, fontsize=16);

```



$$\text{Gini index} = 1 - \sum_{i=1}^N P_i^2$$

In the context of decision trees and machine learning, the Gini index is a measure of impurity used to evaluate the quality of a split in a decision tree. It is often used as a criterion for determining the optimal split when growing a decision tree.

The Gini index for a node in a decision tree is calculated as follows:

$$G = 1 - \sum_{i=1}^c p_i^2$$

Where:

- $G$  is the Gini index for the node.
- $c$  is the number of classes (or categories) in the dataset.
- $p_i$  is the proportion of instances in the  $i$ -th class at the node.

The Gini index ranges from 0 to 1, with lower values indicating less impurity. A node with a Gini index of 0 means that all instances belong to the same class, making it a pure node. A node with a Gini index of 1 means that the distribution of classes is perfectly impure, with an equal proportion of instances from each class.

When building a decision tree, the algorithm evaluates various potential splits in the data and selects the split that minimizes the Gini index, as it aims to create child nodes that are as pure as possible. This process is repeated recursively for each

node in the tree until certain stopping criteria are met, such as reaching a maximum depth or minimum number of samples per node.

### 0.15.2 Decision trees on the blob data

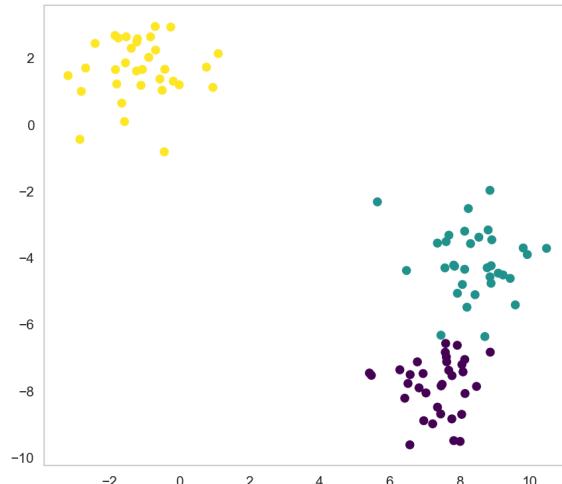
```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import make_blobs
%matplotlib inline
```

```
blob_data, blob_labels = make_blobs(n_samples=100, random_state=2018)
test_pts = pd.DataFrame(blob_data, columns=['x', 'y'])
test_pts['group_id'] = blob_labels
```

```
fig, ax = plt.subplots(1, 2, figsize=(15, 6), dpi=150)

ccolors = plt.cm.BuPu(np.full(3, 0.1))
pd.plotting.table(data=test_pts.sample(10).round(decimals=2), ax=ax[0], loc='center',
                   colColours=ccolors);
ax[0].axis('off');
ax[1].scatter(test_pts.x, test_pts.y, c=test_pts.group_id, cmap='viridis');
```

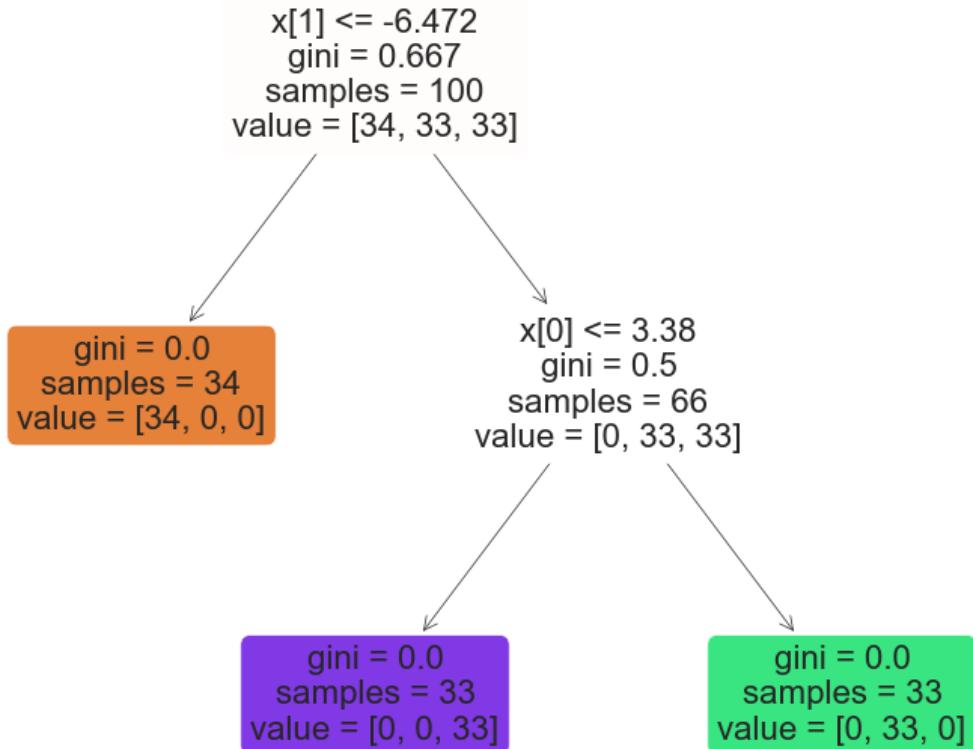
	x	y	group_id
13	6.96	-7.48	0.0
55	8.15	-3.2	1.0
92	7.62	-3.52	1.0
64	6.3	-7.37	0.0
6	-0.56	1.37	2.0
76	5.65	-2.32	1.0
3	-0.69	2.94	2.0
47	7.23	-8.99	0.0
16	9.11	-4.47	1.0
57	7.37	-8.49	0.0



### Train the tree

```
d_tree = DecisionTreeClassifier()
d_tree.fit(test_pts[['x', 'y']],
           test_pts['group_id'])

fig, axes = plt.subplots(nrows = 1, ncols = 1, figsize = (10,8))
tree.plot_tree(d_tree, filled=True, rounded=True, fontsize=16);
```



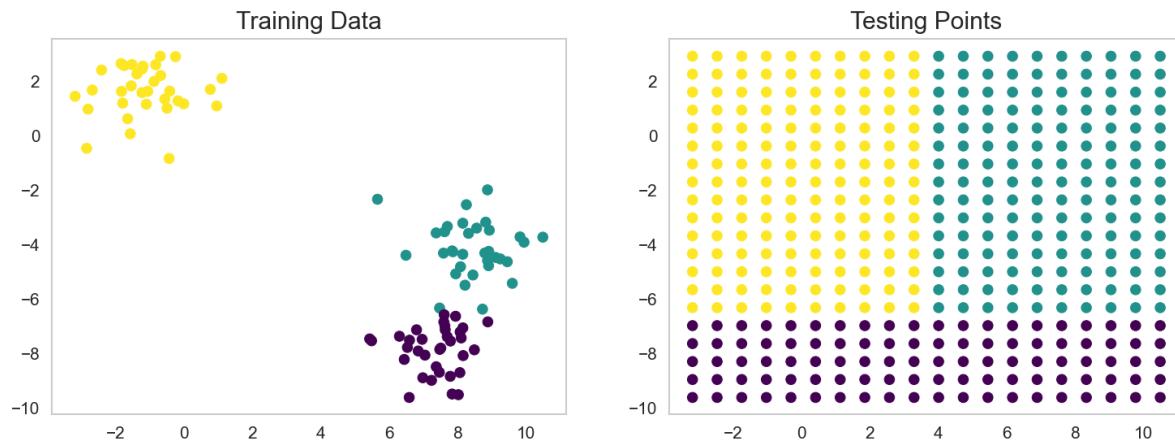
### Let's look at the decision map

```

xx, yy = np.meshgrid(np.linspace(test_pts.x.min(), test_pts.x.max(), 20),
                      np.linspace(test_pts.y.min(), test_pts.y.max(), 20),
                      indexing='ij')
)
grid_pts = pd.DataFrame(dict(x=xx.ravel(), y=yy.ravel()))
grid_pts['predicted_id'] = d_tree.predict(grid_pts[['x', 'y']])
  
```

```

fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 4), dpi=150)
ax1.scatter(test_pts.x, test_pts.y, c=test_pts.group_id, cmap='viridis')
ax1.set_title('Training Data')
ax2.scatter(grid_pts.x, grid_pts.y, c=grid_pts.predicted_id, cmap='viridis')
ax2.set_title('Testing Points');
  
```



### 0.15.3 Random Forests

Forests are basically the idea of taking a number of trees and bringing them together.

So rather than taking a single tree to do the classification, you divide the samples and the features to make different trees and then combine the results. One of the more successful approaches is called **Random Forests** or as a [video](#)

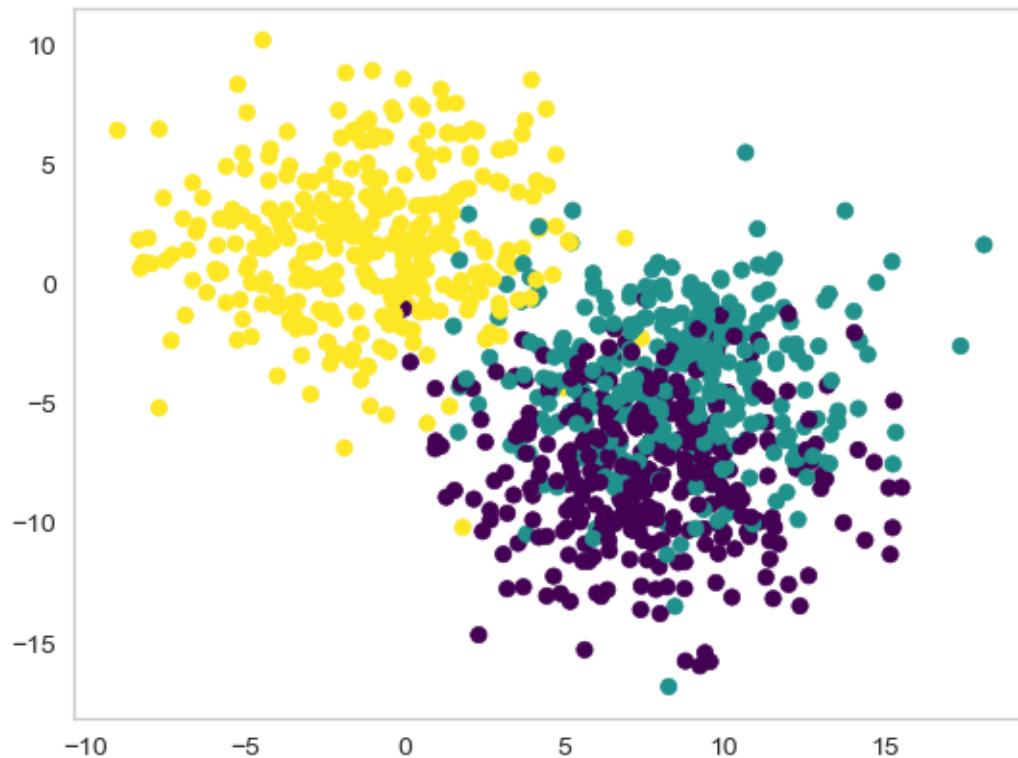


Fig. 1: Random forests train trees on different fractions of the data.

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import make_blobs
%matplotlib inline
```

**Let's make some new blob data**

```
blob_data, blob_labels = make_blobs(n_samples=1000,
                                     cluster_std=3,
                                     random_state=2018)
test_pts = pd.DataFrame(blob_data, columns=['x', 'y'])
test_pts['group_id'] = blob_labels
plt.scatter(test_pts.x, test_pts.y, c=test_pts.group_id, cmap='viridis');
```



## Train a forest

Now we train a tiny forest with five trees.

```
from sklearn.ensemble import RandomForestClassifier
rf_class = RandomForestClassifier(n_estimators=5, random_state=2018)
rf_class.fit(test_pts[['x', 'y']], test_pts['group_id'])

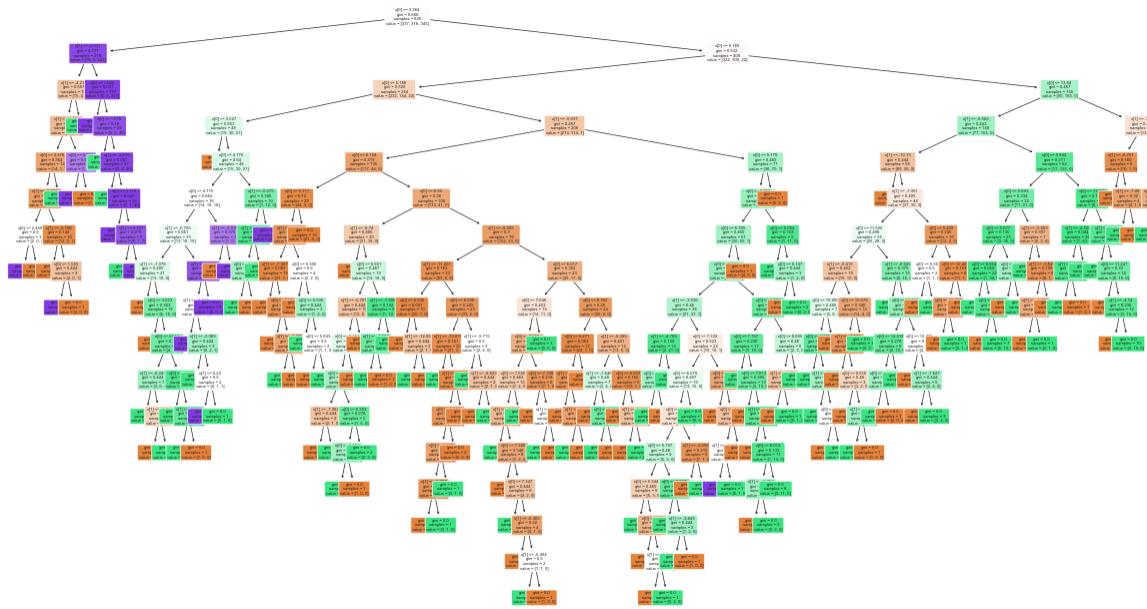
print('Build ', len(rf_class.estimators_), 'decision trees')
```

Build 5 decision trees

## Inspect one tree

Looking at a single tree in the forest

```
fig, axes = plt.subplots(nrows = 1, ncols = 1, figsize = (15,8))
tree.plot_tree(rf_class.estimators_[1], filled=True, rounded=True, fontsize=3);
```



### Looking at the performance of the forest

Here, we'll look at how the forest behaves and some of its trees.

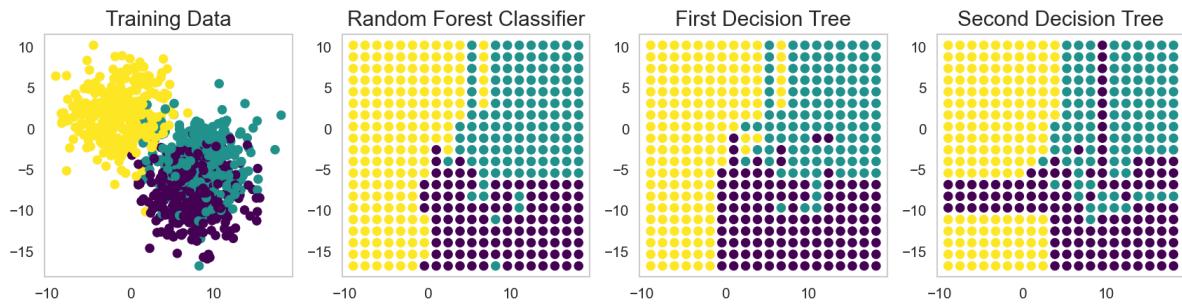
```
xx, yy = np.meshgrid(np.linspace(test_pts.x.min(), test_pts.x.max(), 20),
                     np.linspace(test_pts.y.min(), test_pts.y.max(), 20),
                     indexing='ij')
grid_pts = pd.DataFrame(dict(x=xx.ravel(), y=yy.ravel()))
rnd_forest_classes = rf_class.predict(grid_pts[['x', 'y']]) # Prediction of the forest

rnd_tree1 = rf_class.estimators_[0].predict(grid_pts[['x', 'y']].values) # ↵Prediction of the first tree
rnd_tree2 = rf_class.estimators_[1].predict(grid_pts[['x', 'y']].values) # ↵Prediction of the second tree
```

```
# Visualization
fig, (ax1, ax2, ax3, ax4) = plt.subplots(1, 4, figsize=(14, 3), dpi=150)
ax1.scatter(test_pts.x, test_pts.y, c=test_pts.group_id, cmap='viridis')
ax1.set_title('Training Data')
ax2.scatter(grid_pts.x, grid_pts.y, c=rnd_forest_classes, cmap='viridis')
ax2.set_title('Random Forest Classifier')

ax3.scatter(grid_pts.x, grid_pts.y, c=rnd_tree1, cmap='viridis')
ax3.set_title('First Decision Tree')

ax4.scatter(grid_pts.x, grid_pts.y, c=rnd_tree2, cmap='viridis')
ax4.set_title('Second Decision Tree');
```



## 0.16 Pipelines

We will use the idea of pipelines generically here to refer to the combination of steps that need to be performed to solve a problem.

Pipelines are a technical solution to reduce the amount of coding.

Consists of

- transformations
- predictors

### 0.16.1 Let's return to the blobs

```
blob_data, blob_labels = make_blobs(n_samples=100,
                                     random_state=2018)
test_pts = pd.DataFrame(blob_data, columns=['x', 'y'])
test_pts['group_id'] = blob_labels

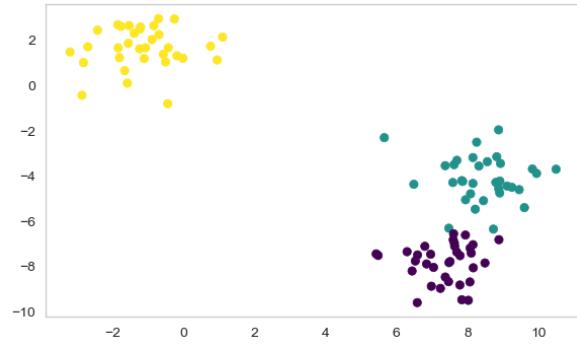
fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(15, 4))
cell_text = []
for row in range(5):
    cell_text.append(test_pts.sample(5).iloc[row])

ax1.table(cellText=cell_text, colLabels=test_pts.columns, loc='center'); ax1.axis('off')

ax2.scatter(test_pts.x, test_pts.y, c=test_pts.group_id, cmap='viridis');

/Users/kaestner/anaconda3/lib/python3.11/site-packages/matplotlib/table.py:809: FutureWarning: Series.__getitem__ treating keys as positions is deprecated. In a future version, integer keys will always be treated as labels (consistent with DataFrame behavior). To access a value by position, use `ser.iloc[pos]`  
text=cellText[row][col],
```

x	y	group_id
8.068840541544834	-7.2086649459478185	0.0
9.819917441984792	-3.708752316916487	1.0
7.696355814943029	-3.3250766781974708	1.0
7.692387397193143	-7.382755171300531	0.0
7.780745328541355	-8.841535234441094	0.0

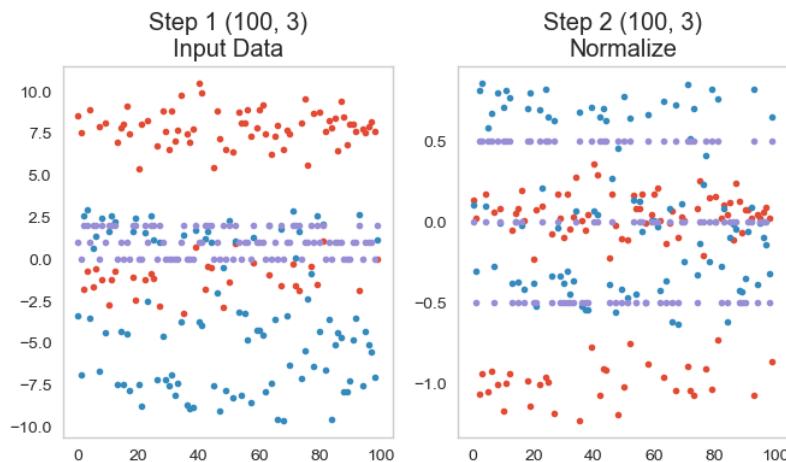


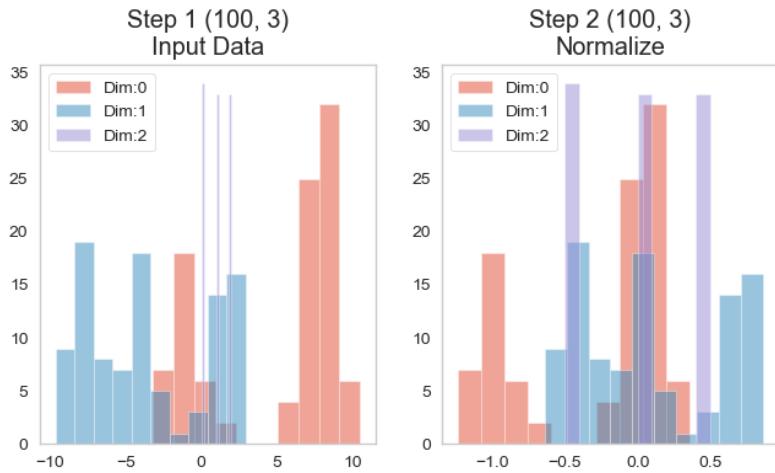
## 0.16.2 A basic pipeline

Our first pipeline has one step that normalizes the data using a RobustScaler

```
from pipe_utils import show_pipe
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import RobustScaler
simple_pipe = Pipeline([('Normalize', RobustScaler())])
simple_pipe.fit(test_pts.values)

show_pipe(simple_pipe, test_pts.values)
show_pipe(simple_pipe, test_pts.values, show_hist=True, show_legend=True)
```



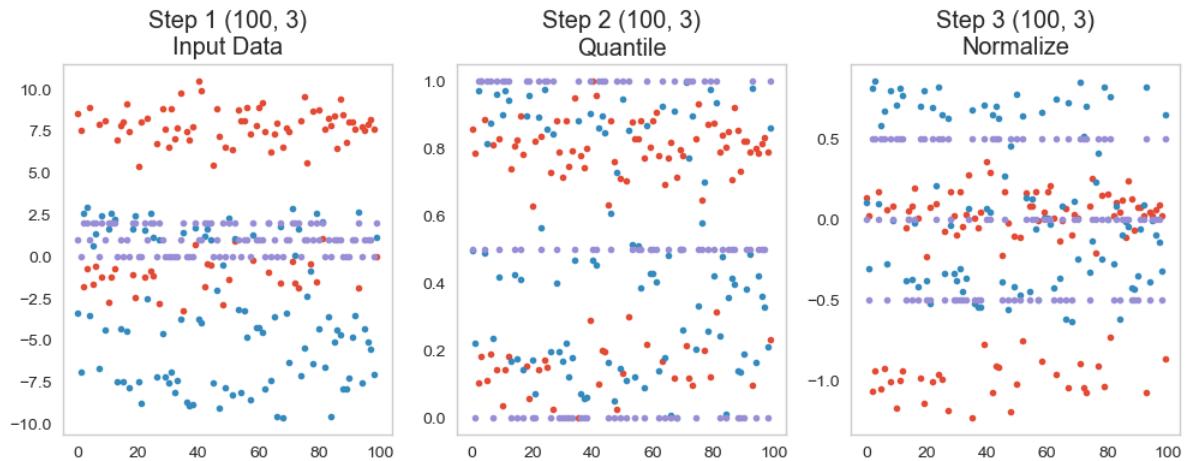


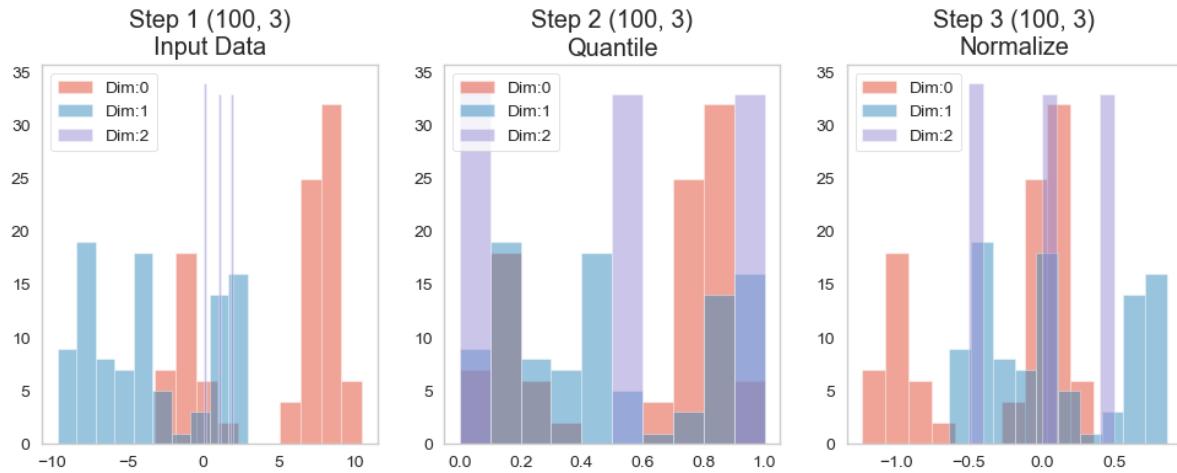
### Adding tasks to the pipeline

Next we populate the pipeline with a QuantileTransformer to force the data into a uniform distribution.

```
from sklearn.preprocessing import QuantileTransformer
longer_pipe = Pipeline([('Quantile', QuantileTransformer(n_quantiles=2)),
                      ('Normalize', RobustScaler())])
longer_pipe.fit(test_pts.values)

show_pipe(longer_pipe, test_pts.values)
show_pipe(longer_pipe, test_pts.values, show_hist=True)
```



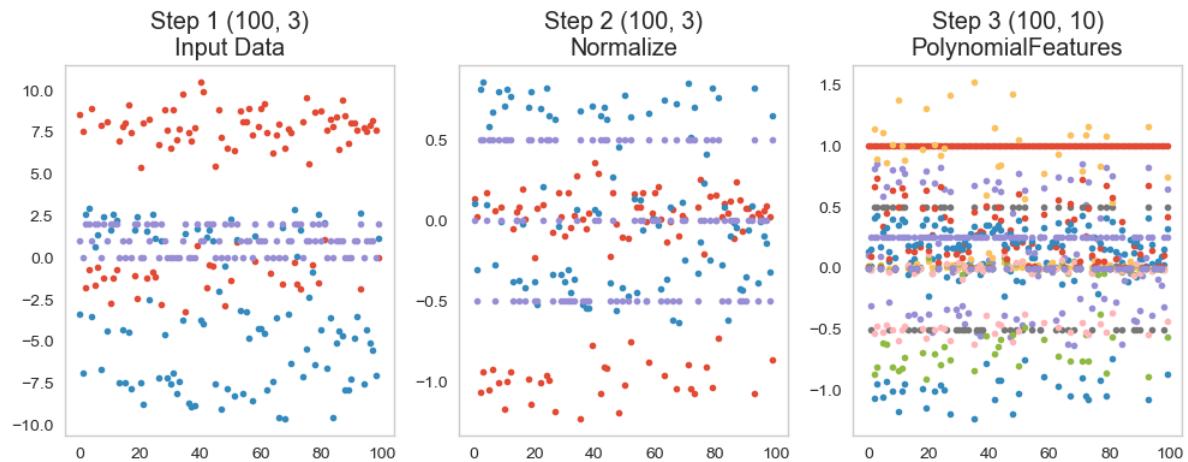


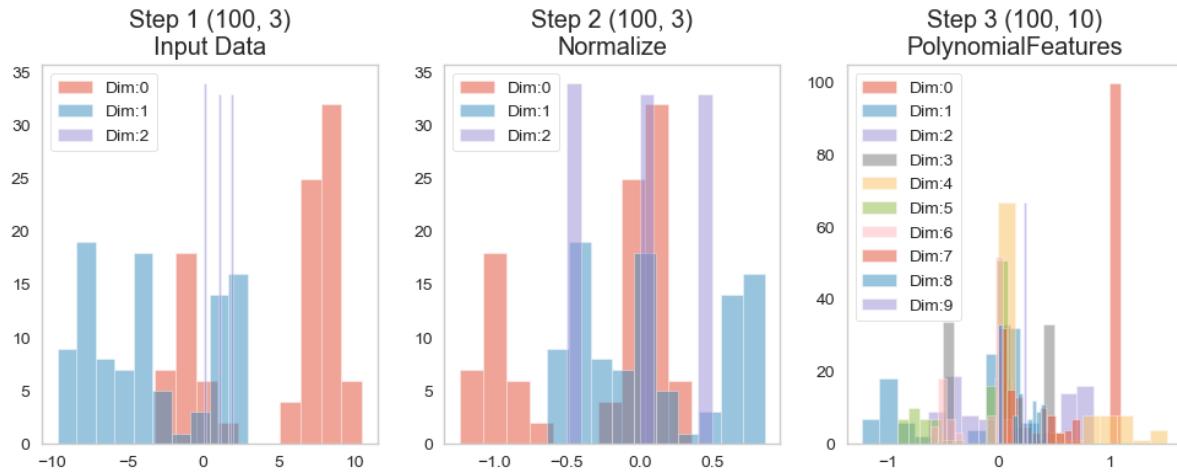
### Adding polynomial features to the pipeline

We saw earlier that it can make sense to add polynomial features to the segmentation.

```
from sklearn.preprocessing import PolynomialFeatures
messy_pipe = Pipeline([
    ('Normalize', RobustScaler()),
    ('PolynomialFeatures', PolynomialFeatures(2))])
messy_pipe.fit(test_pts.values)

show_pipe(messy_pipe, test_pts.values)
show_pipe(messy_pipe, test_pts.values, show_hist=True)
```





## 0.17 Classification using a pipeline

A common problem of putting images into categories.

- The standard problem for this is classifying digits between 0 and 9 (MNIST).
- Fundamentally a classification problem is one where we are taking a large input (images, vectors, ...) and trying to put it into a category.
  - Cats, Dogs
  - Cars, boats
  - etc.

### 0.17.1 A first example with pipeline processing

Let's load some images

- Images of numbers 0 to 9 (MNIST)
- $8 \times 8$  pixels
- 50 Samples

```
from sklearn.datasets import load_digits
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from pipe_utils import show_pipe
%matplotlib inline

# Load the number images
digit_ds = load_digits(return_X_y=False)

# Select the first 50 images and labels
img_data = digit_ds.images[:50]
digit_id = digit_ds.target[:50]
print('Image Data', img_data.shape)
```

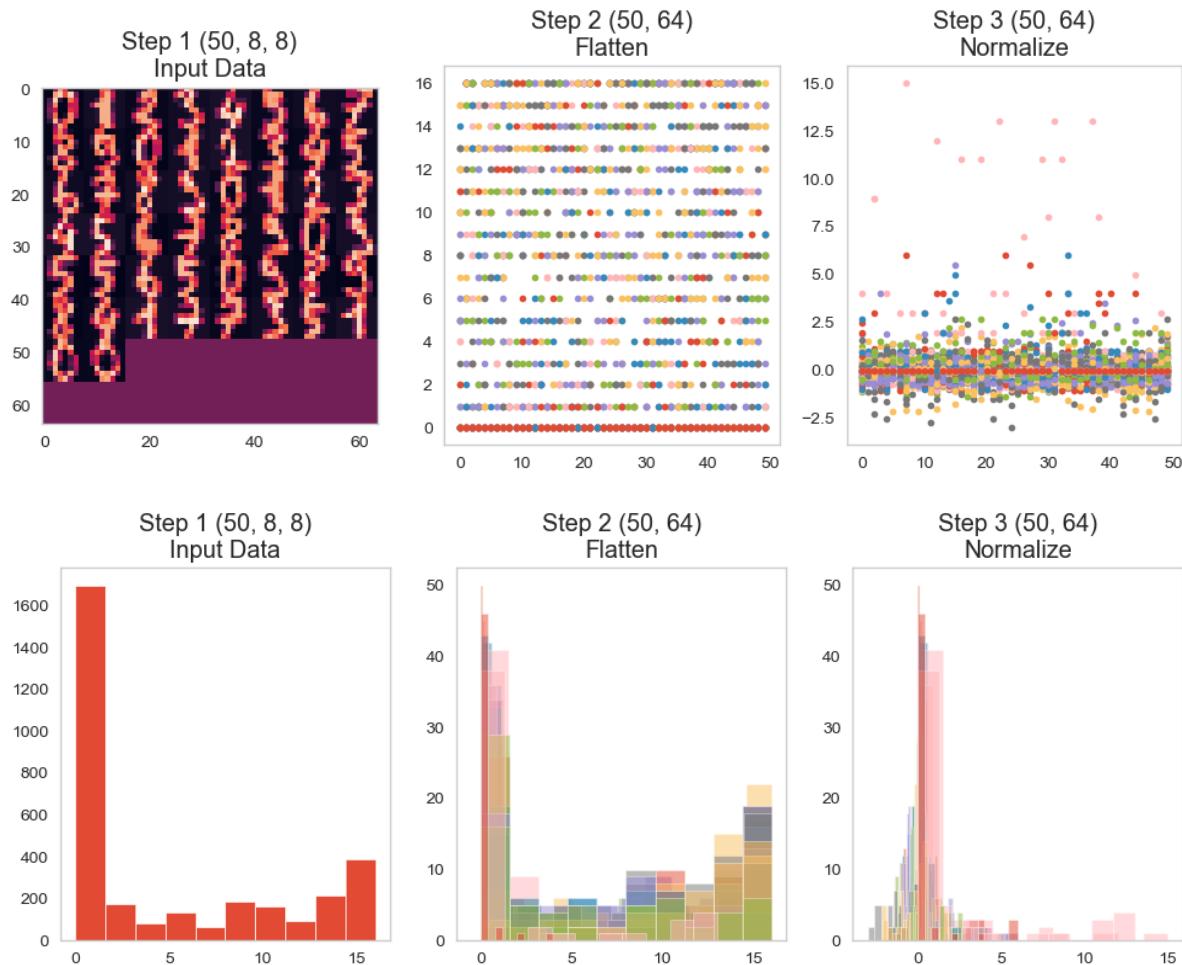
Image Data (50, 8, 8)

### Run a preprocessing pipeline

- Flatten images  $8 \times 8 \rightarrow 1 \times 64$
- Normalize (robust scaling)

```
from pipe_utils import flatten_step
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import RobustScaler
digit_pipe = Pipeline([('Flatten', flatten_step),
                      ('Normalize', RobustScaler())])
digit_pipe.fit(img_data)

show_pipe(digit_pipe, img_data)
show_pipe(digit_pipe, img_data, show_hist=True, show_legend=False)
```



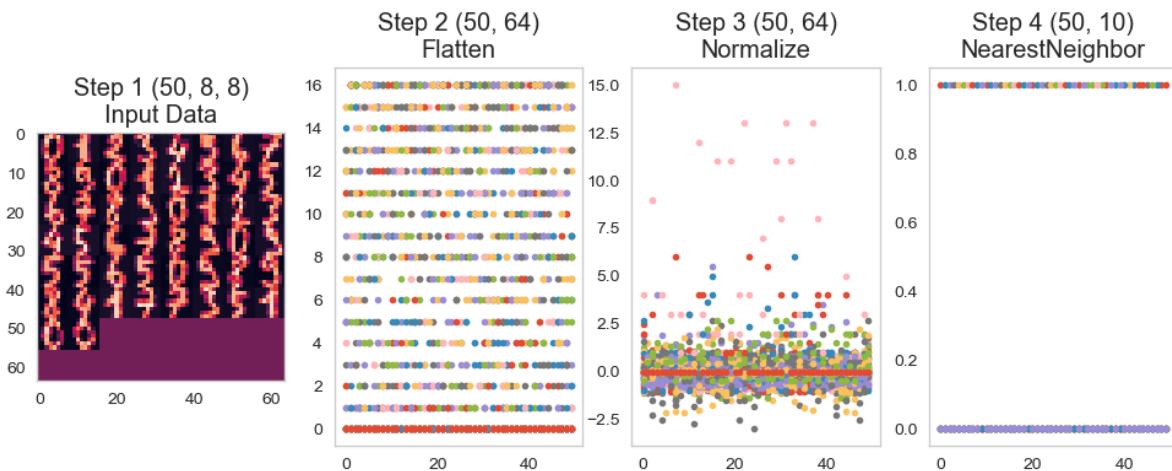
## Add a classifier

- Add a K Nearest Neighbours classifier to the pipeline (K=1)
- Run the fit

```
from sklearn.neighbors import KNeighborsClassifier

digit_class_pipe = Pipeline([('Flatten', flatten_step),
                             ('Normalize', RobustScaler()),
                             ('NearestNeighbor', KNeighborsClassifier(1))])
digit_class_pipe.fit(img_data, digit_id)

show_pipe(digit_class_pipe, img_data, panels_in_row=4)
```



### 0.17.2 Test classifier performance

Let's test with training data

```
from sklearn.metrics import accuracy_score
pred_digit = digit_class_pipe.predict(img_data)

print('{0}% accuracy'.format(100*accuracy_score(digit_id, pred_digit)))
```

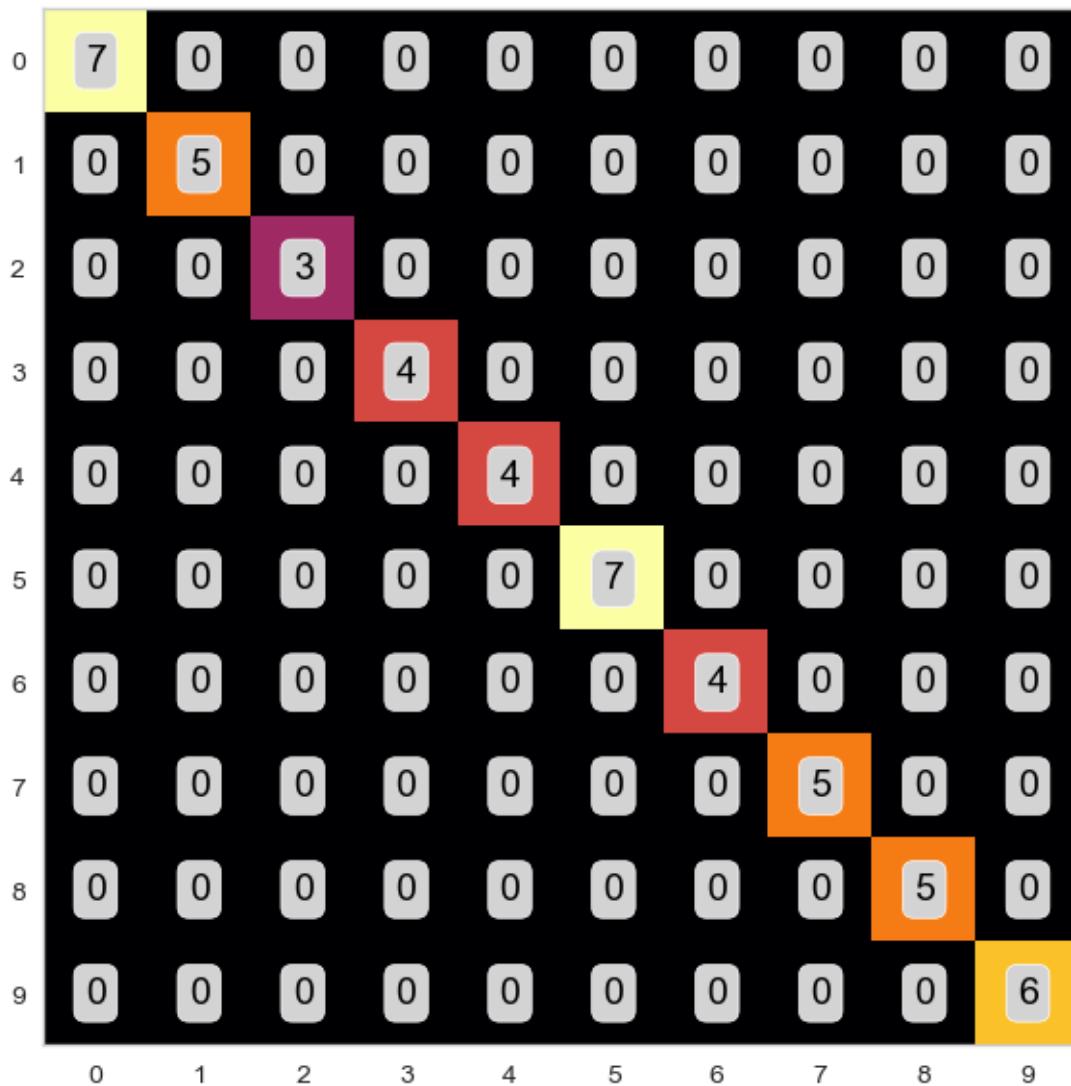
100.0% accuracy

### How about the confusion matrix?

Let's look at the predictions in the confusion matrix.

```
from sklearn.metrics import confusion_matrix
import seaborn as sns
import plotsupport as ps
fig, ax1 = plt.subplots(1, 1, figsize=(8, 7), dpi=100)

ps.heatmap(confusion_matrix(digit_id, pred_digit), precision=0, ax=ax1)
```



### Reporting the classifier output

We can also look at other performance metrics using a `classification_report`

```
from sklearn.metrics import classification_report
print(classification_report(digit_id, pred_digit))
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	7
1	1.00	1.00	1.00	5
2	1.00	1.00	1.00	3
3	1.00	1.00	1.00	4
4	1.00	1.00	1.00	4
5	1.00	1.00	1.00	7
6	1.00	1.00	1.00	4
7	1.00	1.00	1.00	5

(continues on next page)

(continued from previous page)

8	1.00	1.00	1.00	5
9	1.00	1.00	1.00	6
accuracy			1.00	50
macro avg	1.00	1.00	1.00	50
weighted avg	1.00	1.00	1.00	50

### 0.17.3 Wow! We've built an amazing algorithm!

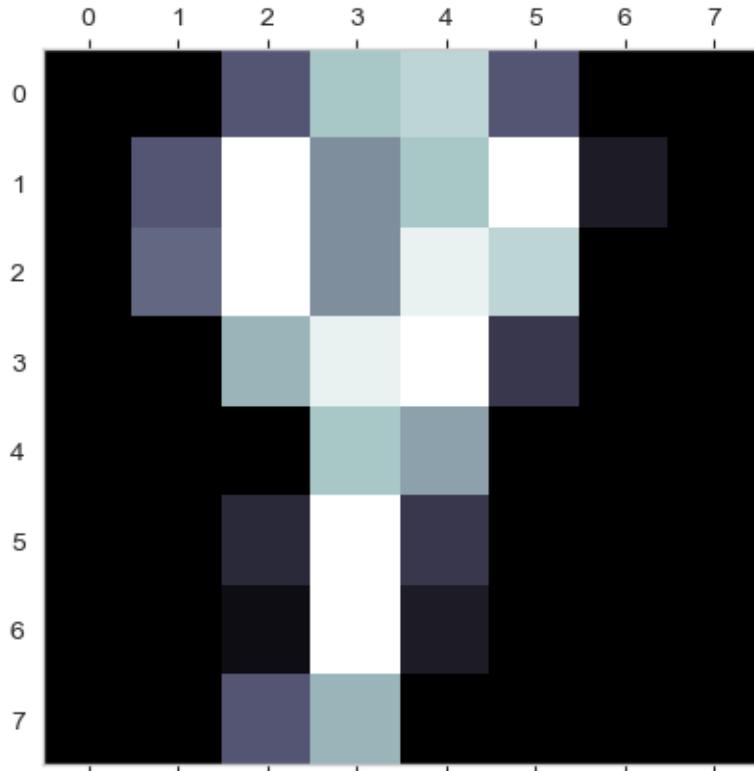
Let's patent it! Call Google!

#### Let's try again

This just too good. Let's make a new try using new, unseen data.

```
test_digit = np.array([[0.,  0.,  6., 12., 13.,  6.,  0.,  0.],
                      [0.,  6., 16.,  9., 12., 16.,  2.,  0.],
                      [0.,  7., 16.,  9., 15., 13.,  0.,  0.],
                      [0.,  0., 11., 15., 16.,  4.,  0.,  0.],
                      [0.,  0.,  0., 12., 10.,  0.,  0.,  0.],
                      [0.,  0.,  3., 16.,  4.,  0.,  0.,  0.],
                      [0.,  0.,  1., 16.,  2.,  0.,  0.,  0.],
                      [0.,  0.,  6., 11.,  0.,  0.,  0.,  0.]])
plt.matshow(test_digit[0], cmap='bone')
print('Prediction:', digit_class_pipe.predict(test_digit))
print('Real Value:', 9)
```

```
Prediction: [7]
Real Value: 9
```



### 0.17.4 Training, Validation, and Testing

Avoid the training “crime” in ML

<https://www.kdnuggets.com/2017/08/dataiku-predictive-model-holdout-cross-validation.html>

## 0.18 Regression using a pipeline

For regression, we can see

- it is very similarly to a classification
  - Predicts the category
- instead of trying to output discrete classes we can output on a continuous scale.
  - Predicts the **actual decimal number**.

```
from sklearn.datasets import load_digits
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import make_blobs
from pipe_utils import show_pipe, flatten_step
%matplotlib inline
```

### 0.18.1 Load the digits data again

This time we load

- 50 image/label pairs for the training
- 450 images/label pairs for the validation

```
digit_ds = load_digits(return_X_y=False)

img_data = digit_ds.images[:50]
digit_id = digit_ds.target[:50]

valid_data = digit_ds.images[50:500]
valid_id = digit_ds.target[50:500]
```

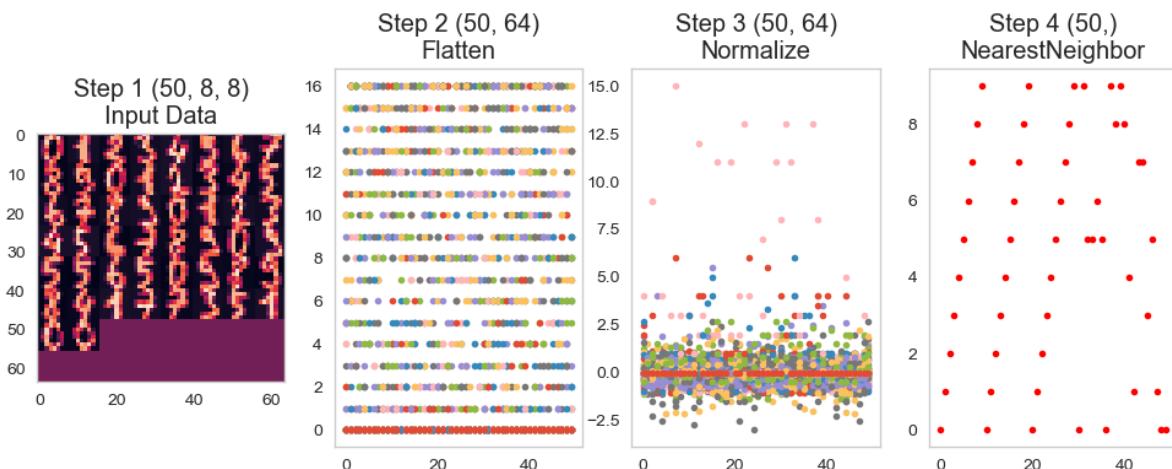
### 0.18.2 Run a KNeighbors regression

In this pipeline we will make a regression based on K-Neighbors.

```
from sklearn.neighbors import KNeighborsRegressor

digit_regress_pipe = Pipeline([('Flatten', flatten_step),
                               ('Normalize', RobustScaler()),
                               ('NearestNeighbor', KNeighborsRegressor(1))])
digit_regress_pipe.fit(img_data, digit_id)

show_pipe(digit_regress_pipe, img_data, panels_in_row=4)
```



## 0.19 Assessment of multi-category data

We can't use accuracy, ROC, precision, recall or any of these factors anymore since we don't have binary / true-or-false conditions we are trying to predict. We know have to go back to some of the initial metrics we covered in the first lectures.

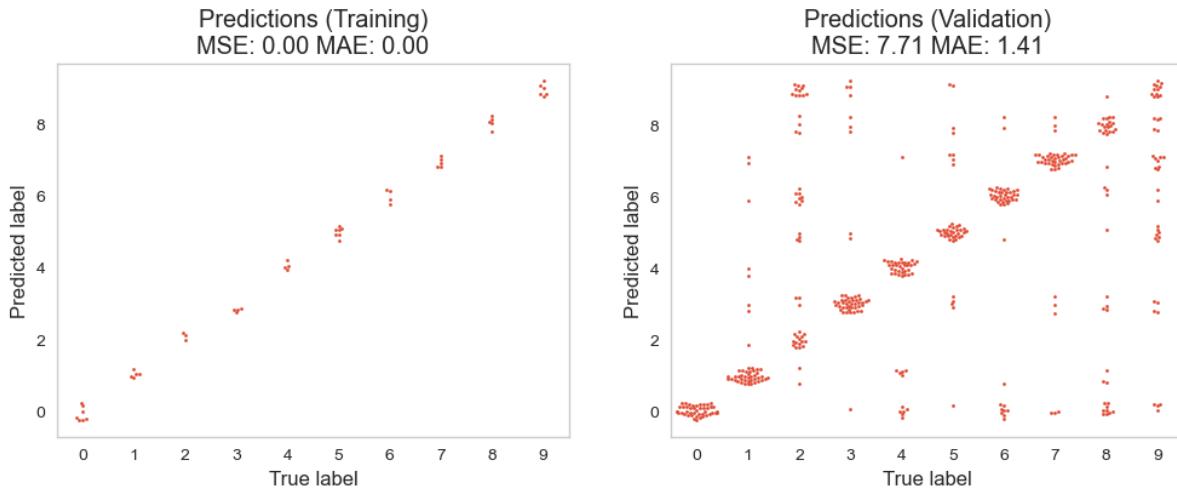
$$MSE = \frac{1}{N} \sum (y_{predicted} - y_{actual})^2$$

$$MAE = \frac{1}{N} \sum |y_{predicted} - y_{actual}|$$

```
fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 4), dpi=100)
pred_train = digit_regress_pipe.predict(img_data)
jitter = lambda x: x+0.25*np.random.uniform(-1, 1, size=x.shape)
sns.swarmplot(x=digit_id, y=jitter(pred_train), ax=ax1, size=2)
ax1.set_title('Predictions (Training)\nMSE: %2.2f MAE: %2.2f' % (np.mean(np.
    np.abs(pred_train-digit_id))), np.mean(np.abs(pred_
    train-digit_id))))
ax1.set_xlabel('True label'), ax1.set_ylabel('Predicted label')

pred_valid = digit_regress_pipe.predict(valid_data)
sns.swarmplot(x=valid_id, y=jitter(pred_valid), ax=ax2, size=2)
ax2.set_title('Predictions (Validation)\nMSE: %2.2f MAE: %2.2f' % (np.mean(np.
    np.abs(pred_valid-valid_id))), np.mean(np.
    abs(pred_valid-valid_id)));
ax2.set_xlabel('True label'), ax2.set_ylabel('Predicted label');
```

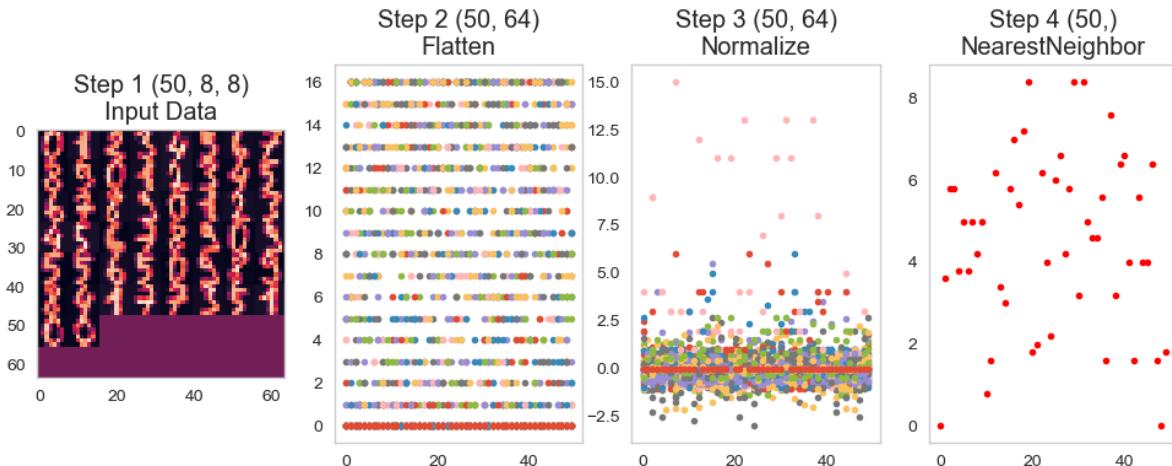
```
/Users/kaestner/anaconda3/lib/python3.11/site-packages/seaborn/_oldcore.py:1119:_
  FutureWarning: use_inf_as_na option is deprecated and will be removed in a_
  future version. Convert inf values to NaN before operating instead.
  with pd.option_context('mode.use_inf_as_na', True):
/Users/kaestner/anaconda3/lib/python3.11/site-packages/seaborn/_oldcore.py:1119:_
  FutureWarning: use_inf_as_na option is deprecated and will be removed in a_
  future version. Convert inf values to NaN before operating instead.
  with pd.option_context('mode.use_inf_as_na', True):
/Users/kaestner/anaconda3/lib/python3.11/site-packages/seaborn/_oldcore.py:1119:_
  FutureWarning: use_inf_as_na option is deprecated and will be removed in a_
  future version. Convert inf values to NaN before operating instead.
  with pd.option_context('mode.use_inf_as_na', True):
/Users/kaestner/anaconda3/lib/python3.11/site-packages/seaborn/_oldcore.py:1119:_
  FutureWarning: use_inf_as_na option is deprecated and will be removed in a_
  future version. Convert inf values to NaN before operating instead.
  with pd.option_context('mode.use_inf_as_na', True):
/Users/kaestner/anaconda3/lib/python3.11/site-packages/seaborn/_oldcore.py:1119:_
  FutureWarning: use_inf_as_na option is deprecated and will be removed in a_
  future version. Convert inf values to NaN before operating instead.
  with pd.option_context('mode.use_inf_as_na', True):
/Users/kaestner/anaconda3/lib/python3.11/site-packages/seaborn/categorical.
  py:3544: UserWarning: 9.1% of the points cannot be placed; you may want to_
  decrease the size of the markers or use stripplot.
  warnings.warn(msg, UserWarning)
```



### 0.19.1 Increasing neighbor count

```
digit_regress_pipe = Pipeline([('Flatten', flatten_step), ('Normalize', RobustScaler()), ('NearestNeighbor', KNeighborsRegressor(5))])
digit_regress_pipe.fit(img_data, digit_id)

show_pipe(digit_regress_pipe, img_data, panels_in_row=4)
```



```
fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 4))
pred_train = digit_regress_pipe.predict(img_data)

sns.swarmplot(x=digit_id, y=jitter(pred_train), ax=ax1, size=2)
ax1.set_title('Predictions (Training)\nMSE: %2.2f MAE: %2.2f' % (np.mean(np.square(pred_train-digit_id)), np.mean(np.abs(pred_train-digit_id))))
ax1.set_xlabel('True label'), ax1.set_ylabel('Predicted label')

pred_valid = digit_regress_pipe.predict(valid_data)
sns.swarmplot(x=valid_id, y=jitter(pred_valid), ax=ax2, size=2)
```

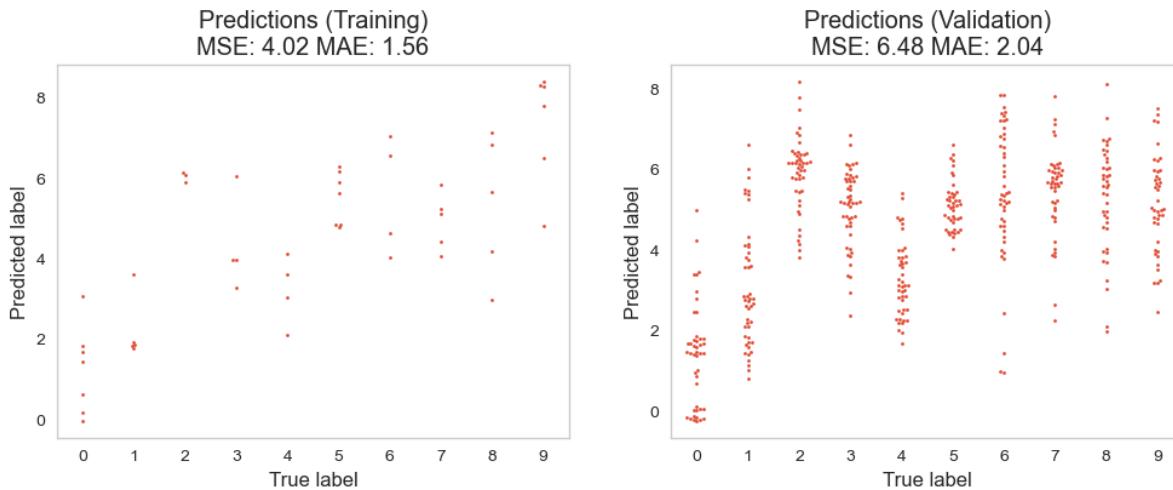
(continues on next page)

(continued from previous page)

```
ax2.set_title('Predictions (Validation)\nMSE: %2.2f MAE: %2.2f' % (np.mean(np.
    np.square(pred_valid-valid_id)),
    np.mean(np.
    np.abs(pred_valid-valid_id))));

ax2.set_xlabel('True label'), ax2.set_ylabel('Predicted label');
```

```
/Users/kaestner/anaconda3/lib/python3.11/site-packages/seaborn/_oldcore.py:1119:_
  FutureWarning: use_inf_as_na option is deprecated and will be removed in a_
  future version. Convert inf values to NaN before operating instead.
  with pd.option_context('mode.use_inf_as_na', True):
/Users/kaestner/anaconda3/lib/python3.11/site-packages/seaborn/_oldcore.py:1119:_
  FutureWarning: use_inf_as_na option is deprecated and will be removed in a_
  future version. Convert inf values to NaN before operating instead.
  with pd.option_context('mode.use_inf_as_na', True):
/Users/kaestner/anaconda3/lib/python3.11/site-packages/seaborn/_oldcore.py:1119:_
  FutureWarning: use_inf_as_na option is deprecated and will be removed in a_
  future version. Convert inf values to NaN before operating instead.
  with pd.option_context('mode.use_inf_as_na', True):
/Users/kaestner/anaconda3/lib/python3.11/site-packages/seaborn/_oldcore.py:1119:_
  FutureWarning: use_inf_as_na option is deprecated and will be removed in a_
  future version. Convert inf values to NaN before operating instead.
  with pd.option_context('mode.use_inf_as_na', True):
/Users/kaestner/anaconda3/lib/python3.11/site-packages/seaborn/_oldcore.py:1119:_
  FutureWarning: use_inf_as_na option is deprecated and will be removed in a_
  future version. Convert inf values to NaN before operating instead.
  with pd.option_context('mode.use_inf_as_na', True):
```



## 0.20 Segmentation (Pixel Classification)

**Previously** Predict something based on the information in the image

- Single class (classification)
- Values (regression)

**Segmentation** Now we want to change problem:

- instead of assigning a single class for each image,
- we want a class or value for each pixel.

This requires that we restructure the problem.

### 0.20.1 Where segmentation fails: Mitochondria Segmentation in EM

- The mitochondria are visible and humans easily spot them
- Other structures have same gray levels
- SNR is not ideal

- 
- A simple threshold is insufficient to finding the mitochondria structures
  - Other filtering techniques are unlikely to magically fix this problem

### 0.20.2 Let's try some methods to segment the mitochondria image

#### 1. Decision trees

- DecisionTreeRegressor
- DecisionTreeRegressor with position

#### 2. Random forests

- Random forest + KMeans
- Random forest + polynomials
- Random forest + Filters

#### 3. Linear regression

- Neighborhood

#### 4. Nearest neighbor

- KNeighborsRegressor

#### 5. U-Net

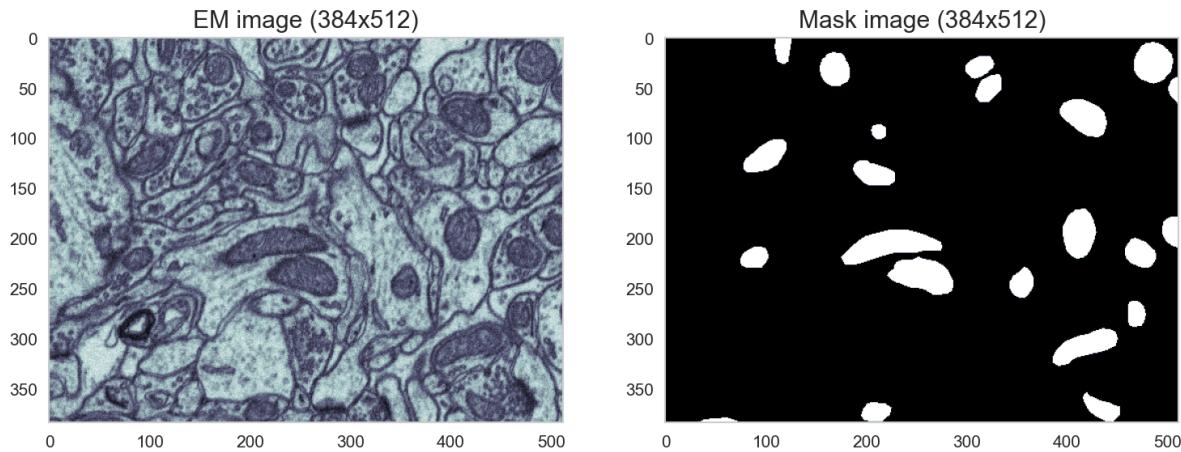
### 0.20.3 Preparing the mitochondria image and mask

First we need image data for:

- Training
- Validation

```
cell_img = (imread("data/em_image.png") [:, :, ::2]) / 255.0
cell_seg = imread("data/em_image_seg.png",
                  as_gray=True) [:, :, ::2] > 0
np.random.seed(2018)
```

```
fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 8), dpi=150)
ax1.imshow(cell_img, cmap='bone'); ax1.set_title("EM image ({0}x{1})".format(cell_img.shape[0], cell_img.shape[1]))
ax2.imshow(cell_seg, cmap='bone'); ax2.set_title("Mask image ({0}x{1})".format(cell_seg.shape[0], cell_seg.shape[1]));
```



### Training and validation data

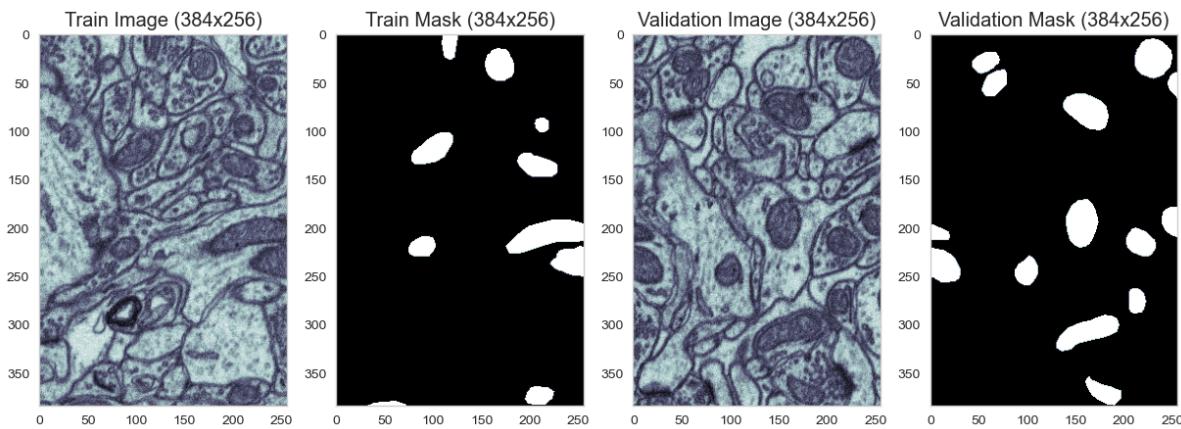
We only have one image available...

**Solution:** Split it into two parts!

```
train_img, valid_img = cell_img[:, :256], cell_img[:, 256:]
train_mask, valid_mask = cell_seg[:, :256], cell_seg[:, 256:]
```

```
fig, ((ax1, ax2, ax3, ax4)) = plt.subplots(1, 4, figsize=(15, 5), dpi=100)
ax1.imshow(train_img, cmap='bone'); ax1.set_title('Train Image ({0}x{1})'.
    format(train_img.shape[0], train_img.shape[1]))
ax2.imshow(train_mask, cmap='bone'); ax2.set_title('Train Mask ({0}x{1})'.
    format(train_mask.shape[0], train_mask.shape[1]))

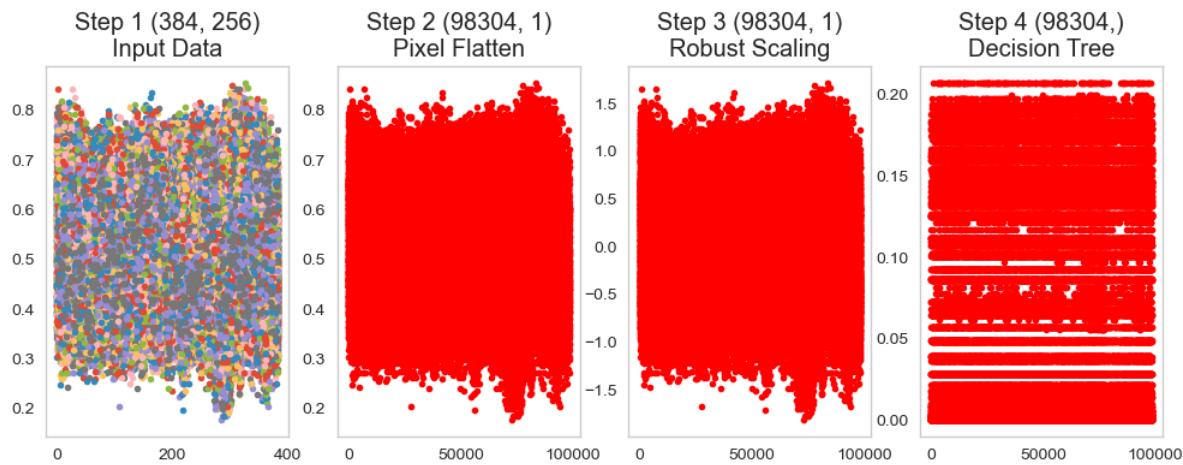
ax3.imshow(valid_img, cmap='bone'); ax3.set_title('Validation Image ({0}x{1})'.
    format(valid_img.shape[0], valid_img.shape[1]))
ax4.imshow(valid_mask, cmap='bone'); ax4.set_title('Validation Mask ({0}x{1})'.
    format(valid_mask.shape[0], valid_mask.shape[1]))
```



### 0.20.4 Try a Regression tree

```
from pipe_utils import px_flatten_func, fit_img_pipe
px_flatten_step = FunctionTransformer(px_flatten_func, validate=False)
rf_seg_model = Pipeline([('Pixel Flatten', px_flatten_step),
                        ('Robust Scaling', RobustScaler()),
                        ('Decision Tree', DecisionTreeRegressor())
                       ])

pred_func = fit_img_pipe(rf_seg_model, train_img, train_mask)
show_pipe(rf_seg_model, train_img, panels_in_row=4)
show_tree(rf_seg_model.steps[-1][1]);
```



### Segmentation results from the regression tree

```
fig, ((ax1, ax5, ax2), (ax3, ax6, ax4)) = plt.subplots(2, 3, figsize=(12, 8), dpi=150)
ax1.imshow(train_img, cmap='bone')
ax1.set_title('Train Image')

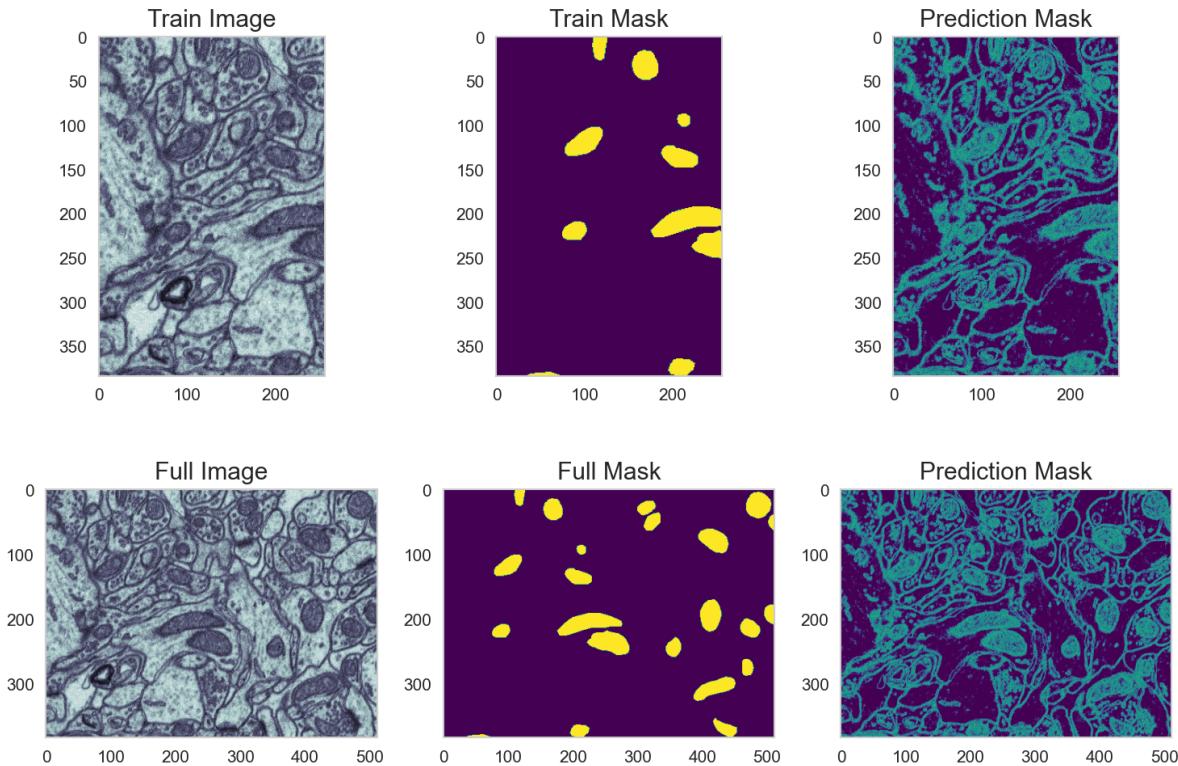
ax5.imshow(train_mask, cmap='viridis'); ax5.set_title('Train Mask')

ax2.imshow(pred_func(train_img)[:, :, 0], cmap='viridis', vmin=0, vmax=0.3); ax2.set_
           title('Prediction Mask')

ax3.imshow(cell_img, cmap='bone'); ax3.set_title('Full Image')

ax6.imshow(cell_seg, cmap='viridis'); ax6.set_title('Full Mask')

ax4.imshow(pred_func(cell_img)[:, :, 0], cmap='viridis', vmin=0, vmax=0.3); ax4.set_
           title('Prediction Mask');
```

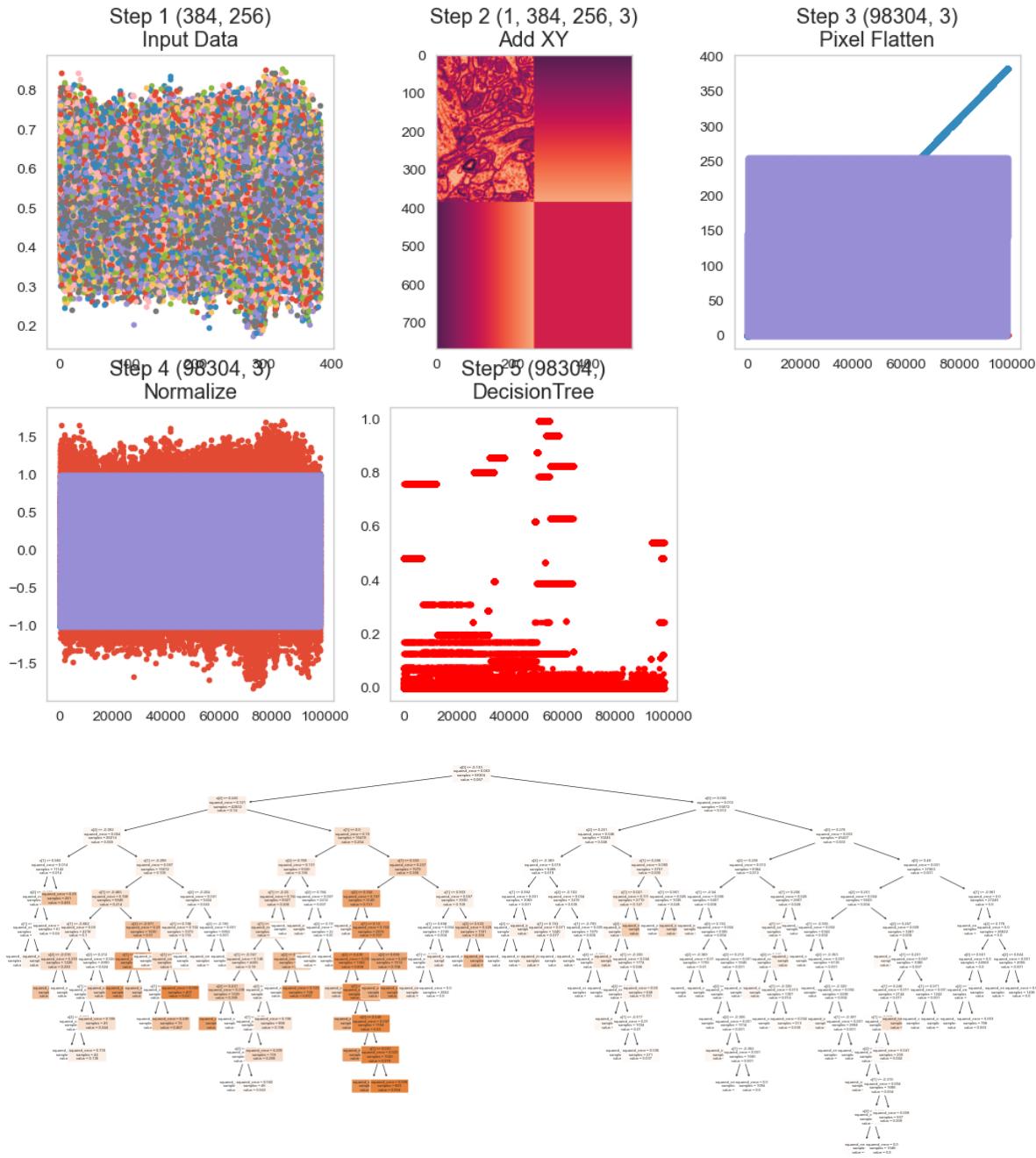


### 0.20.5 Regression tree with *position information*

```
from pipe_utils import xy_step

rf_xyseg_model = Pipeline([('Add XY', xy_step),
                           ('Pixel Flatten', px_flatten_step),
                           ('Normalize', RobustScaler()),
                           ('DecisionTree', DecisionTreeRegressor(
                               min_samples_split=1000))
                          ])

pred_func = fit_img_pipe(rf_xyseg_model, train_img, train_mask)
show_pipe(rf_xyseg_model, train_img)
plt.figure(figsize=(15, 6))
tree.plot_tree(rf_xyseg_model.steps[-1][1], fontsize=3, filled=True, rounded=True);
```



Did the segmentation performance improve?

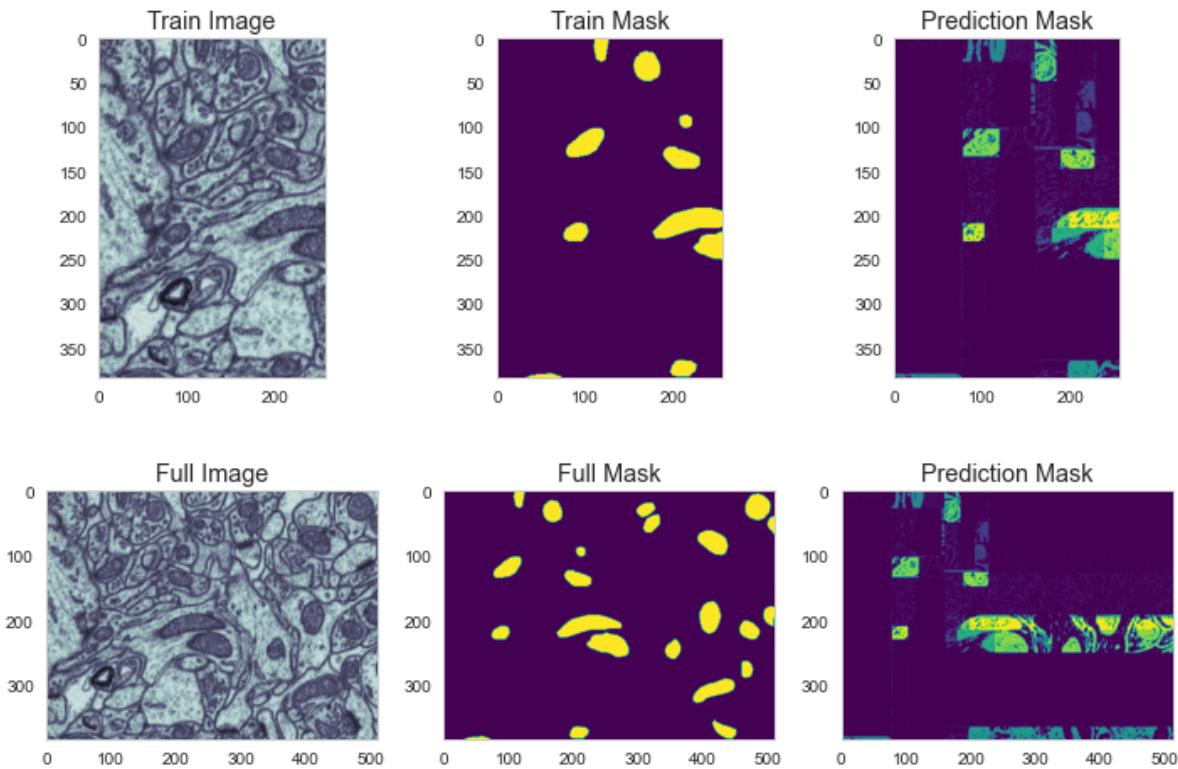
```

fig, ((ax1, ax5, ax2), (ax3, ax6, ax4)) = plt.subplots(2, 3, figsize=(12, 8), dpi=72)
ax1.imshow(train_img, cmap='bone'); ax1.set_title('Train Image')
ax5.imshow(train_mask, cmap='viridis'); ax5.set_title('Train Mask')
ax2.imshow(pred_func(train_img)[:, :, 0], cmap='viridis', vmin=0, vmax=1); ax2.set_
    title('Prediction Mask')
ax3.imshow(cell_img, cmap='bone'); ax3.set_title('Full Image')
ax6.imshow(cell_seg, cmap='viridis'); ax6.set_title('Full Mask')
    
```

(continues on next page)

(continued from previous page)

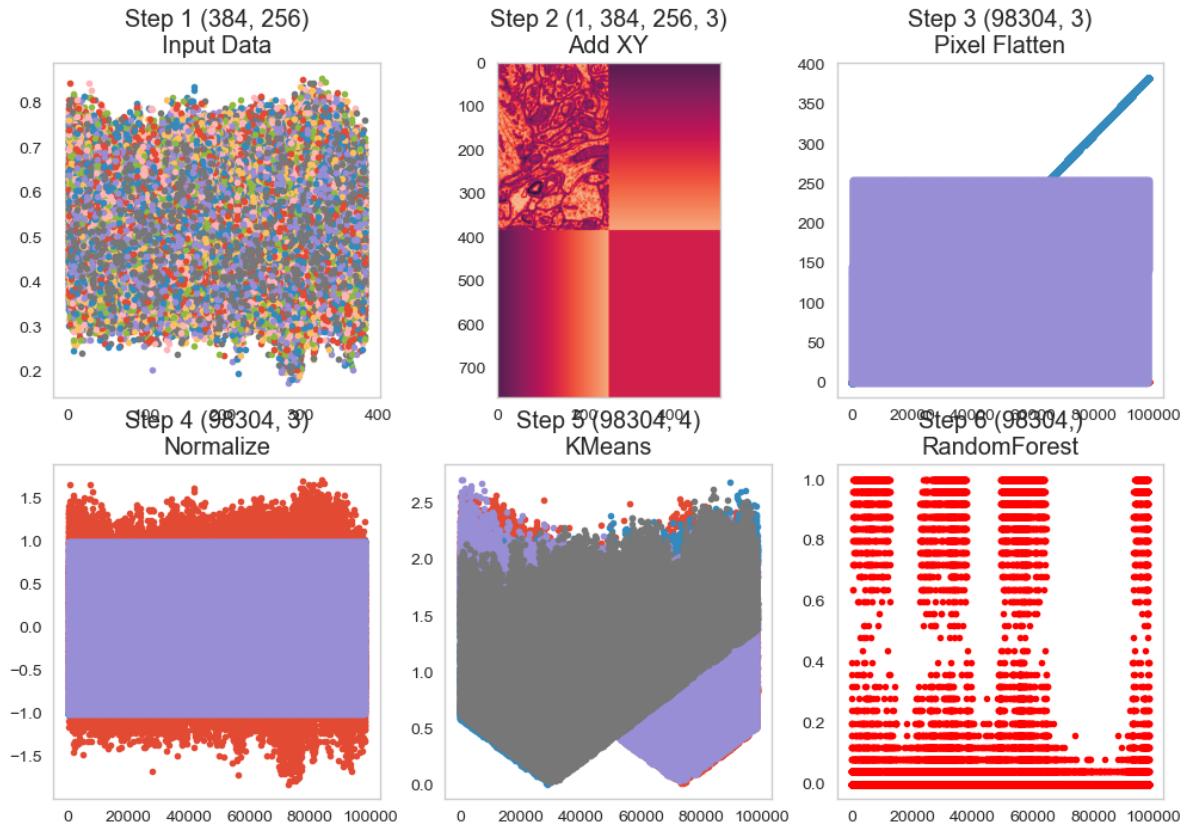
```
ax4.imshow(pred_func(cell_img)[:, :, 0], cmap='viridis', vmin=0, vmax=1); ax4.set_
    title('Prediction Mask');
```



### 0.20.6 Combine K-Means and Random Forest Regression

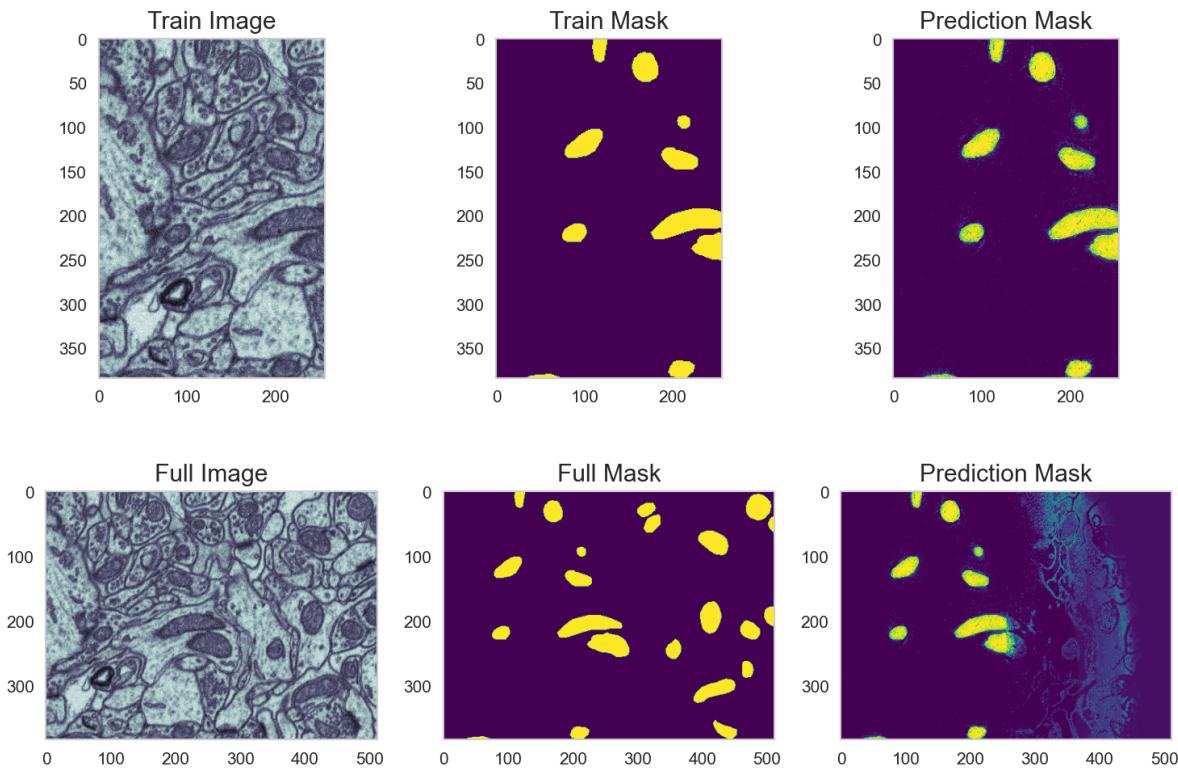
```
from sklearn.cluster import KMeans
rf_xyseg_k_model = Pipeline([('Add XY', xy_step),
                             ('Pixel Flatten', px_flatten_step),
                             ('Normalize', RobustScaler()),
                             ('KMeans', KMeans(4, n_init='auto')),
                             ('RandomForest', RandomForestRegressor(n_estimators=25))
                            ])

pred_func = fit_img_pipe(rf_xyseg_k_model, train_img, train_mask)
show_pipe(rf_xyseg_k_model, train_img)
```



### Did it improve now?

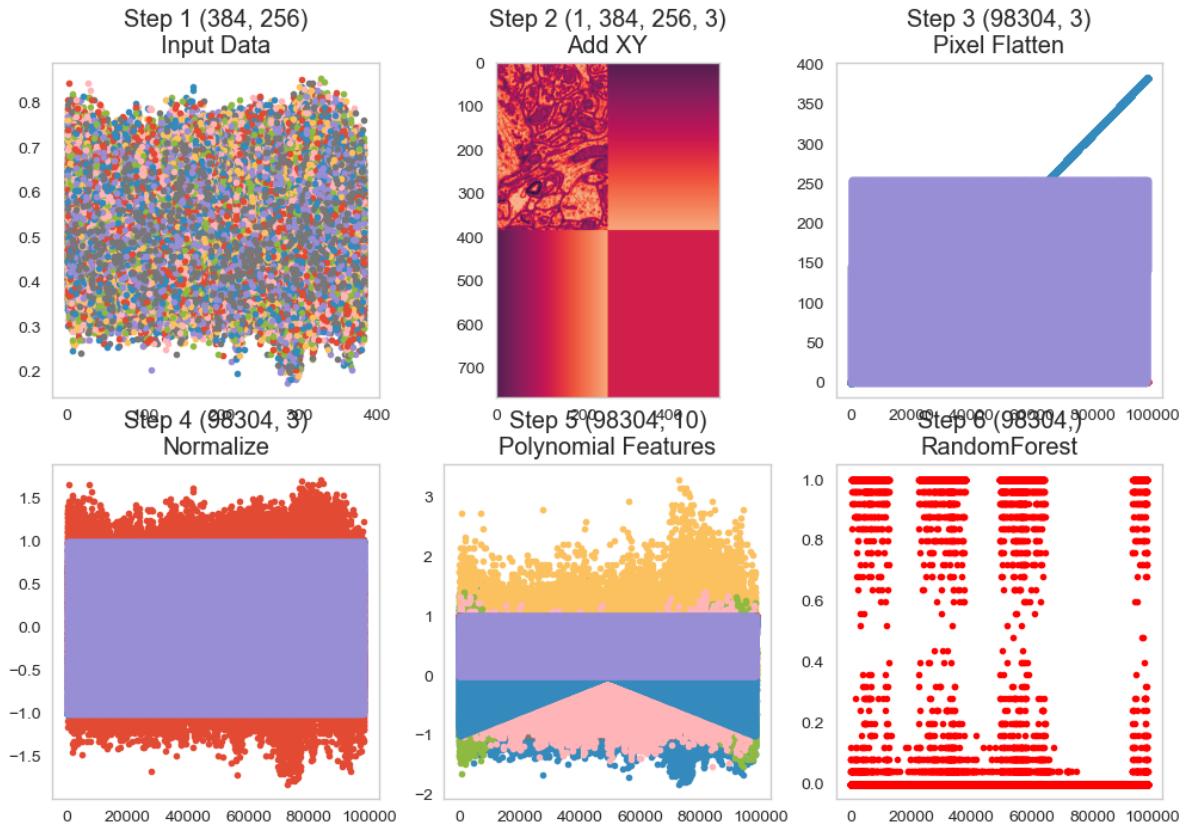
```
fig, ((ax1, ax5, ax2), (ax3, ax6, ax4)) = plt.subplots(2, 3, figsize=(12, 8), dpi=150)
ax1.imshow(train_img, cmap='bone'); ax1.set_title('Train Image')
ax5.imshow(train_mask, cmap='viridis'); ax5.set_title('Train Mask')
ax2.imshow(pred_func(train_img)[:, :, 0], cmap='viridis', vmin=0, vmax=1); ax2.set_
    title('Prediction Mask')
ax3.imshow(cell_img, cmap='bone'); ax3.set_title('Full Image')
ax6.imshow(cell_seg, cmap='viridis'); ax6.set_title('Full Mask')
ax4.imshow(pred_func(cell_img)[:, :, 0], cmap='viridis', vmin=0, vmax=1); ax4.set_
    title('Prediction Mask');
```



### 0.20.7 Trying polynomial features

```
from sklearn.preprocessing import PolynomialFeatures
from pipe_utils import add_xy_coord
xy_step = FunctionTransformer(add_xy_coord, validate=False)
rf_xyseg_py_model = Pipeline([('Add XY', xy_step),
                             ('Pixel Flatten', px_flatten_step),
                             ('Normalize', RobustScaler()),
                             ('Polynomial Features', PolynomialFeatures(2)),
                             ('RandomForest', RandomForestRegressor(n_estimators=25))])

pred_func = fit_img_pipe(rf_xyseg_py_model, train_img, train_mask)
show_pipe(rf_xyseg_py_model, train_img)
```



### What happens with this complicated pipeline?

We...

- Added XY position
- Normalized
- Added polynomial features
- Used a random forest regressor with 25 trees

```
fig, ((ax1, ax5, ax2), (ax3, ax6, ax4)) = plt.subplots(2, 3, figsize=(12, 8), dpi=150)
ax1.imshow(train_img, cmap='bone')
ax1.set_title('Train Image')

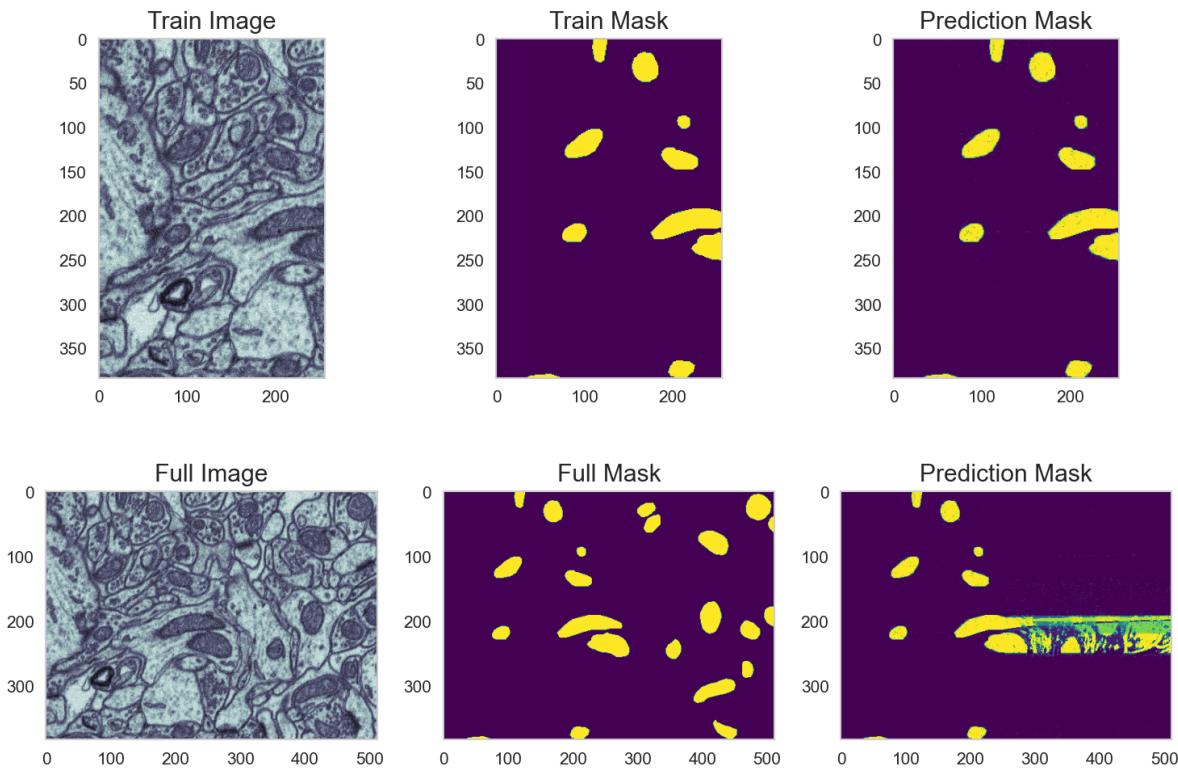
ax5.imshow(train_mask, cmap='viridis'); ax5.set_title('Train Mask')

ax2.imshow(pred_func(train_img)[:, :, 0], cmap='viridis', vmin=0, vmax=1); ax2.set_
    title('Prediction Mask')

ax3.imshow(cell_img, cmap='bone'); ax3.set_title('Full Image')

ax6.imshow(cell_seg, cmap='viridis'); ax6.set_title('Full Mask')

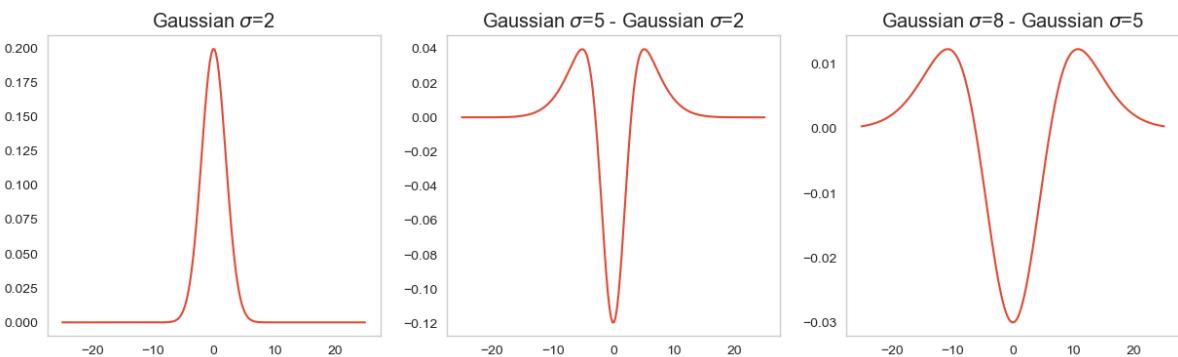
ax4.imshow(pred_func(cell_img)[:, :, 0], cmap='viridis', vmin=0, vmax=1); ax4.set_
    title('Prediction Mask');
```



### 0.20.8 Adding smarter features

Here we add images with weighted neighborhood information using filters based on Gaussians

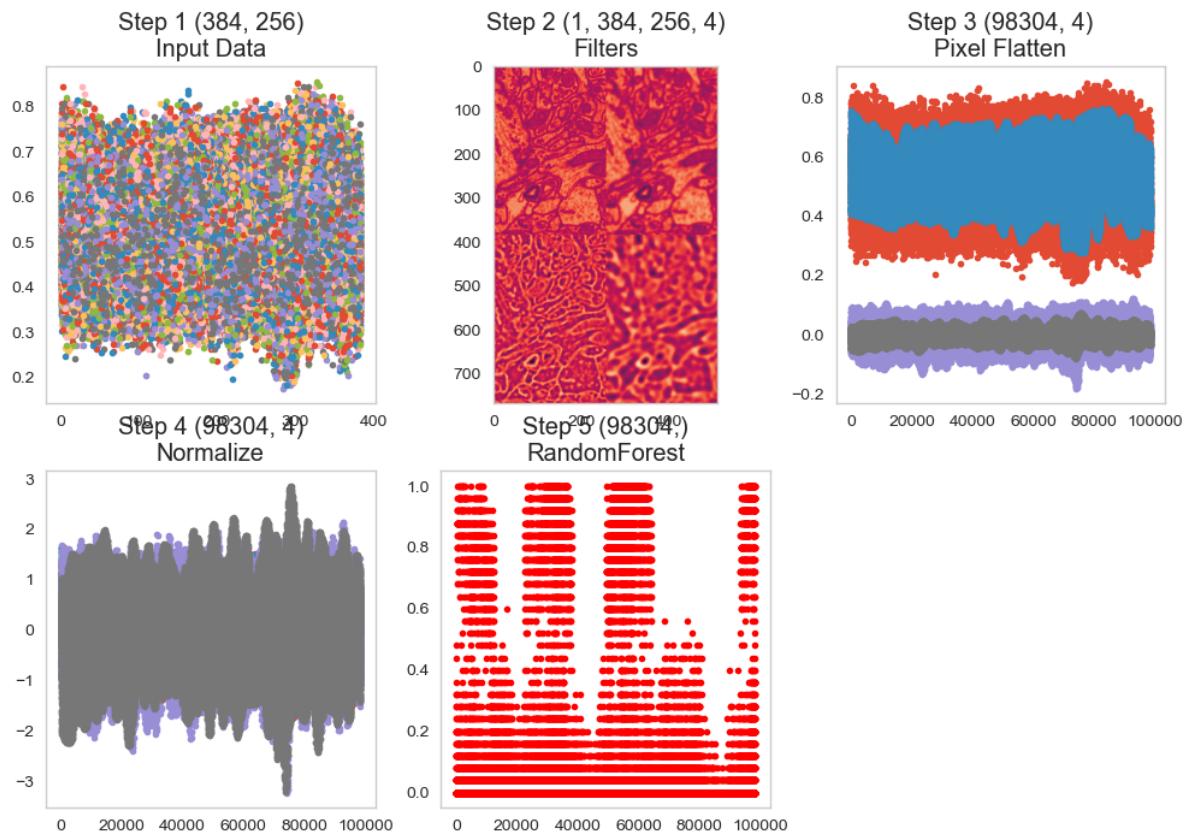
```
import scipy.stats as stats
x=np.linspace(-25,25,200)
fig, (ax1,ax2,ax3) = plt.subplots(1,3,figsize=[15,4])
ax1.plot(x,stats.norm.pdf(x,0,2)); ax1.set_title("Gaussian $\sigma=2");
ax2.plot(x,stats.norm.pdf(x,0,5)-stats.norm.pdf(x,0,2)), ax2.set_title("Gaussian $\sigma=5 - Gaussian $\sigma=2");
ax3.plot(x,stats.norm.pdf(x,0,8)-stats.norm.pdf(x,0,5)), ax3.set_title("Gaussian $\sigma=8 - Gaussian $\sigma=5");
```



## Pipeline with filters

```
from pipe_utils import filter_step
rf_filterseg_model = Pipeline([('Filters', filter_step),
                               ('Pixel Flatten', px_flatten_step),
                               ('Normalize', RobustScaler()),
                               ('RandomForest', RandomForestRegressor(n_
                               estimators=25))
                             ])

pred_func = fit_img_pipe(rf_filterseg_model, train_img, train_mask)
show_pipe(rf_filterseg_model, train_img)
```



## Results with features based on filters

```
fig, ((ax1, ax5, ax2), (ax3, ax6, ax4)) = plt.subplots( 2, 3, figsize=(12, 8),_
                                         dpi=150)
ax1.imshow(train_img, cmap='bone'); ax1.set_title('Train Image')

ax5.imshow(train_mask, cmap='viridis'); ax5.set_title('Train Mask')

ax2.imshow(pred_func(train_img)[:, :, 0], cmap='viridis', vmin=0, vmax=1); ax2.set_-
title('Prediction Mask')

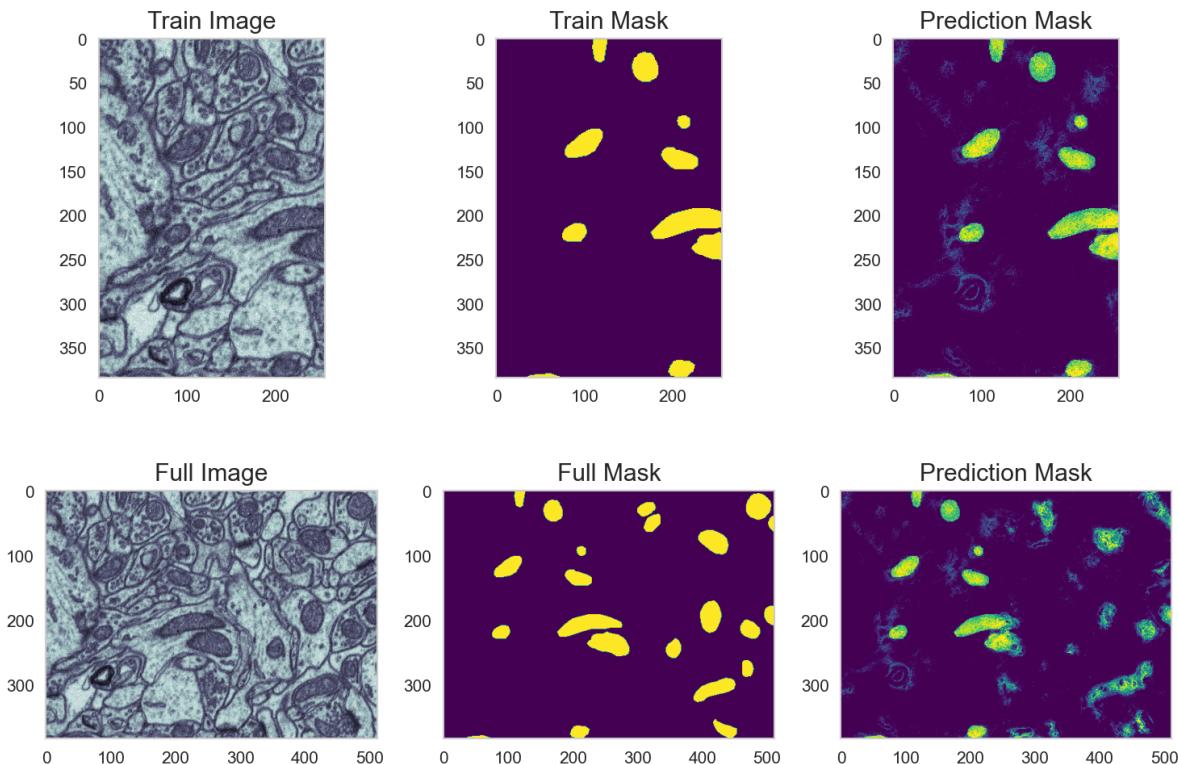
ax3.imshow(cell_img, cmap='bone'); ax3.set_title('Full Image')
```

(continues on next page)

(continued from previous page)

```
ax6.imshow(cell_seg, cmap='viridis'); ax6.set_title('Full Mask')

ax4.imshow(pred_func(cell_img)[:, :, 0], cmap='viridis', vmin=0, vmax=1); ax4.set_
title('Prediction Mask');
```



### 0.20.9 Using the Neighborhood

We can also include the whole neighborhood

- shifting the image in x and y by  $\pm 1$  pixel.
- Gives nine feature images

For the first example we will then use **linear regression** so we can see the exact coefficients that result.

#### Some code to create the neighborhood

```
from sklearn.preprocessing import FunctionTransformer

def add_neighborhood(in_x, x_steps=3, y_steps=3):
    if len(in_x.shape) == 2: x = np.expand_dims(np.expand_dims(in_x, 0), -1)
    elif len(in_x.shape) == 3: x = np.expand_dims(in_x, -1)
    elif len(in_x.shape) == 4: x = in_x
    else:
        raise ValueError('Cannot work with images with dimensions {}'.format(in_x.
shape))
```

(continues on next page)

(continued from previous page)

```

n_img, x_dim, y_dim, c_dim = x.shape
out_imgs = []
for i in range(-x_steps, x_steps+1):
    for j in range(-y_steps, y_steps+1):
        out_imgs += [np.roll(np.roll(x, axis=1, shift=i), axis=2, shift=j)]
return np.concatenate(out_imgs, -1)

def neighbor_step(x_steps=3, y_steps=3):
    return FunctionTransformer(
        lambda x: add_neighborhood(x, x_steps, y_steps),
        validate=False)

```

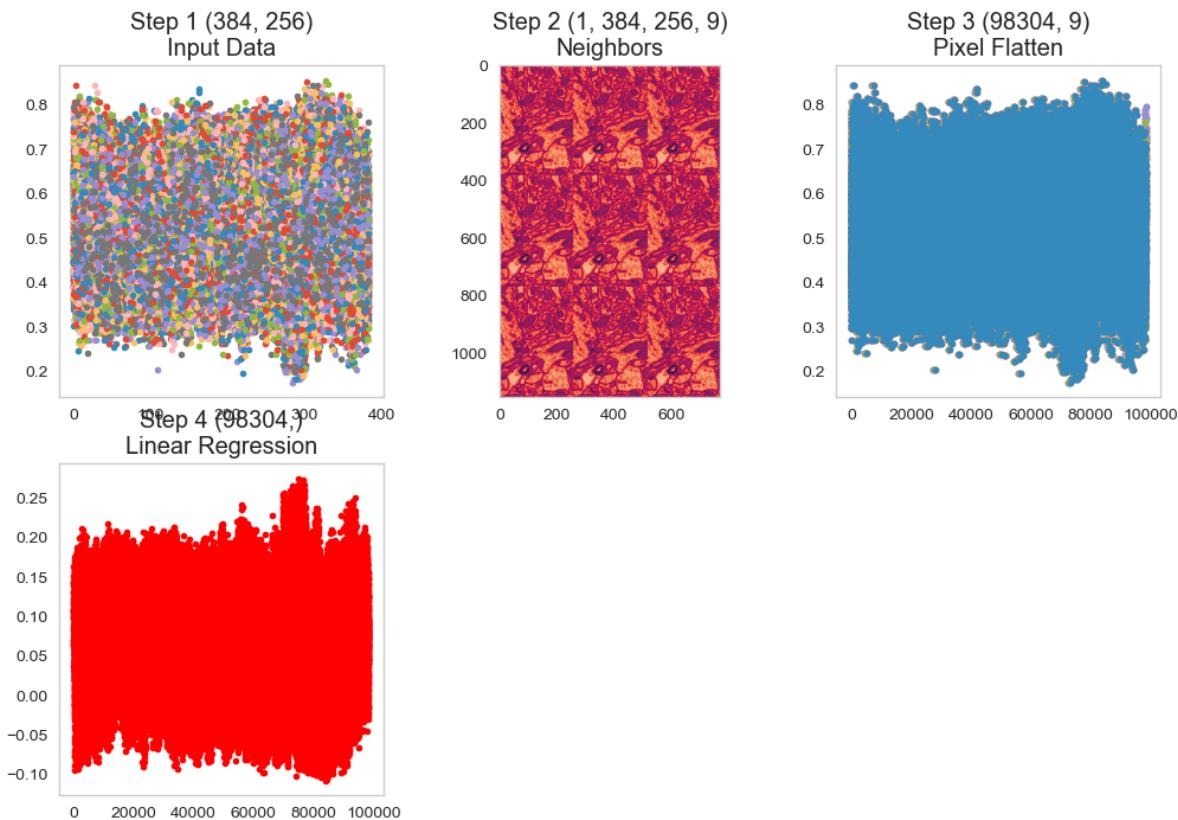
## Pipeline with neighborhood features

```

from sklearn.linear_model import LinearRegression
linreg_neighorseg_model = Pipeline([('Neighbors', neighbor_step(1, 1)),
                                    ('Pixel Flatten', px_flatten_step),
                                    ('Linear Regression', LinearRegression())
                                   ])

pred_func = fit_img_pipe(linreg_neighorseg_model, train_img, train_mask)
show_pipe(linreg_neighorseg_model, train_img)

```



### Result of neighborhood regression

```
fig, ((ax1, ax5, ax2), (ax3, ax6, ax4)) = plt.subplots(2, 3, figsize=(12, 8), dpi=150)
ax1.imshow(train_img, cmap='bone') ; ax1.set_title('Train Image')

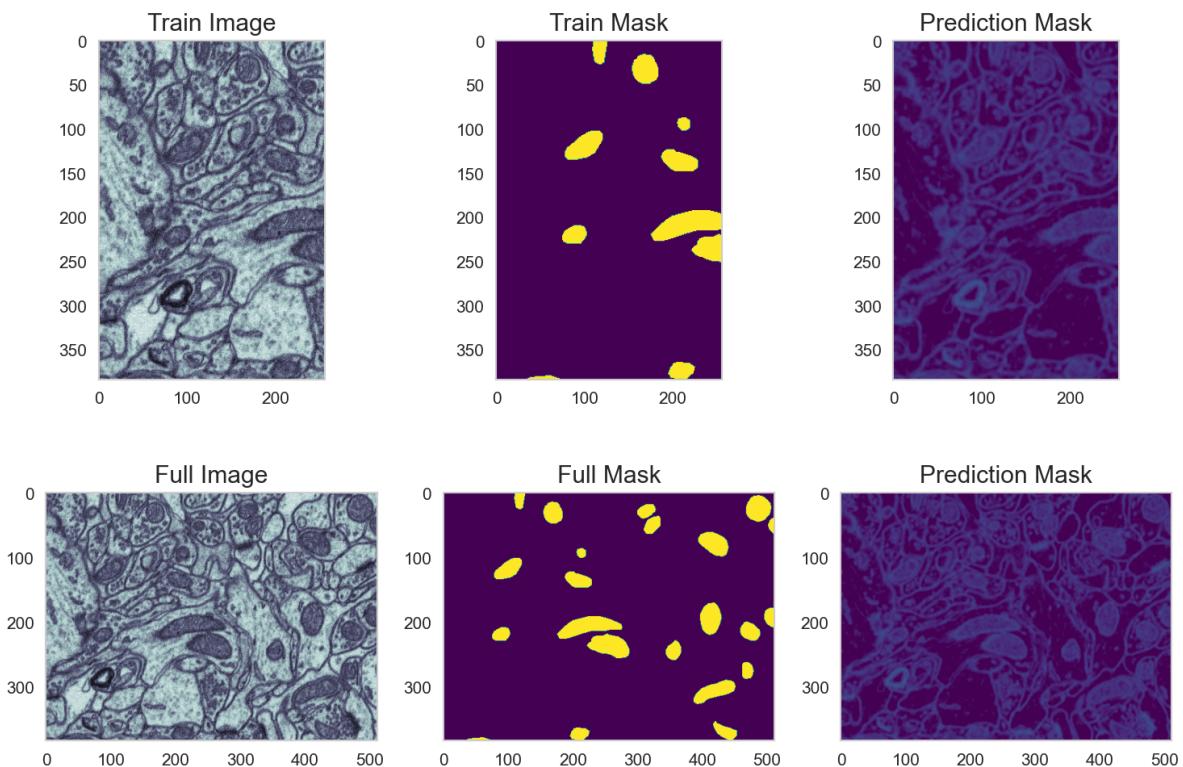
ax5.imshow(train_mask, cmap='viridis'); ax5.set_title('Train Mask')

ax2.imshow(pred_func(train_img)[:, :, 0], cmap='viridis', vmin=0, vmax=1); ax2.set_
    title('Prediction Mask')

ax3.imshow(cell_img, cmap='bone') ; ax3.set_title('Full Image')

ax6.imshow(cell_seg, cmap='viridis') ; ax6.set_title('Full Mask')

ax4.imshow(pred_func(cell_img)[:, :, 0], cmap='viridis', vmin=0, vmax=1); ax4.set_
    title('Prediction Mask');
```



### Why Linear Regression?

We choose linear regression so we could get easily understood coefficients.

The model fits  $\vec{m}$  and  $b$  to the  $\vec{x}_{i,j}$  points in the image  $I(i,j)$  to match the  $y_{i,j}$  output in the segmentation as closely as possible  $y_{i,j} = \vec{m} \cdot \vec{x}_{i,j} + b$  For a 3x3 case this looks like  $\vec{x}_{i,j} = [I(i-1, j-1), I(i-1, j), I(i-1, j+1) \dots I(i+1, j-1), I(i+1, j), I(i+1, j+1)]$

```
m = linreg_neighborhoodseg_model.steps[-1][1].coef_
b = linreg_neighborhoodseg_model.steps[-1][1].intercept_
```

(continues on next page)

(continued from previous page)

```
print('M: [{:0.4}, {:0.4}, {:0.4}, {:0.4}, \033[1m{:0.4}\033[0m, {:0.4}, {:0.4}, {:0.4}, {:0.4}]\n'.format(m[0],m[1],m[2],m[3],m[4],m[5],m[6],m[7],m[8]))
print('b: {:0.4}'.format(b))
```

```
M: [-0.1501, -0.05296, -0.1412, -0.02355, 0.06657, -0.04512, -0.1015, -0.03607, -0.\n ↵1819]\nb: 0.4213
```

## Convolution

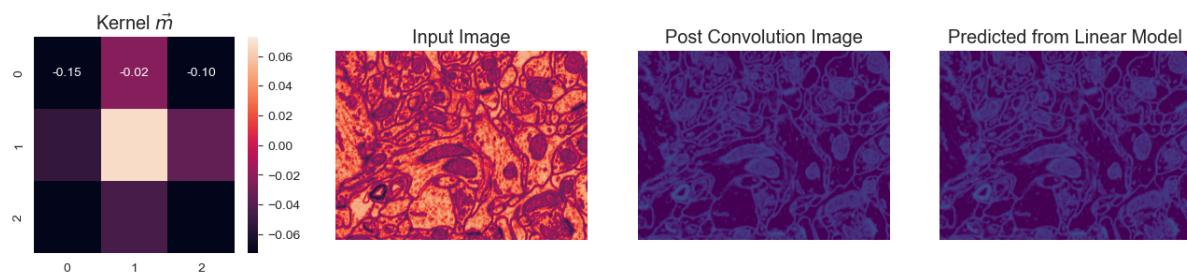
The steps we have here make up a convolution.

- What we have effectively done is use linear regression
- to learn which coefficients we should use in a convolutional kernel to get the best results

```
from scipy.ndimage import convolve
m_mat = m.reshape((3, 3)).T
fig, (ax1, ax2, ax3, ax4) = plt.subplots(1, 4, figsize=(16, 3))
sns.heatmap(m_mat,
            annot=True,
            ax=ax1, fmt='2.2f',
            vmin=-m_mat.std(),
            vmax=m_mat.std())
ax1.set_title(r'Kernel $\vec{m}$')
ax2.imshow(cell_img)
ax2.set_title('Input Image')
ax2.axis('off')

ax3.imshow(convolve(cell_img, m_mat)+b,
           vmin=0,
           vmax=1,
           cmap='viridis')
ax3.set_title('Post Convolution Image')
ax3.axis('off')

ax4.imshow(pred_func(cell_img)[:, :, 0],
           cmap='viridis', vmin=0, vmax=1)
ax4.set_title('Predicted from Linear Model')
ax4.axis('off');
```

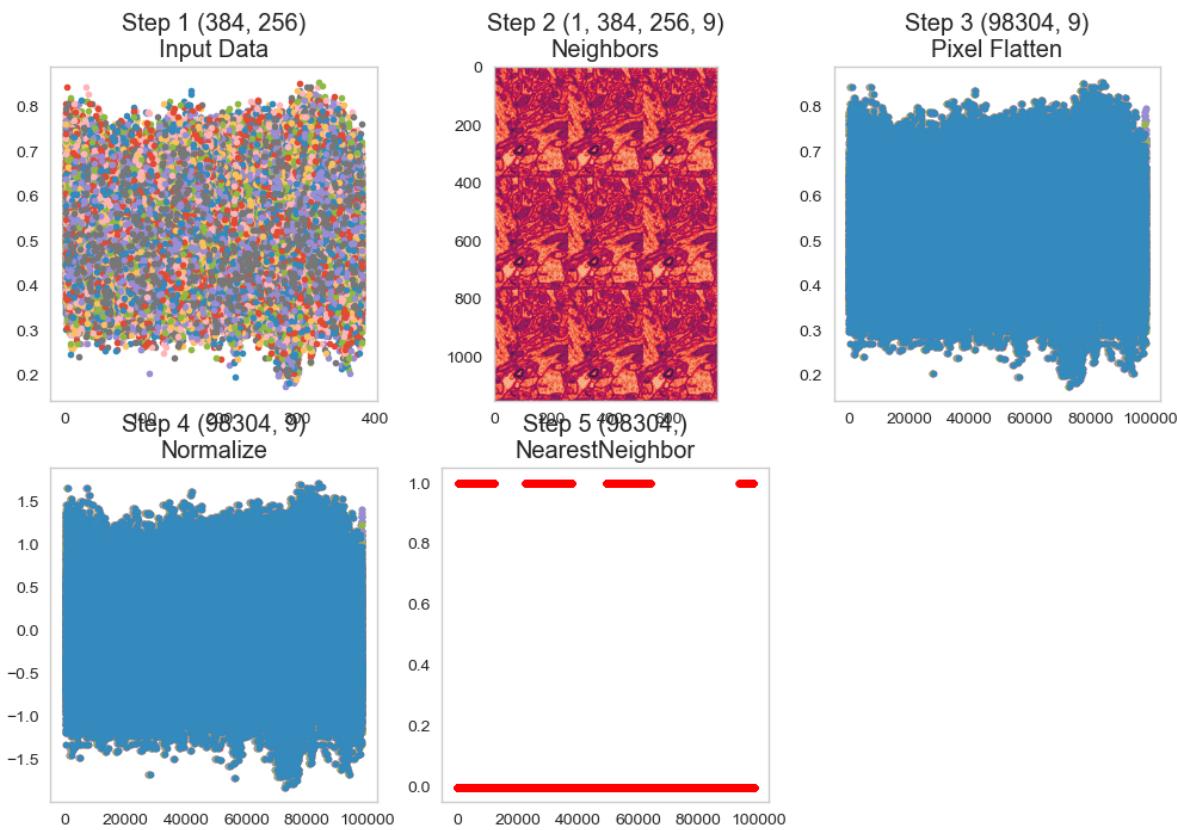


### 0.20.10 Nearest Neighbor

We can also use the neighborhood and nearest neighbor, this means for each pixel and its surrounds we find the pixel in the training set that looks most similar

```
nn_neighborhoodseg_model = Pipeline([('Neighbors', neighbor_step(1, 1)),
                                     ('Pixel Flatten', px_flatten_step),
                                     ('Normalize', RobustScaler()),
                                     ('NearestNeighbor', KNeighborsRegressor(n_neighbors=1))
                                    ])

pred_func = fit_img_pipe(nn_neighborhoodseg_model, train_img, train_mask)
show_pipe(nn_neighborhoodseg_model, train_img)
```



### Results of the nearest neighbor

```
fig, ((ax1, ax5, ax2), (ax3, ax6, ax4)) = plt.subplots(
    2, 3, figsize=(12, 8), dpi=150)
ax1.imshow(train_img, cmap='bone'); ax1.set_title('Train Image')

ax5.imshow(train_mask, cmap='viridis'); ax5.set_title('Train Mask')

ax2.imshow(pred_func(train_img)[:, :, 0], cmap='viridis', vmin=0, vmax=1); ax2.set_
           title('Prediction Mask')

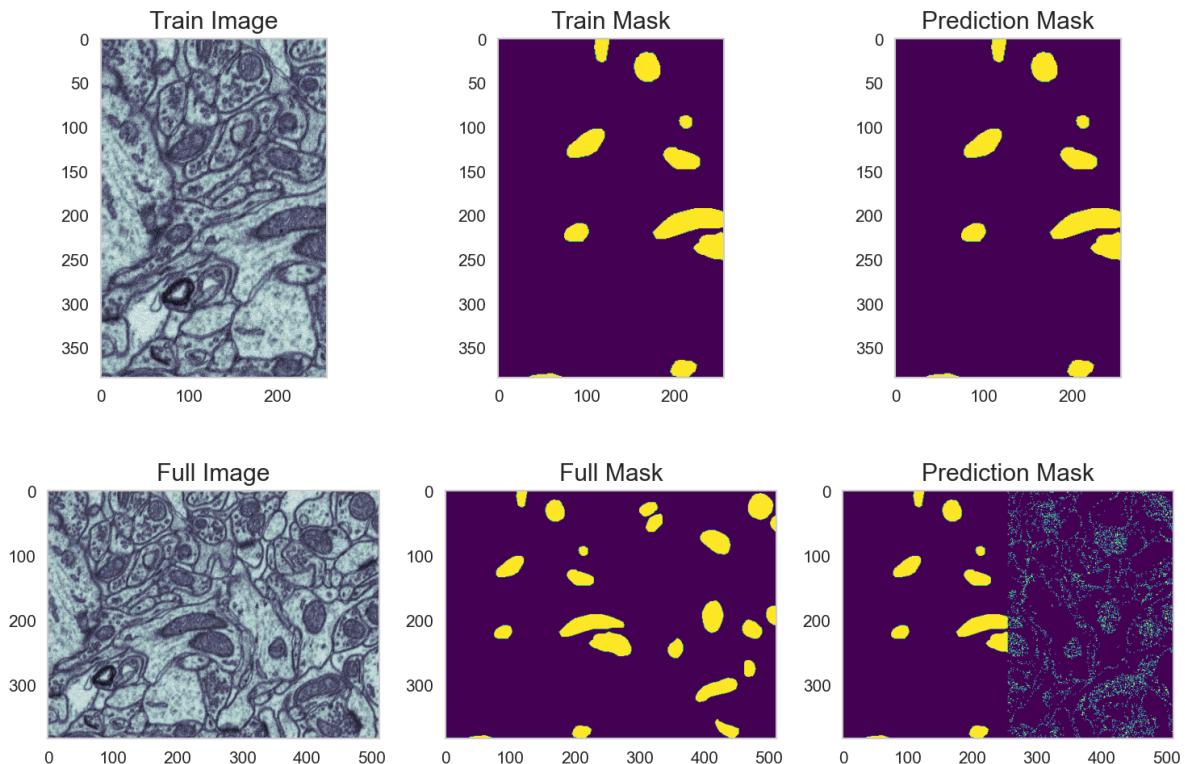
ax3.imshow(cell_img, cmap='bone'); ax3.set_title('Full Image')
```

(continues on next page)

(continued from previous page)

```
ax6.imshow(cell_seg, cmap='viridis'); ax6.set_title('Full Mask')

ax4.imshow(pred_func(cell_img)[:, :, 0], cmap='viridis', vmin=0, vmax=1)
ax4.set_title('Prediction Mask');
```



### 0.20.11 Summarizing the pipeline segmentations

- We have seen that a pipeline can be efficiently used for segmentation tasks.
  - The pipeline is easy to configure
  - It is easy to read
- Adding generated features can help improving the performance
- None of the method performed convincingly
  - Some failed on validation data
  - Other failed on all data, including training data!

## 0.21 Deep learning with a U-Net

The last approach we will briefly cover is the idea of [U-Net](#) a landmark paper from 2015 that dominates MICCAI submissions and contest winners today. A nice overview of the techniques is presented by [Vladimir Iglovikov](#) a winner of a recent Kaggle competition on masking images of cars [slides](#)

U-Net Diagram

### 0.21.1 Let's build a small U-Net

```
# import numpy as np
# import skimage.io as io
# import matplotlib.pyplot as plt
# import tensorflow as tf
# from tensorflow.keras.models import Model
# from tensorflow.keras.models import Sequential
# from tensorflow.keras.layers import Input, Conv2D, MaxPooling2D, UpSampling2D,
#     concatenate
# import tensorflow.keras.losses as losses
# from IPython.display import SVG
# # from keras.utils.vis_utils import model_to_dot
# from tensorflow.keras.utils import model_to_dot
# base_depth = 32

# t_unet = Sequential([
#     Conv2D(base_depth, kernel_size=(3, 3), padding='same', activation='relu',
#     input_shape=(None, None, 1)),
#     Conv2D(base_depth, kernel_size=(3, 3), padding='same', activation='relu'),
#     MaxPooling2D((2, 2)),
#     Conv2D(base_depth*2, kernel_size=(3, 3), padding='same', activation='relu'),
#     Conv2D(base_depth*2, kernel_size=(3, 3), padding='same', activation='relu'),
#     MaxPooling2D((2, 2)),
#     Conv2D(base_depth*4, kernel_size=(3, 3), padding='same', activation='relu'),
#     Conv2D(base_depth*4, kernel_size=(3, 3), padding='same', activation='relu'),
#     UpSampling2D((2, 2)),
#     Conv2D(base_depth*2, kernel_size=(3, 3), padding='same', activation='relu'),
#     Conv2D(base_depth*2, kernel_size=(3, 3), padding='same', activation='relu'),
#     UpSampling2D((2, 2)),
#     Conv2D(base_depth, kernel_size=(3, 3), padding='same', activation='relu'),
#     Conv2D(base_depth, kernel_size=(3, 3), padding='same', activation='relu'),
#     Conv2D(1, kernel_size=(1, 1), padding='same', activation='sigmoid')
# ], name='SmallCNN')
```

```
import numpy as np
import skimage.io as io
import matplotlib.pyplot as plt
import tensorflow as tf
from tensorflow.keras.models import Model
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Input, Conv2D, MaxPooling2D, UpSampling2D,
    concatenate
import tensorflow.keras.losses as losses

inputs = Input((None, None, 1))
```

(continues on next page)

(continued from previous page)

```

base_depth = 64

conv1 = Conv2D(base_depth, (3, 3), activation='relu', padding='same')(inputs)
conv1 = Conv2D(base_depth, (3, 3), activation='relu', padding='same')(conv1)
pool1 = MaxPooling2D(pool_size=(2, 2))(conv1)

conv2 = Conv2D(base_depth*2, (3, 3), activation='relu', padding='same')(pool1)
conv2 = Conv2D(base_depth*2, (3, 3), activation='relu', padding='same')(conv2)
pool2 = MaxPooling2D(pool_size=(2, 2))(conv2)

conv3 = Conv2D(base_depth*4, (3, 3), activation='relu', padding='same')(pool2)
conv3 = Conv2D(base_depth*4, (3, 3), activation='relu', padding='same')(conv3)

up4 = concatenate([UpSampling2D(size=(2, 2))(conv3), conv2], axis=3)
conv4 = Conv2D(base_depth*2, (3, 3), activation='relu', padding='same')(up4)
conv4 = Conv2D(base_depth*2, (3, 3), activation='relu', padding='same')(conv4)

up5 = concatenate([UpSampling2D(size=(2, 2))(conv4), conv1], axis=3)
conv5 = Conv2D(base_depth, (3, 3), activation='relu', padding='same')(up5)
conv5 = Conv2D(base_depth, (3, 3), activation='relu', padding='same')(conv5)

conv6 = Conv2D(1, (1, 1), activation='sigmoid')(conv5)

t_unet = Model(inputs=[inputs], outputs=[conv6])

```

```

-----
ImportError                                     Traceback (most recent call last)
Cell In[76], line 4
      2 import skimage.io as io
      3 import matplotlib.pyplot as plt
----> 4 import tensorflow as tf
      5 from tensorflow.keras.models import Model
      6 from tensorflow.keras.models import Sequential

File ~/anaconda3/lib/python3.11/site-packages/tensorflow/__init__.py:45
    42 from tensorflow.python import tf2 as _tf2
    43 _tf2.enable()
---> 45 from tensorflow._api.v2 import __internal__
    46 from tensorflow._api.v2 import __operators__
    47 from tensorflow._api.v2 import audio

File ~/anaconda3/lib/python3.11/site-packages/tensorflow/_api/v2/__internal__/__init__.py:8
    3 """Public API for tf._api.v2.__internal__ namespace
    4 """
    6 import sys as _sys
---> 8 from tensorflow._api.v2.__internal__ import autograph
    9 from tensorflow._api.v2.__internal__ import decorator
   10 from tensorflow._api.v2.__internal__ import dispatch

File ~/anaconda3/lib/python3.11/site-packages/tensorflow/_api/v2/__internal__/autograph/__init__.py:8
    3 """Public API for tf._api.v2.__internal__.autograph namespace
    4 """
    6 import sys as _sys

```

(continues on next page)

(continued from previous page)

```

----> 8 from tensorflow.python.autograph.core.ag_ctx import control_status_ctx #_
  ↵line: 34
  9 from tensorflow.python.autograph.impl.api import tf_convert

File ~/anaconda3/lib/python3.11/site-packages/tensorflow/python/autograph/core/ag_#
  ↵ctx.py:21
  18 import inspect
  19 import threading
---> 21 from tensorflow.python.autograph.utils import ag_logging
  22 from tensorflow.python.util.tf_export import tf_export
  25 stacks = threading.local()

File ~/anaconda3/lib/python3.11/site-packages/tensorflow/python/autograph/utils/_#
  ↵init_.py:17
  1 # Copyright 2016 The TensorFlow Authors. All Rights Reserved.
  2 #
  3 # Licensed under the Apache License, Version 2.0 (the "License");
(...).
  13 # limitations under the License.
  14 #_
  =====
  15 """Utility module that contains APIs usable in the generated code."""
---> 17 from tensorflow.python.autograph.utils.context_managers import control_
  ↵dependency_on_returns
  18 from tensorflow.python.autograph.utils.misc import alias_tensors
  19 from tensorflow.python.autograph.utils.tensor_list import dynamic_list_
  ↵append

File ~/anaconda3/lib/python3.11/site-packages/tensorflow/python/autograph/utils/
  ↵context_managers.py:19
  15 """Various context managers."""
  17 import contextlib
---> 19 from tensorflow.python.framework import ops
  20 from tensorflow.python.ops import tensor_array_ops
  23 def control_dependency_on_returns(return_value):

File ~/anaconda3/lib/python3.11/site-packages/tensorflow/python/framework/ops.py:28
  25 import types
  26 from typing import cast, TypeVar, Any, AnyStr, NoReturn, Optional, Pattern,
  ↵ Union, ContextManager
---> 28 from absl import app
  29 import numpy as np
  30 from numpy import typing as npt

ImportError: cannot import name 'app' from 'absl' (unknown location)

```

### A network summary

```
t_unet.summary()
```

```
Model: "model_14"
```

Layer (type)	Output Shape	Param #	Connected to
--------------	--------------	---------	--------------

(continues on next page)

(continued from previous page)

input_15 (InputLayer)	[ (None, None, None, 1) ]	0	[]
conv2d_211 (Conv2D) ↳ 15[0][0]'	(None, None, None, 64)	640	['input_
conv2d_212 (Conv2D) ↳ 211[0][0]'	(None, None, None, 64)	36928	['conv2d_
max_pooling2d_38 (MaxPooling2D) ↳ 212[0][0]'	(None, None, None, 64)	0	['conv2d_
conv2d_213 (Conv2D) ↳ pooling2d_38[0][0]'	(None, None, None, 128)	73856	['max_
conv2d_214 (Conv2D) ↳ 213[0][0]'	(None, None, None, 128)	147584	['conv2d_
max_pooling2d_39 (MaxPooling2D) ↳ 214[0][0]'	(None, None, None, 128)	0	['conv2d_
conv2d_215 (Conv2D) ↳ pooling2d_39[0][0]'	(None, None, None, 256)	295168	['max_
conv2d_216 (Conv2D) ↳ 215[0][0]'	(None, None, None, 256)	590080	['conv2d_
up_sampling2d_38 (UpSampling2D) ↳ 216[0][0]'	(None, None, None, 256)	0	['conv2d_
concatenate_32 (Concatenat e) ↳ 214[0][0]'	(None, None, None, 384)	0	'up_
conv2d_217 (Conv2D) ↳ 32[0][0]'	(None, None, None, 128)	442496	'conv2d_
conv2d_218 (Conv2D) ↳ 217[0][0]'	(None, None, None, 128)	147584	['conv2d_
up_sampling2d_39 (UpSampling2D) ↳ 218[0][0]'	(None, None, None, 128)	0	['conv2d_
concatenate_33 (Concatenat e) ↳ 212[0][0]'	(None, None, None, 192)	0	'up_
conv2d_219 (Conv2D) ↳ 33[0][0]'	(None, None, None, 64)	110656	'conv2d_
conv2d_220 (Conv2D) ↳ 219[0][0]'	(None, None, None, 64)	36928	['concatenate_

(continues on next page)

(continued from previous page)

```
conv2d_221 (Conv2D)           (None, None, None, 1)      65      ['conv2d_220[0][0]']
```

```
=====
Total params: 1881985 (7.18 MB)
Trainable params: 1881985 (7.18 MB)
Non-trainable params: 0 (0.00 Byte)
```

### A schematic rendering of the network

```
from IPython.display import SVG
from tensorflow.keras.utils import model_to_dot

dot_mod = model_to_dot(t_unet, show_shapes=False, show_layer_names=False, dpi=50)
dot_mod.set_rankdir('LR')
SVG(dot_mod.create_svg())
```

<IPython.core.display.SVG object>

### 0.21.2 New training data

We need to reduce the training image (to save time)

```
cell_img = ((io.imread("data/em_image.png") [::2, ::2])/255.0).astype(float)
cell_seg = (io.imread("data/em_image_seg.png", as_gray=True) [::2, ::2] > 0).
    astype(float)

train_img, valid_img = cell_img[:256, 50:250], cell_img[:, 256:]
train_mask, valid_mask = cell_seg[:256, 50:250], cell_seg[:, 256:]

# add channels and sample dimensions
def prep_img(x, n=1): return (
    prep_mask(x, n=n)-train_img.mean()) / train_img.std()

def prep_mask(x, n=1): return np.stack([np.expand_dims(x, -1)]*n, 0)

print('Training', train_img.shape, train_mask.shape)
print('Validation Data', valid_img.shape, valid_mask.shape)
```

Training (256, 200) (256, 200)  
Validation Data (384, 256) (384, 256)

```
from matplotlib.patches import Rectangle
from matplotlib.spines import Spine

r_train = Rectangle((50, 0), 200, 256, fc='none', ec='limegreen', lw=3)
```

(continues on next page)

(continued from previous page)

```

r_validate = Rectangle((256,0),cell_img.shape[1]-258,cell_img.shape[0]-1,fc='none',ec=
    ↪'magenta',lw=3)
fig, ax = plt.subplots(2, 3, figsize=(15, 6) )
ax=ax.ravel()
cmap='gray'

def coloraxes(ax,color) :
    for child in ax.get_children():
        if isinstance(child, Spine):
            child.set_color(color)
            child.set_linewidth(3)

ax[0].imshow(train_img, cmap=cmap)
ax[0].set_title('Train Image')
coloraxes(ax[0],'limegreen')

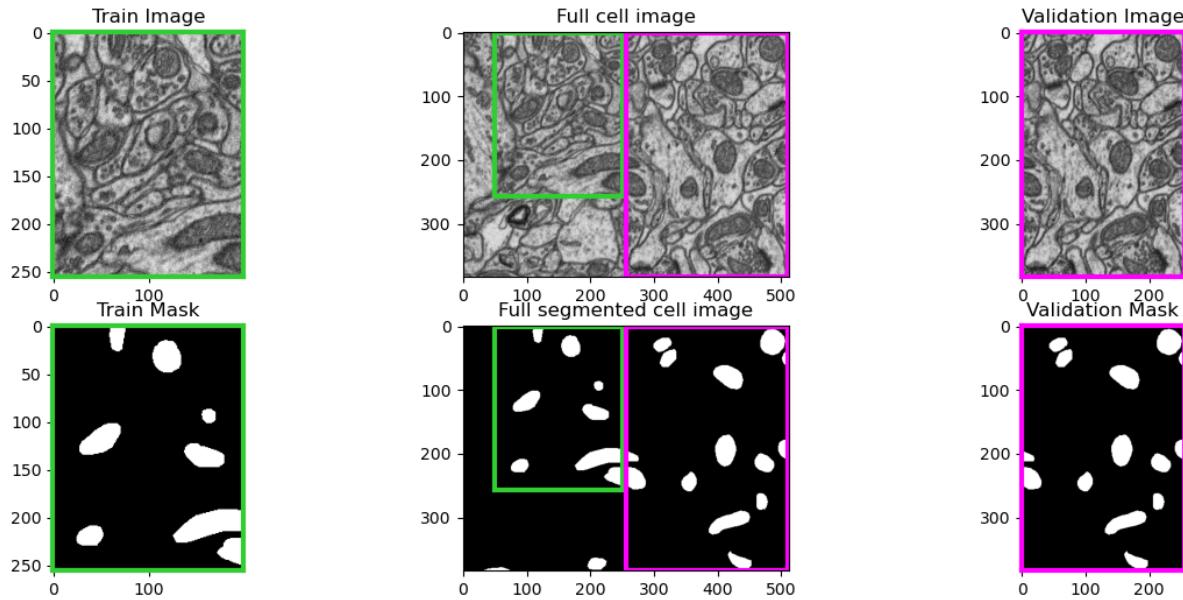
ax[3].imshow(train_mask, cmap=cmap)
ax[3].set_title('Train Mask')
coloraxes(ax[3],'limegreen')

ax[1].imshow(cell_img,cmap=cmap)
ax[1].set_title('Full cell image')
ax[1].add_patch(r_train)
ax[1].add_patch(r_validate)

r_train      = Rectangle((50,0),200,256,fc='none',ec='limegreen',lw=3)
r_validate = Rectangle((256,0),cell_img.shape[1]-258,cell_img.shape[0]-1,fc='none',ec=
    ↪'magenta',lw=3)
ax[4].imshow(cell_seg,cmap=cmap)
ax[4].set_title('Full segmented cell image')
ax[4].add_patch(r_train)
ax[4].add_patch(r_validate)

ax[2].imshow(valid_img, cmap=cmap)
ax[2].set_title('Validation Image')
coloraxes(ax[2],'magenta')
ax[5].imshow(valid_mask, cmap=cmap)
ax[5].set_title('Validation Mask');
coloraxes(ax[5],'magenta')

```



### 0.21.3 Results from Untrained Model

- We can make predictions with an untrained model (default parameters)
- but we clearly do not expect them to be very good

```
unet_pred = t_unet.predict(prep_img(cell_img)) [0, :, :, 0];
```

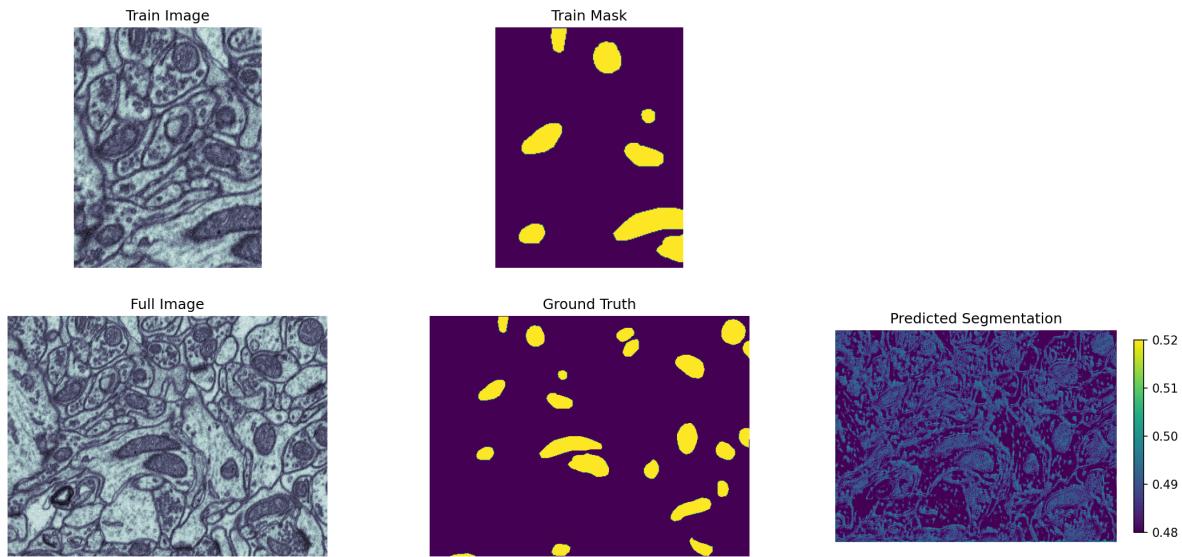
```
1/1 [=====] - 0s 217ms/step
```

```
fig, m_axs = plt.subplots(2, 3,
                        figsize=(18, 8), dpi=150)
for c_ax in m_axs.flatten():
    c_ax.axis('off')
((ax1, ax2, _), (ax3, ax4, ax5)) = m_axs
ax1.imshow(train_img, cmap='bone')
ax1.set_title('Train Image')
ax2.imshow(train_mask, cmap='viridis')
ax2.set_title('Train Mask')

ax3.imshow(cell_img, cmap='bone')
ax3.set_title('Full Image')

ax4.imshow(cell_seg,
           cmap='viridis')
ax4.set_title('Ground Truth');

a5=ax5.imshow(unet_pred,
              cmap='viridis', vmin=0.48,vmax=0.52)
ax5.set_title('Predicted Segmentation')
fig.colorbar(a5,ax=ax5,shrink=0.8);
```



Note here that the predictions all are around 0.5, *i.e.*, close to a random guess.

#### 0.21.4 A general note on the following demo

This is a very bad way to train a model;

- the loss function is poorly chosen,
- the optimizer can be improved the learning rate can be changed,
- the training and validation data **should not** come from the same sample (and **definitely** not the same measurement).

The goal is to be aware of these techniques and have a feeling for how they can work for complex problems

#### Training conditions

- Loss function - MAE
- Optimizer - Stochastic Gradient Decent and specifically the Adam optimizer.
- 20 Epochs (training iterations)
- Metrics
  1. Binary accuracy (percentage of pixels correct classified)  $BA = \frac{1}{N} \sum_i (f_i == g_i)$
  2. Mean absolute error

Another popular metric is the Dice score  $DSC = \frac{2|X \cap Y|}{|X| + |Y|} = \frac{2TP}{2TP + FP + FN}$

### Let's train the model

```
opt = tf.keras.optimizers.legacy.Adam(learning_rate=1e-4)

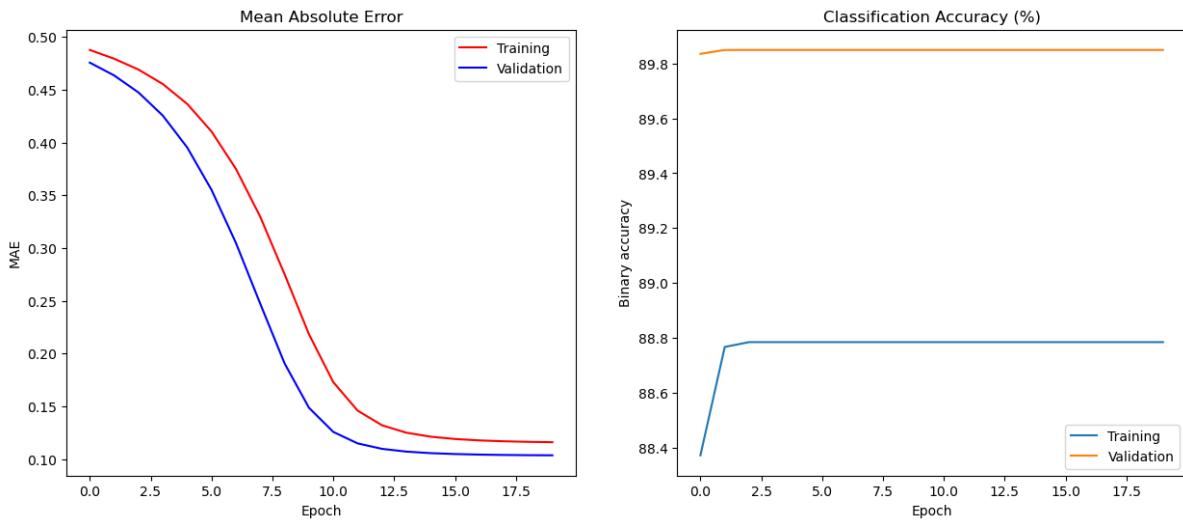
t_unet.compile(
    optimizer=opt,
    # we use a simple loss metric of mean-squared error to optimize
    loss="MSE",
    # we keep track of the number of pixels correctly classified and the mean_
    ↵absolute error as well
    metrics=['binary_accuracy', 'mae'])

loss_history = t_unet.fit(prep_img(train_img),
                          prep_mask(train_mask),
                          validation_data=(prep_img(valid_img),
                                           prep_mask(valid_mask)),
                          epochs=20, verbose=0)
```

```
fig, (ax1, ax2) = plt.subplots(1, 2,
                               figsize=(15, 6))
ax1.plot(loss_history.epoch,
          loss_history.history['mae'], 'r-', label='Training')
ax1.plot(loss_history.epoch,
          loss_history.history['val_mae'], 'b-', label='Validation')
ax1.set_title('Mean Absolute Error')
ax1.set_xlabel('Epoch')
ax1.set_ylabel('MAE')
ax1.legend()

ax2.plot(loss_history.epoch,
          100*np.array(loss_history.history['binary_accuracy']), '-',
          label='Training')
ax2.plot(loss_history.epoch,
          100*np.array(loss_history.history['val_binary_accuracy']), '-',
          label='Validation')

ax2.set_xlabel('Epoch')
ax2.set_ylabel('Binary accuracy')
ax2.set_title('Classification Accuracy (%)')
ax2.legend();
```



## Prediction results

```
unet_train = t_unet.predict(prep_img(train_img)) [0, :, :, 0]
unet_pred = t_unet.predict(prep_img(cell_img)) [0, :, :, 0]
```

```
1/1 [=====] - 0s 152ms/step
1/1 [=====] - 1s 788ms/step
```

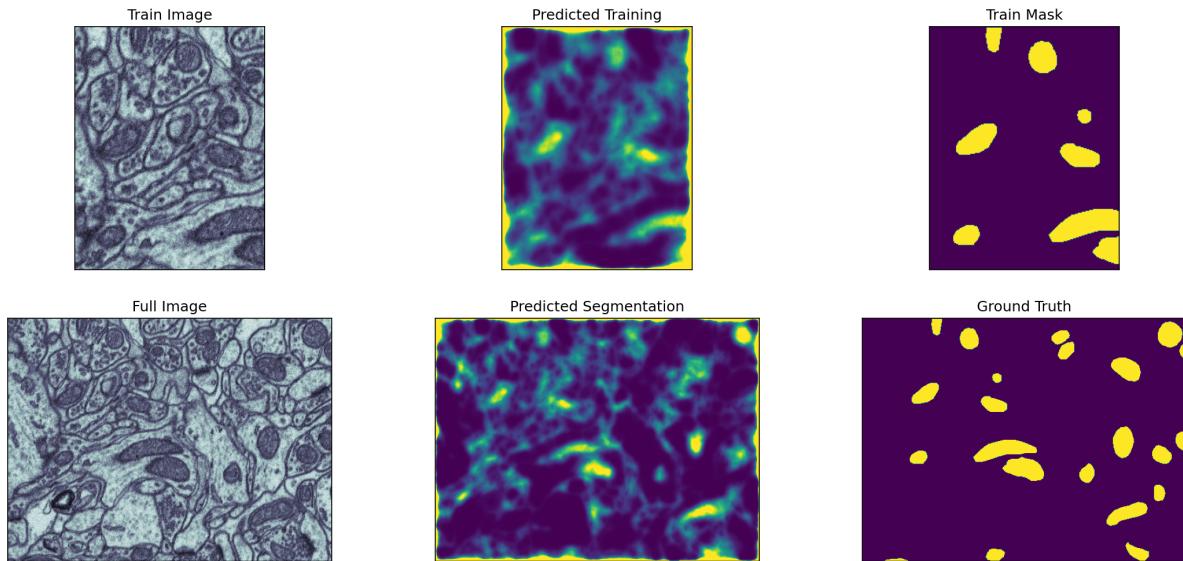
```
fig, m_axs = plt.subplots(2, 3, figsize=(18, 8), dpi=150)
for c_ax in m_axs.flatten():
    c_ax.set(xticks=[], yticks=[])
((ax1, ax15, ax2), (ax3, ax4, ax5)) = m_axs
ax1.imshow(train_img, cmap='bone')
ax1.set_title('Train Image')
vmin=0
vmax=0.01
ax15.imshow(unet_train, cmap='viridis', vmin=vmin, vmax=vmax)
ax15.set_title('Predicted Training')
ax2.imshow(train_mask, cmap='viridis')
ax2.set_title('Train Mask')

ax3.imshow(cell_img, cmap='bone')
ax3.set_title('Full Image')

ax4.imshow(unet_pred,
           cmap='viridis', vmin=vmin, vmax=vmax)
ax4.set_title('Predicted Segmentation')

ax5.imshow(cell_seg,
           cmap='viridis')
ax5.set_title('Ground Truth');
```

```
Text(0.5, 1.0, 'Ground Truth')
```



### 0.21.5 Overfitting

Having a model with 470,000 free parameters means that it is quite easy to overfit the model by training for too long.

Overfitting is when:

- The model has gotten very good at the training data
- but hasn't generalized to other kinds of problems

**Consequence:** The model starts to perform worse on regions that aren't exactly the same as the training.

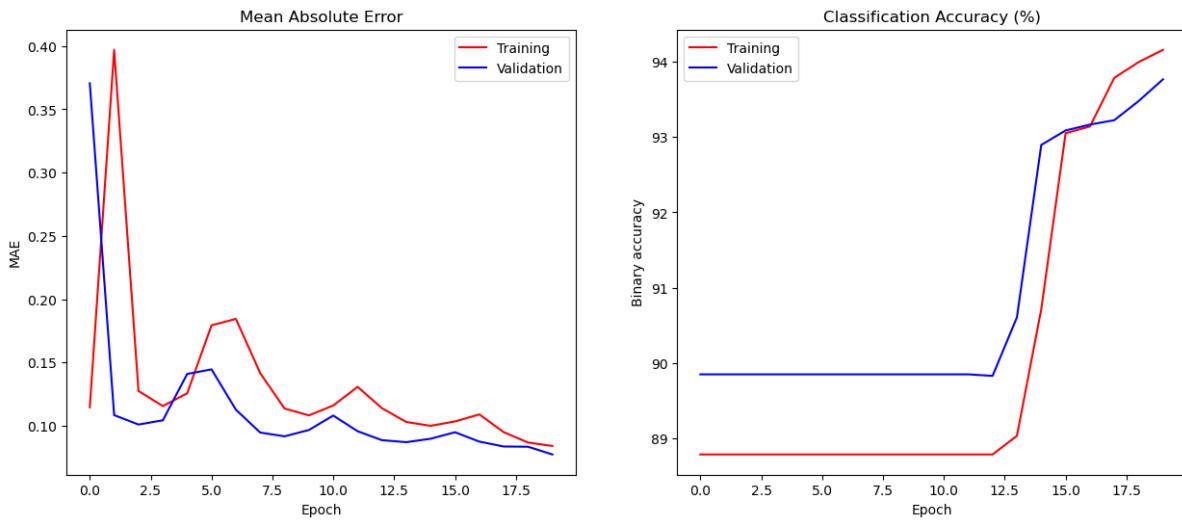
```
t_unet.compile(
    # we use a simple loss metric of mean-squared error to optimize
    loss='mse',
    optimizer='Adam',
    # we keep track of the number of pixels correctly classified and the mean_
    ↪absolute error as well
    metrics=['binary_accuracy', 'mae']
)

loss_history = t_unet.fit(prep_img(train_img),
                          prep_mask(train_mask),
                          validation_data=(prep_img(valid_img),
                                           prep_mask(valid_mask)),
                          epochs=20, verbose=0)
```

## Loss history

```
fig, (ax1, ax2) = plt.subplots(1, 2,
                             figsize=(15, 6))
ax1.plot(loss_history.epoch,
          loss_history.history['mae'], 'r-', label='Training')
ax1.plot(loss_history.epoch,
          loss_history.history['val_mae'], 'b-', label='Validation')
ax1.set_title('Mean Absolute Error')
ax1.set_xlabel('Epoch')
ax1.set_ylabel('MAE')
ax1.legend()

ax2.plot(loss_history.epoch,
          100*np.array(loss_history.history['binary_accuracy']), 'r-', label='Training')
ax2.plot(loss_history.epoch,
          100*np.array(loss_history.history['val_binary_accuracy']), 'b-', label='Validation')
ax2.set_xlabel('Epoch')
ax2.set_ylabel('Binary accuracy')
ax2.set_title('Classification Accuracy (%)')
ax2.legend();
```



## Prediction results

```
unet_train = t_unet.predict(prep_img(train_img))[0, :, :, 0]
unet_pred = t_unet.predict(prep_img(cell_img))[0, :, :, 0]
```

```
1/1 [=====] - 0s 167ms/step
1/1 [=====] - 0s 176ms/step
```

```
fig, m_axs = plt.subplots(2, 3,
                        figsize=(18, 8), dpi=150)
for c_ax in m_axs.flatten():
```

(continues on next page)

(continued from previous page)

```

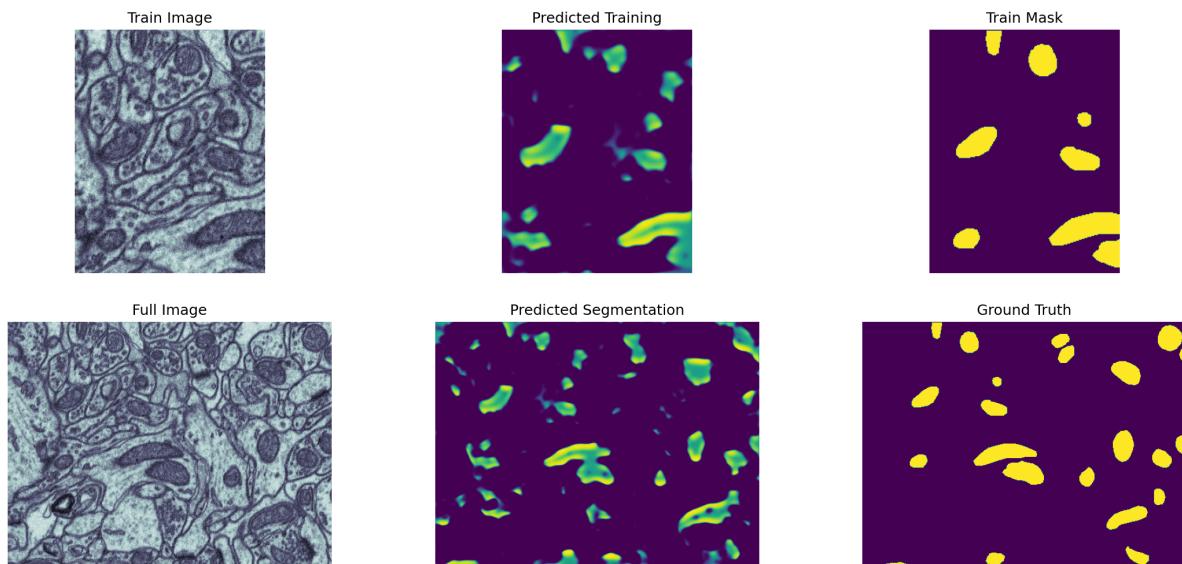
c_ax.axis('off')
((ax1, ax15, ax2), (ax3, ax4, ax5)) = m_axs
ax1.imshow(train_img, cmap='bone')
ax1.set_title('Train Image')
ax15.imshow(unet_train, cmap='viridis', vmin=0, vmax=1)
ax15.set_title('Predicted Training')
ax2.imshow(train_mask, cmap='viridis')
ax2.set_title('Train Mask')

ax3.imshow(cell_img, cmap='bone')
ax3.set_title('Full Image')

ax4.imshow(unet_pred,
           cmap='viridis', vmin=0, vmax=1)
ax4.set_title('Predicted Segmentation')

ax5.imshow(cell_seg,
           cmap='viridis')
ax5.set_title('Ground Truth');

```



## 0.22 Summary

- Concepts of supervised segmentation
- Supervised Classification
- Classification vs. Segmentation
  - Nearest neighbour
  - Trees
  - Random Forests
- Regression vs Classification
- Training, validation

- Deep learning