
Quantitative Big Imaging - Advanced segmentation

Anders Kaestner

Mar 31, 2022

CONTENTS

1	Advanced segmentation	3
1.1	Today's Outline	3
1.2	Literature / Useful References	3
1.3	Load some needed modules	3
1.4	Previously on QBI	4
1.5	Different types of segmentation	4
2	Where segmentation fails	7
2.1	Typical image features that makes life harder	7
2.2	Inconsistent illumination on real data	8
2.3	Apply a treshold	10
2.4	Where segmentation fails: Canaliculi	11
2.5	Where segmentation fails: Brain Cortex	11
3	Apply thresholds	13
3.1	Automated Threshold Selection	14
3.2	Automated Methods	15
3.3	Histogram-based Methods	15
3.4	Histogram-Methods	16
3.5	Try All Thresholds	16
3.6	Pitfalls of different thresholding methods	19
3.7	Realistic Approaches for Dealing with these Shortcomings	19
3.8	Hysteresis Thresholding	20
3.9	Goldilocks Situation	21
3.10	Multiple thresholds	22
4	More Complicated Images	25
4.1	Feature Vectors	25
4.2	Let's create a feature table	26
4.3	Inspect the features - Intensity vs. IsText	26
4.4	What does the ROC curve look like?	27
5	Adding Information	29
5.1	Testing with enhanced edges	30
5.2	Histogram of enhanced image	30
5.3	Checking the ROC performance	31
5.4	Why does the second filter perform better?	32
6	Clustering / Classification (Unsupervised)	35
6.1	Example data	35
6.2	K-Means Algorithm	36

6.3	What vector space do we have?	38
6.4	Add spatial information to k-means	38
6.5	When can clustering be used on images?	43
7	Quad trees	45
7.1	Principle	45
7.2	Quad tree example	45
8	Superpixels	47
8.1	Why use superpixels	47
8.2	Super pixels	48
8.3	Merging super pixels	49
8.4	Segmentation using super pixels	50
9	Probabilistic Models of Segmentation	51
10	Summary	53
10.1	Histogram based thresholding	53
10.2	Clustering - K-means	53
10.3	Similar region segmentations	53
11	Supervised Segmentation Approaches	55
11.1	Overview	55
11.2	Reading Material	55
12	Basic Methods Overview	57
12.1	Training	57
12.2	Predicting	57
13	Classification	59
13.1	Lets create some data...	59
13.2	Nearest Neighbor (or K Nearest Neighbors)	60
13.3	Text example	61
13.4	Nearest neighbor predictions	62
13.5	Let's come back to the blob data	62
13.6	Training the model using one neighbor	63
13.7	Stabilizing Results - Increase number of neighbors	64
14	Linear Regression	65
14.1	We need data to fit...	65
14.2	Regression on blob data	66
14.3	Train the regression model	67
15	Decision trees	69
15.1	Create a decision tree classifier	69
15.2	Decision trees on the blob data	70
15.3	Random Forests	73
16	Pipelines	77
16.1	Let's the return to the blobs	77
16.2	A basic pipeline	78
17	Classification using a pipeline	81
17.1	A first example with pipeline processing	81
17.2	Test classifier performance	83
17.3	Wow! We've built an amazing algorithm!	85

17.4	Training, Validation, and Testing	86
18	Regression using a pipeline	87
18.1	Load the digits data again	87
18.2	Run a KNeighbors regression	88
19	Assessment of multi-category data	89
19.1	Increasing neighbor count	90
20	Segmentation (Pixel Classification)	93
20.1	Where segmentation fails: Mitochondria Segmentation in EM	93
20.2	Let's try some methods to segment the mitochondria image	93
20.3	Preparing the mitochondria image and mask	94
20.4	Try a Regression tree	95
20.5	Regression tree with <i>position information</i>	96
20.6	Combine K-Means and Random Forest Regression	98
20.7	Trying polynomial features	100
20.8	Adding smarter features	102
20.9	Using the Neighborhood	104
20.10	Nearest Neighbor	108
20.11	Summarizing the pipeline segmentations	110
21	Deep learning with a U-Net	111
21.1	Let's build a small U-Net	111
21.2	New training data	115
21.3	Results from Untrained Model	116
21.4	A general note on the following demo	117
21.5	Overfitting	120
22	Summary	123

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.

ADVANCED SEGMENTATION

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

1.2 Literature / Useful References

- Superpixels

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
```

(continues on next page)

(continued from previous page)

```
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
```

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

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.

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

Woodland Encounter Bev Doolittle

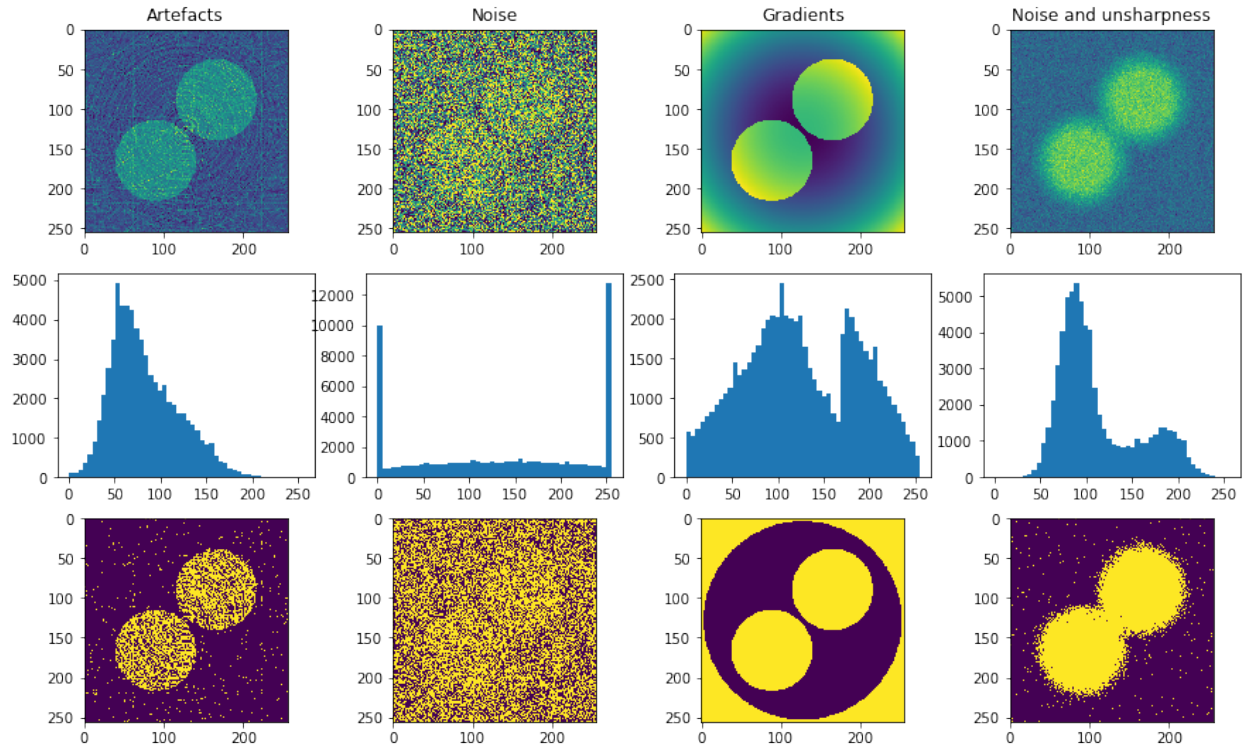
2.1 Typical image features that makes life harder

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])
    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')
```



2.2 Inconsistent illumination on real data

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

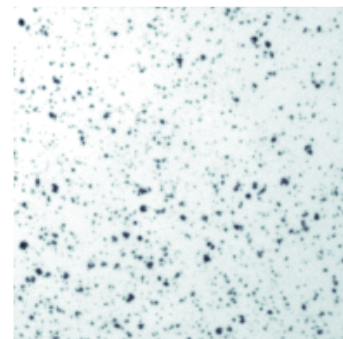
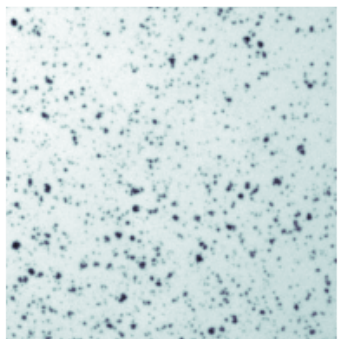
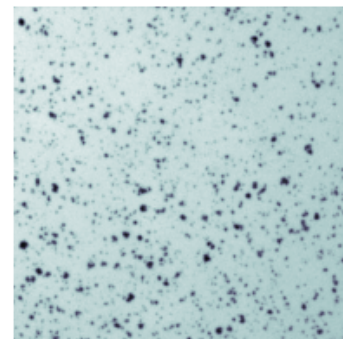
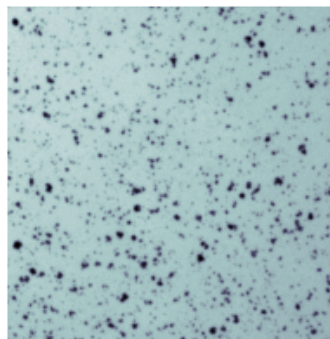
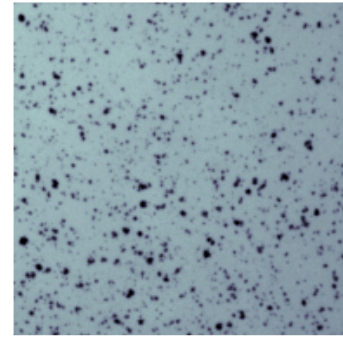
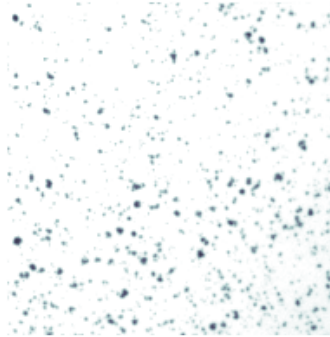
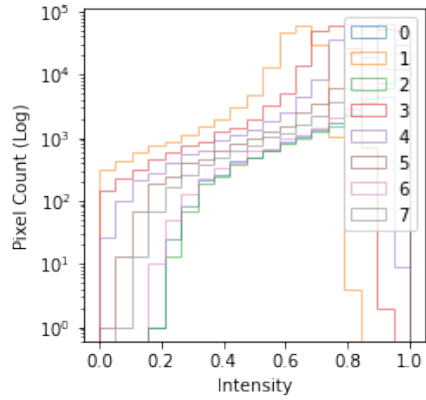
In this particular example of a cell colony, there is a random bias added to the images.

```
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 = 0.5, histtype = 'step')
    c_ax.imshow(c_img, cmap = 'bone', vmin = 0, vmax = 1)
    c_ax.axis('off')
```

(continues on next page)

(continued from previous page)

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

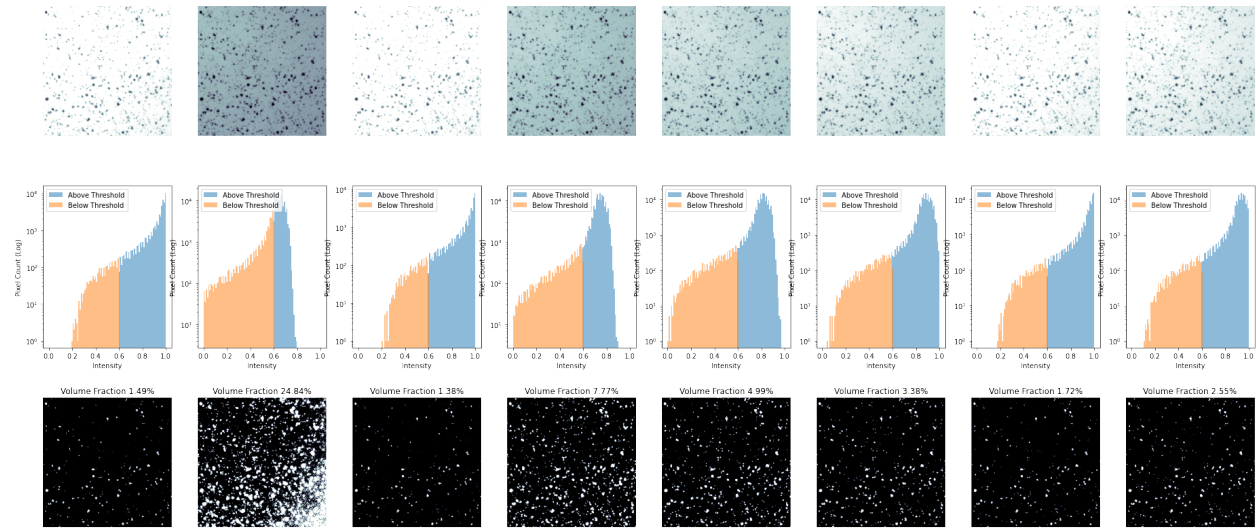


2.3 Apply a threshold

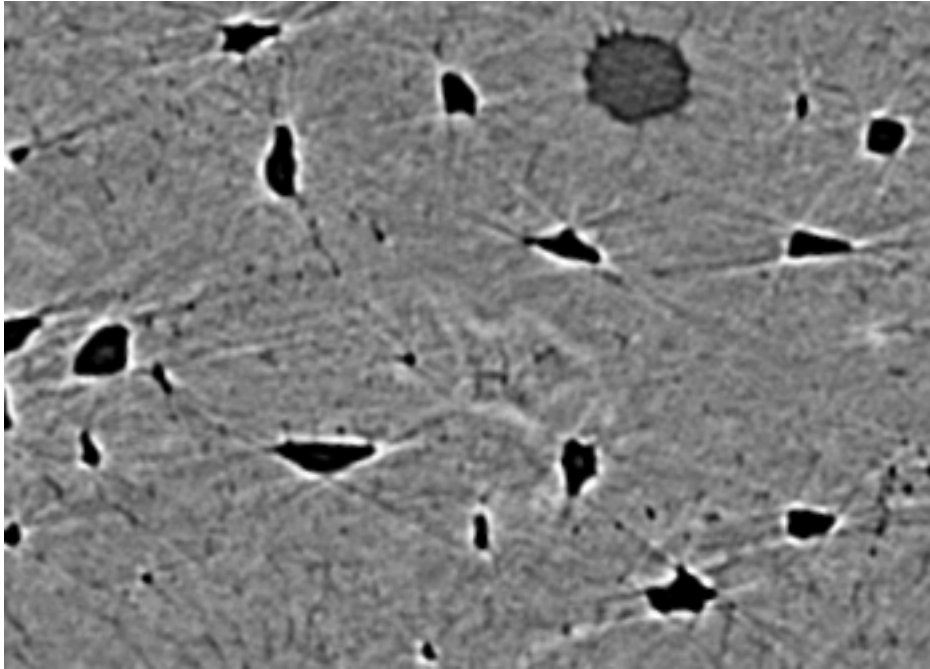
When we apply a constant threshold of 0.6 to the images with different illuminations we just looked at you will see that there are great differences in how the cells are represented.

```
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',
            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.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');
```



2.4 Where segmentation fails: Canaliculi



2.4.1 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.

2.5 Where segmentation fails: Brain 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

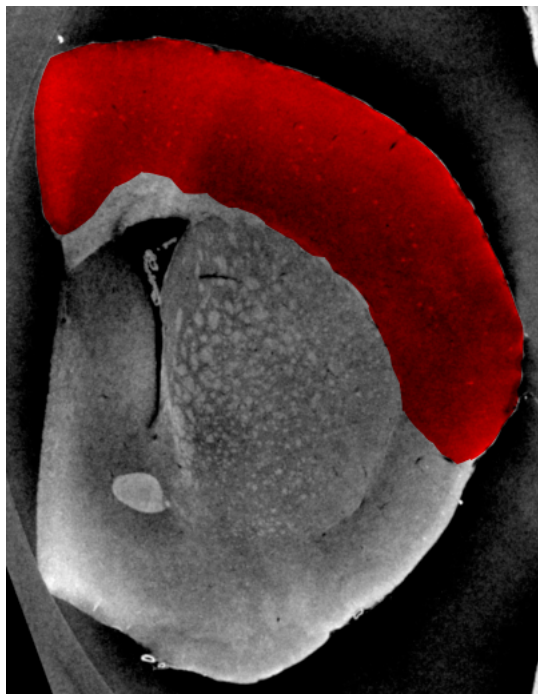


Fig. 2.2: Brain image with masked cortex.

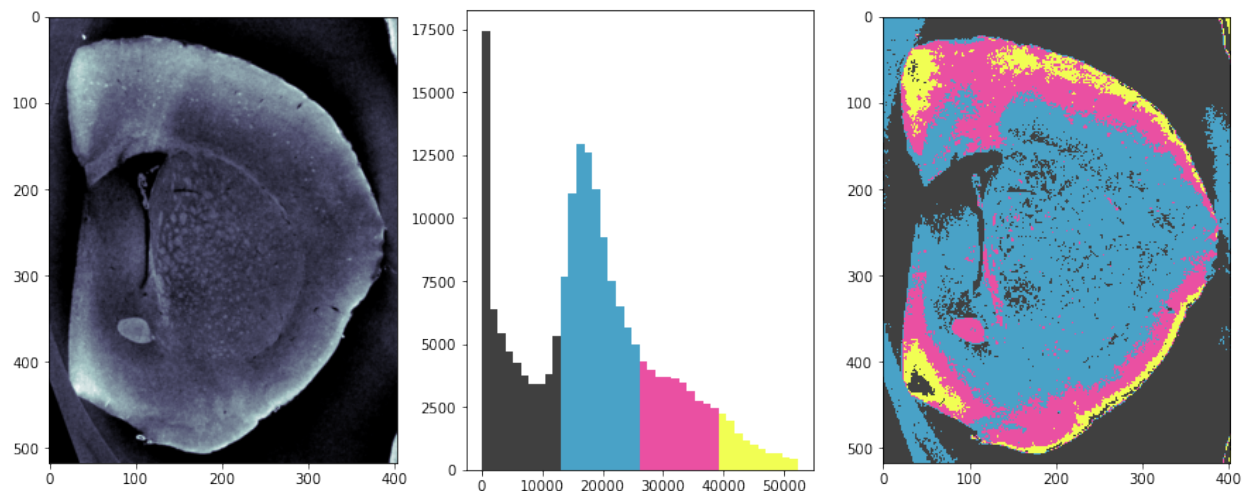
APPLY THRESHOLDS

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

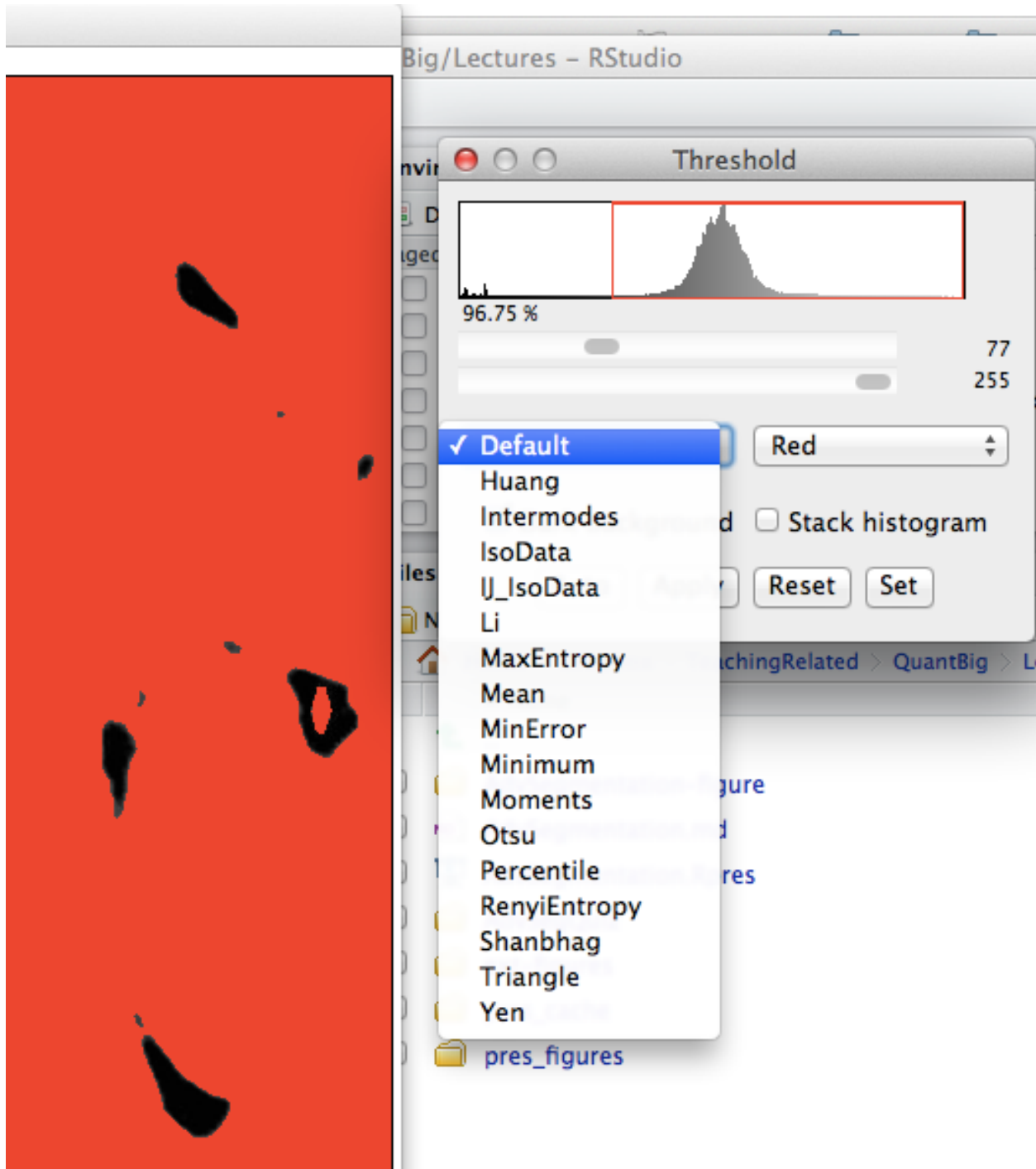
```
cortex_img = imread("figures/cortex.png")
np.random.seed(2018)
fig, (ax1, ax2, ax3) = plt.subplots(1, 3,
                                   figsize = (15, 6))
ax1.imshow(cortex_img, cmap = 'bone')

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');
```



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.

3.2 Automated Methods

3.2.1 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, Intermodos, etc

3.2.2 Image based Methods

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

3.2.3 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

3.3 Histogram-based Methods

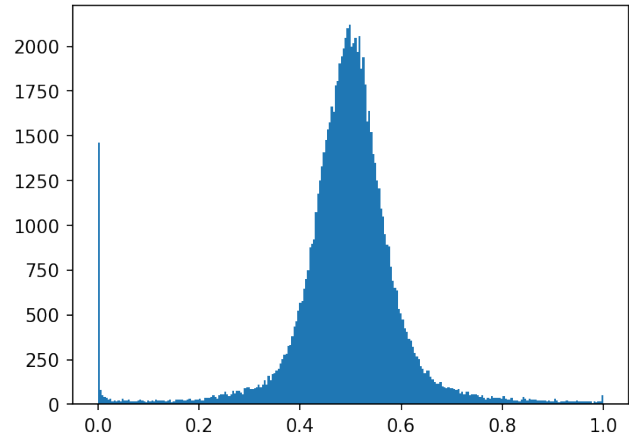
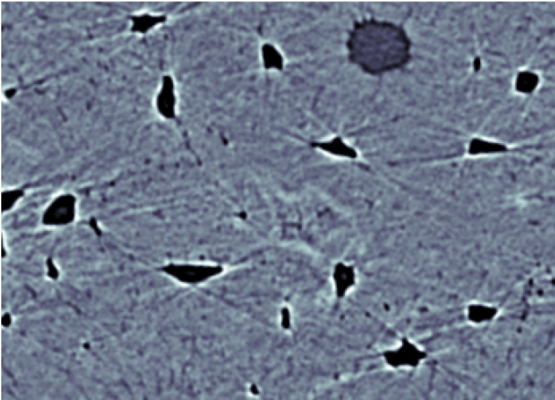
Taking a typical image of a bone slice, we can examine the variations in calcification density in the image

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/bonefiltslice.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);
```



3.4 Histogram-Methods

3.4.1 Intermodes

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

3.4.2 Otsu

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

3.4.3 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 $\text{thresh} = (\text{avg}(bg) + \text{avg}(obj))/2$

3.5 Try All Thresholds

- `opencv2`
- `scikit-image`

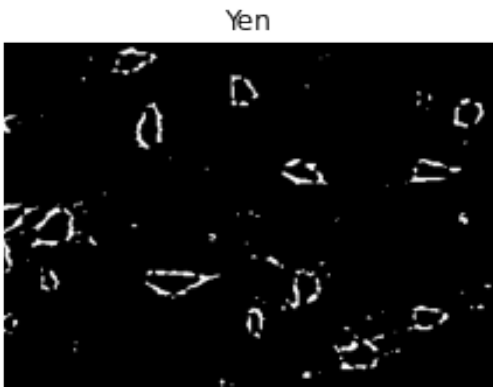
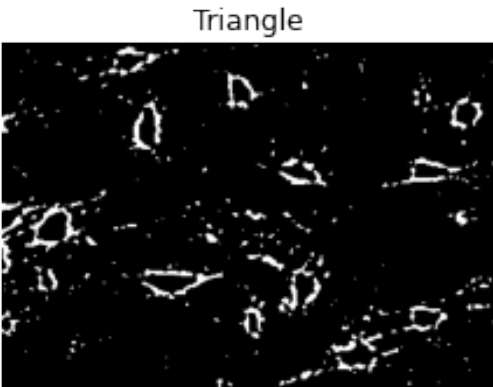
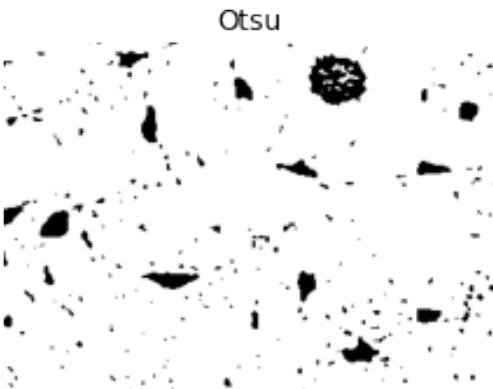
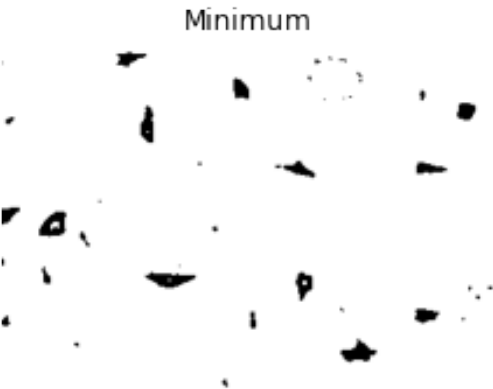
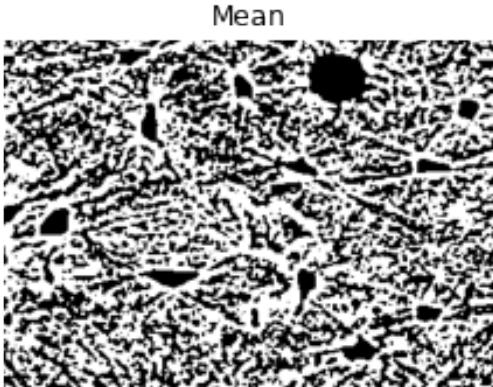
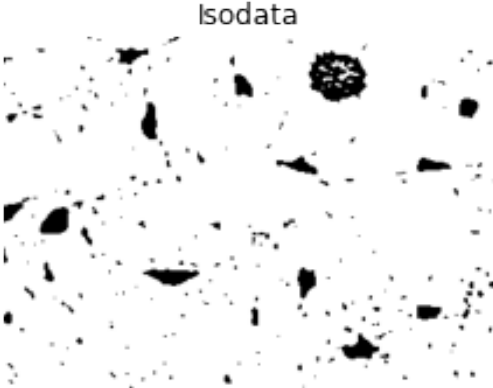
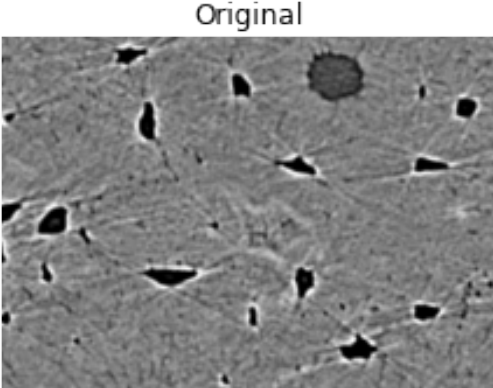
There are many methods and they can be complicated to implement yourself.

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/bonefiltslice.png")/255.0

from skimage.filters import try_all_threshold

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



3.6 Pitfalls of different thresholding methods

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

3.6.1 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

3.6.2 Image-based

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

3.6.3 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.

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

3.7.1 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

3.7.2 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.

3.7.3 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).

3.8 Hysteresis Thresholding

For some images a single threshold does not work

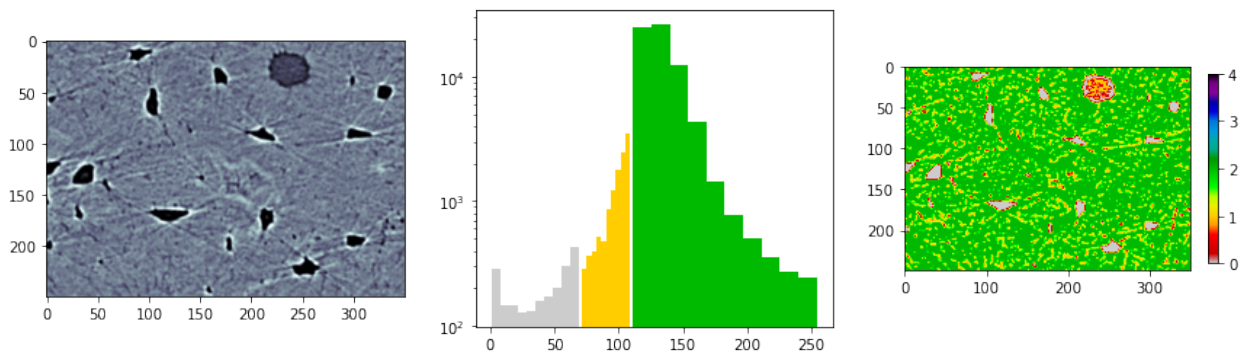
- large structures are very clearly defined
- smaller structures are difficult to differentiate (see [partial volume effect](#))

ImageJ Source

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

```
bone_img = imread("figures/bonefiltslice.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))
plt.colorbar(th_ax, shrink=0.6);
```



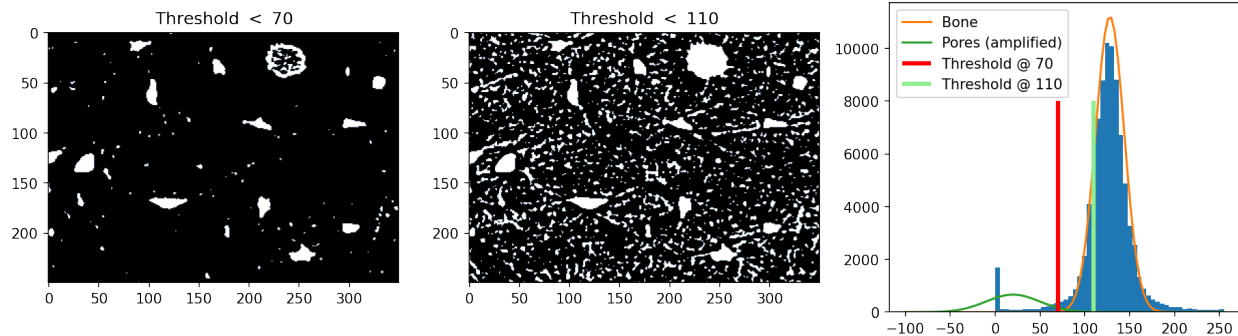
3.9 Goldilocks Situation

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 < thresh_vals[1], cmap = 'bone')
ax1.set_title('Threshold < $ %d' % (thresh_vals[1]))
ax2.imshow(bone_img < thresh_vals[2], cmap = 'bone')
ax2.set_title('Threshold < $ %d' % (thresh_vals[2]));
import scipy.stats as stats
ax3.hist(bone_img.ravel(), bins=50);
x = np.linspace(-100, 255, 100)
ax3.plot(x, 4.5e5 * stats.norm.pdf(x, 128, 16), label="Bone")
ax3.plot(x, 5e4 * stats.norm.pdf(x, 20, 30), label="Pores (amplified)")
ax3.vlines([70], ymin=0, ymax=8000, color='r', label="Threshold @ 70", lw=3)
ax3.vlines([110], ymin=0, ymax=8000, color='lightgreen', label="Threshold @ 110", lw=3)
ax3.legend();
```



3.9.1 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
```

3.9.2 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 < thresh_vals[1]
step_list['Remove Small Objects'] = opening(step_list['Strict Threshold'], disk(1))
step_list['Looser Threshold'] = bone_img < thresh_vals[2]
step_list['Both Thresholds'] = 1.0 * step_list['Looser Threshold'] + 1.0 * step_list[
    ↪ 'Remove Small Objects']
```

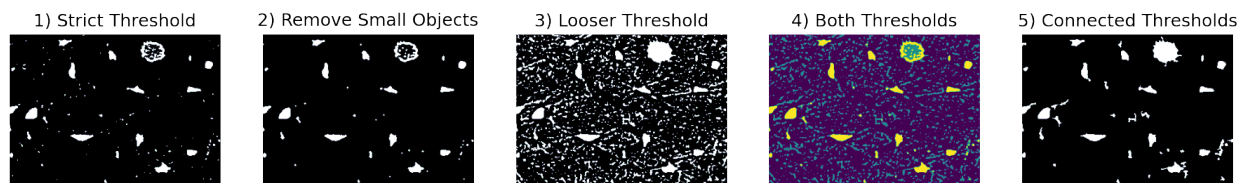
Step 2 Combine the different images $I_{ConnectedThresholds} = \delta(I_{NoSmallObjects}) \cap I_{LooseThresh}$

```
# 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' % i + 's' % (i, c_title)); c_ax.axis('off');
```



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!

3.10.1 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.

3.10.2 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 of Gaussian filter to increase the robustness for low SNR.

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. $\$I(x, y) = \hat{f}(x, y)\$$

4.1 Feature Vectors

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

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.

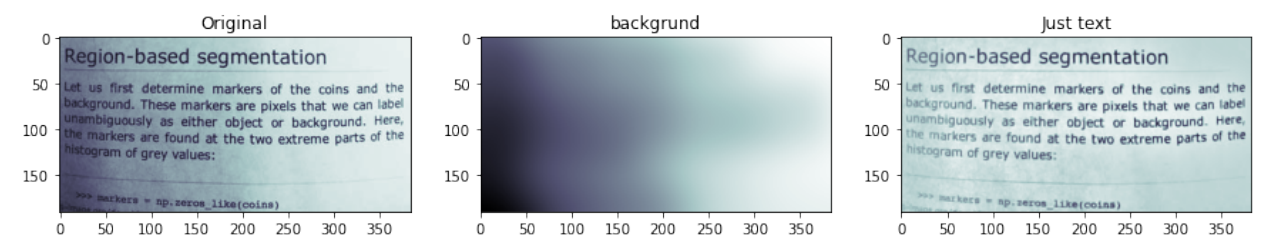
```

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

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

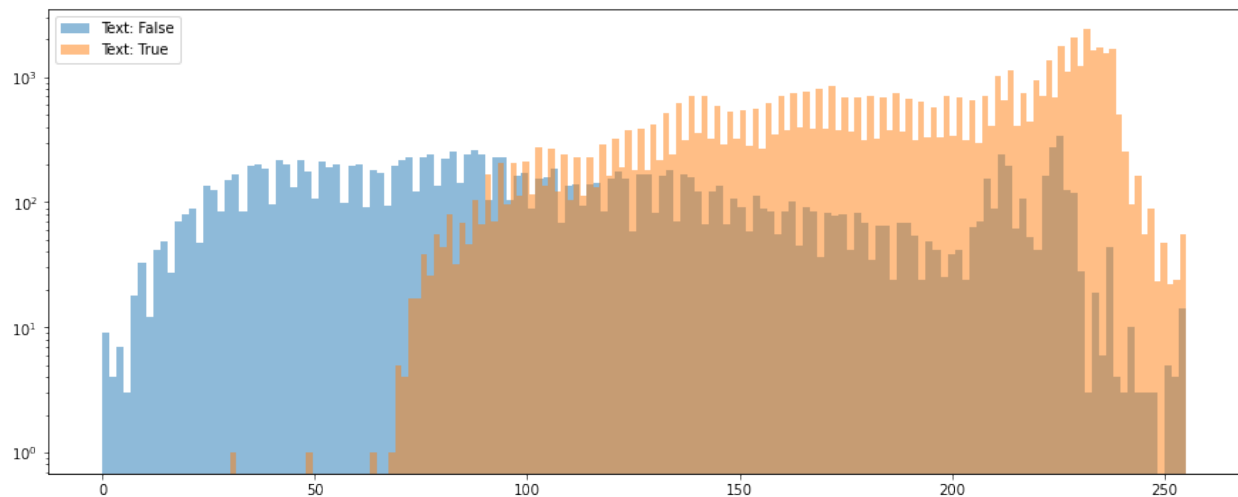
4.2 Let's create a feature table

```
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

4.3 Inspect the features - Intensity vs. IsText

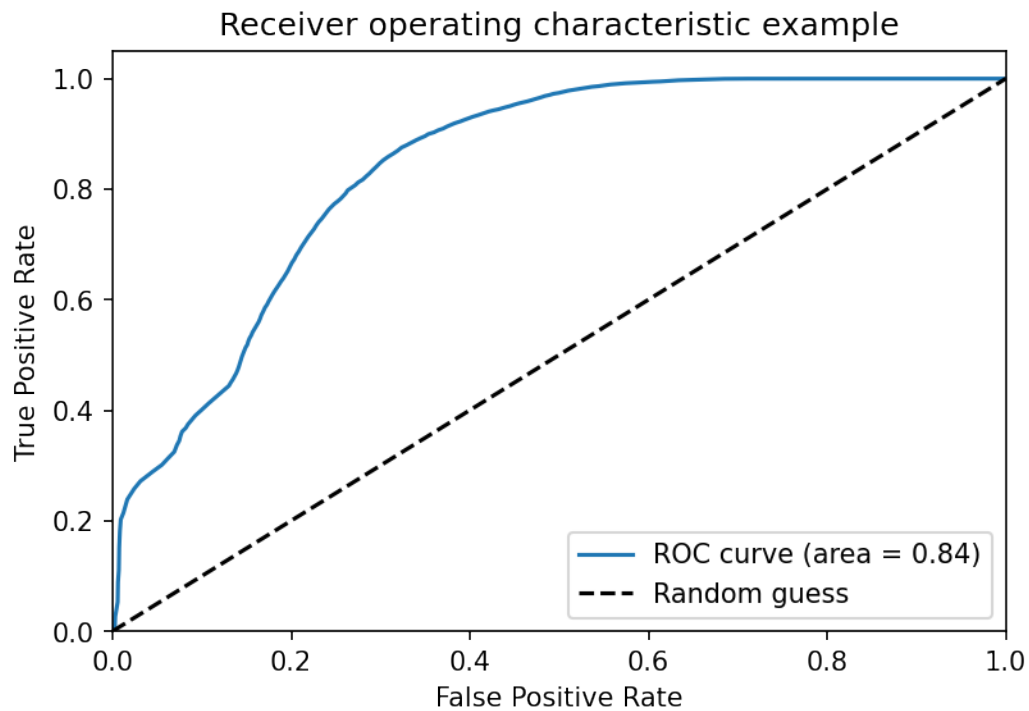
```
fig, ax1 = plt.subplots(1,1,figsize=(15,6))
for c_cat, c_df in page_table.groupby(['is_text']):
    ax1.hist(c_df['intensity'],
            label = 'Text: {}'.format(c_cat),
            alpha = 0.5, bins=150)
ax1.set_yscale("log", nonpositive='clip')
ax1.legend();
```



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,dpi=150)
ax.plot(fpr, tpr, label='ROC curve (area = {0:0.2})'.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
↪");
```



ADDING INFORMATION

Here we can improve the results by adding information.

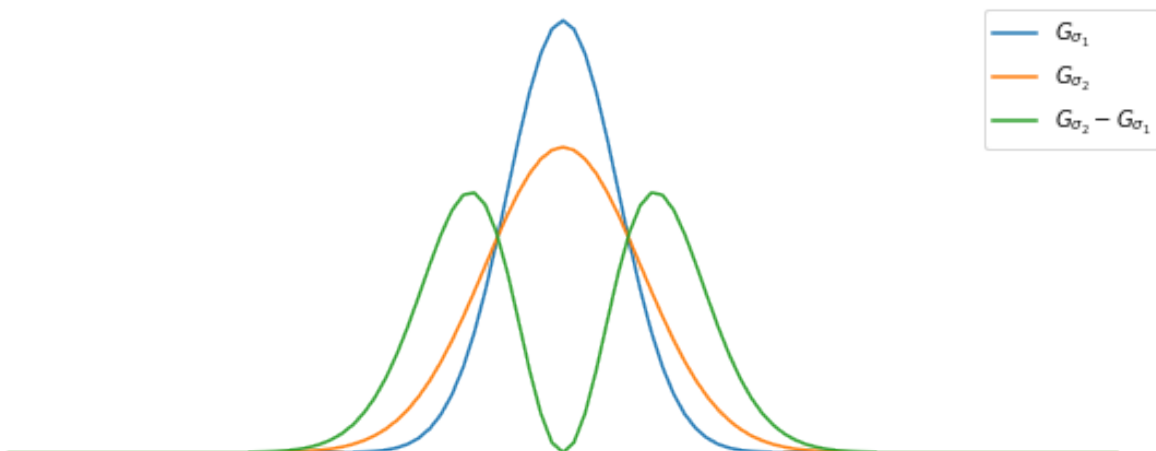
As we discussed in the third lecture (enhancement), edge-enhancing filters can be very useful for classifying images.

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");
```



5.1 Testing with enhanced edges

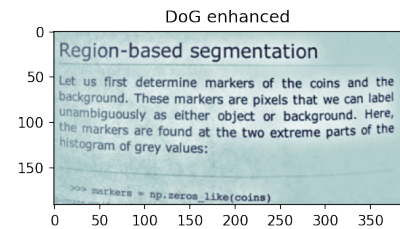
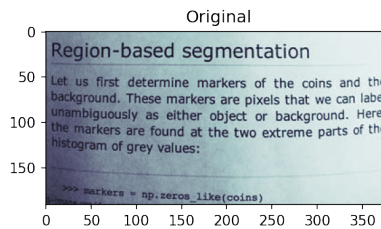
```
def dog_filter(in_img, sig_1, sig_2):
    return gaussian(in_img, sig_1) - gaussian(in_img, sig_2)
```

```
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');
```

	x	y	intensity	is_text	edges
31063	343	80	234	True	0.1
17677	13	46	115	True	0.01
11181	45	29	110	False	-0.04
55242	330	143	229	True	0.0
17891	227	46	218	True	0.05
1345	193	3	202	True	0.02
34692	132	90	157	True	0.05
23682	258	61	241	True	0.07

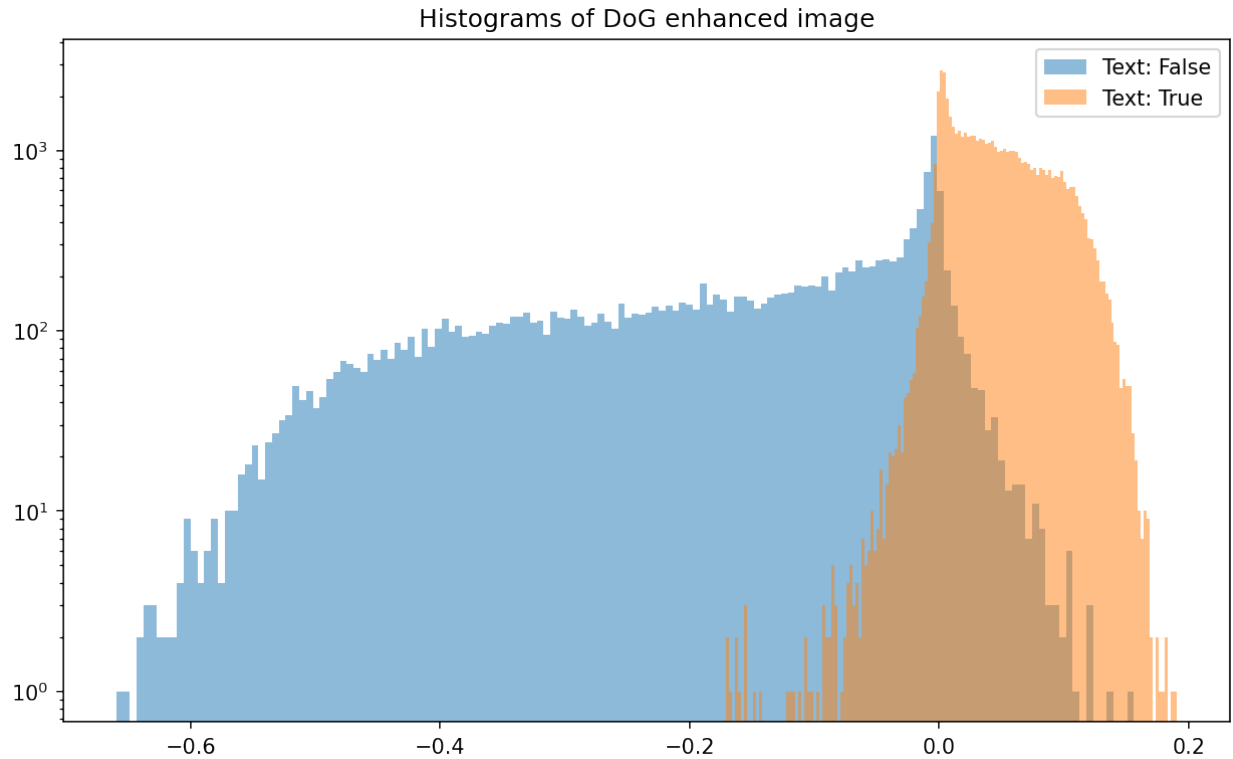


5.2 Histogram of enhanced image

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



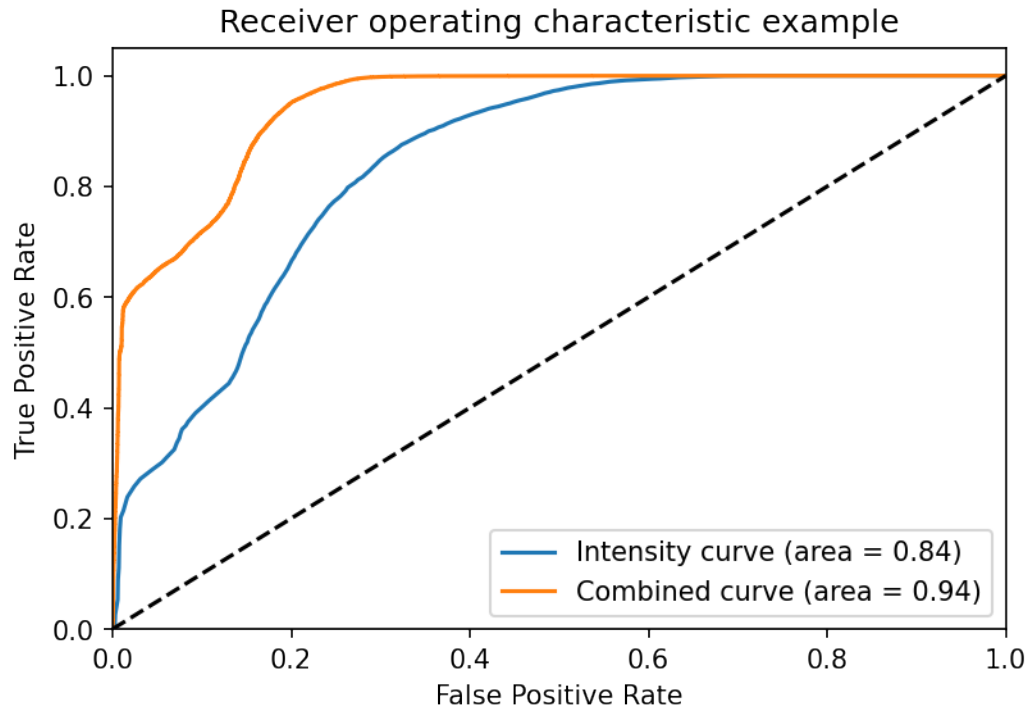
We see here that the text pixels have been compressed into a narrow interval with some tail toward the lower intensities.

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,dpi=150)
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");
```



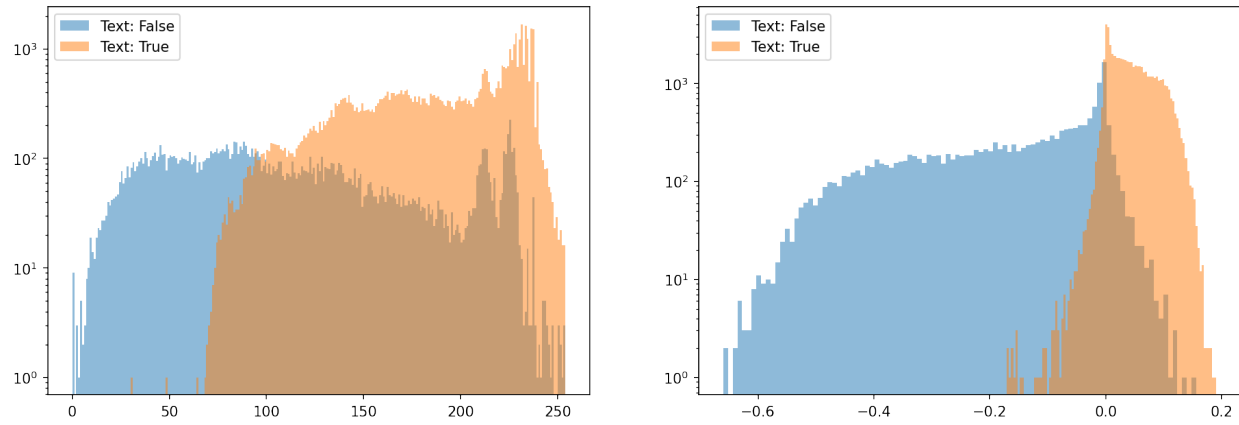
We can see that it is doing better already by looking at the curves. Also the AUC confirms this which is great. We have made an improvement thanks to the preprocessing of the data.

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), dpi=150)
for c_cat, c_df in page_table.groupby(['is_text']):
    ax2.hist(c_df['edges'],
            label = 'Text: {}'.format(c_cat),
            alpha = 0.5, bins=100)
ax2.set_yscale("log", nonpositive='clip')
ax2.legend();

for c_cat, c_df in page_table.groupby(['is_text']):
    ax1.hist(c_df['intensity'],
            np.arange(255),
            label = 'Text: {}'.format(c_cat),
            alpha = 0.5)
ax1.set_yscale("log", nonpositive='clip')
ax1.legend();
```



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.

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

6.1 Example data

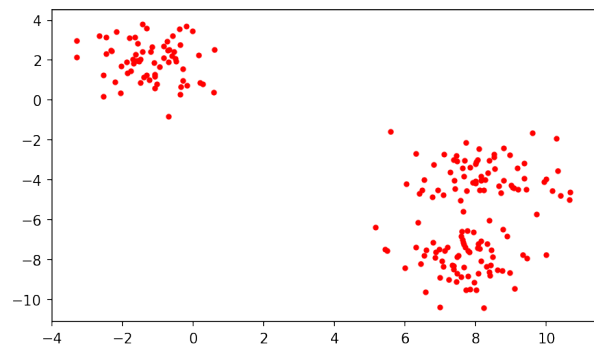
With clustering methods you aim to group data points together into a limited number of clusters. Here, we start to look at an example where each data point has two values. The test data is generated using the `make_blobs` function.

```
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



6.1.1 First clustering attempt

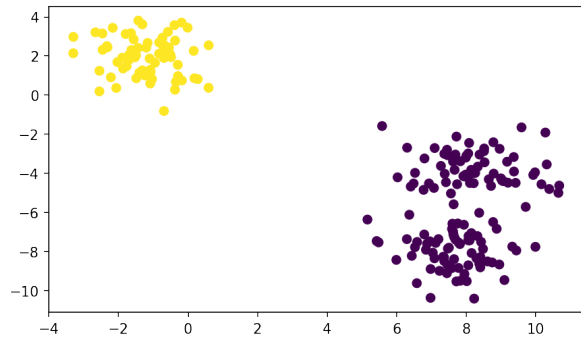
The generated data set has two obvious clusters, but if look closer it is even possible to identify three clusters. We will now use this data to try the k-means algorithm.

```
km = KMeans(n_clusters=2, random_state=2018)
n_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))

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)	7.936151356509988	-6.630712020263951	0.0
(0, 172)	8.062238824854578	-8.70500225647179	0.0
(0, 125)	7.370485145512019	-8.489127187418553	0.0
(0, 105)	6.295207189761017	-7.366814493077998	0.0
(0, 195)	7.720958193819376	-2.1207408297319548	0.0
(1, 72)	-1.0861012238651107	1.2970874865192412	1.0
(1, 129)	-2.166615328033237	3.4441213119929324	1.0
(1, 197)	-2.0652721320971414	0.36827461896146896	1.0
(1, 118)	-1.0774532002573662	1.1922372983022358	1.0
(1, 198)	-0.012441330856637345	3.451248344986084	1.0



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 $\|p_i - c_j\|_2$.

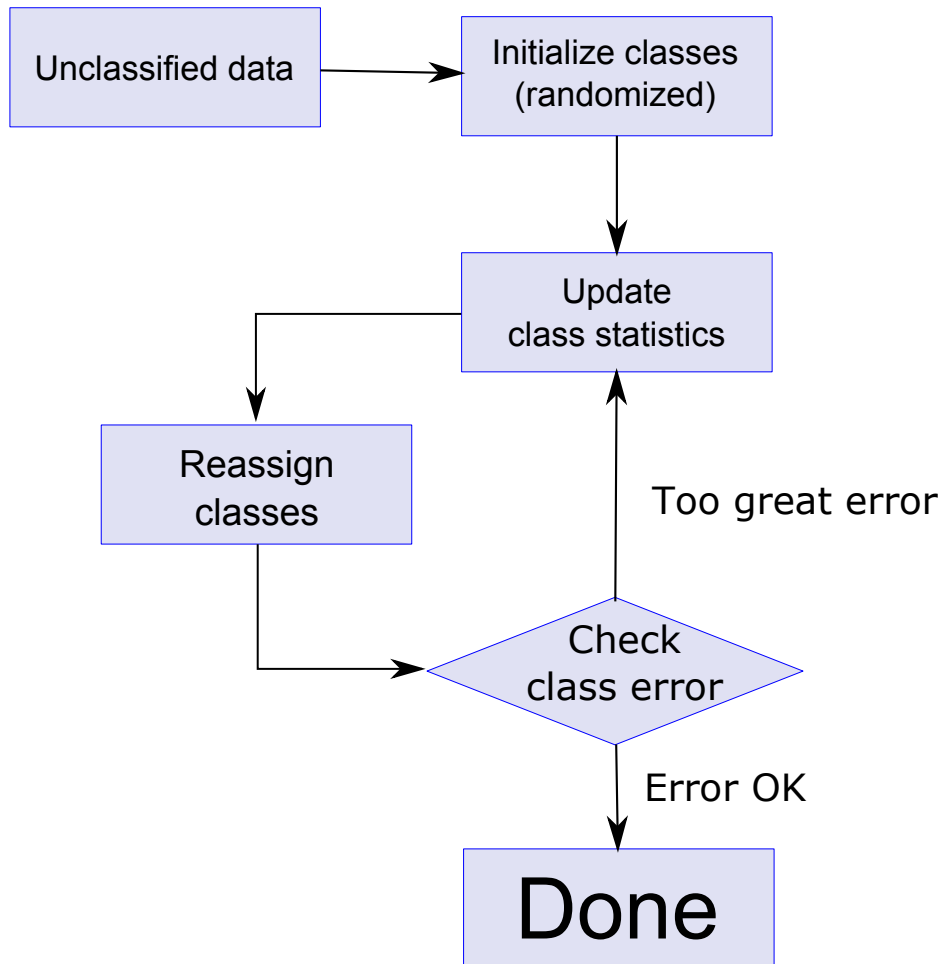


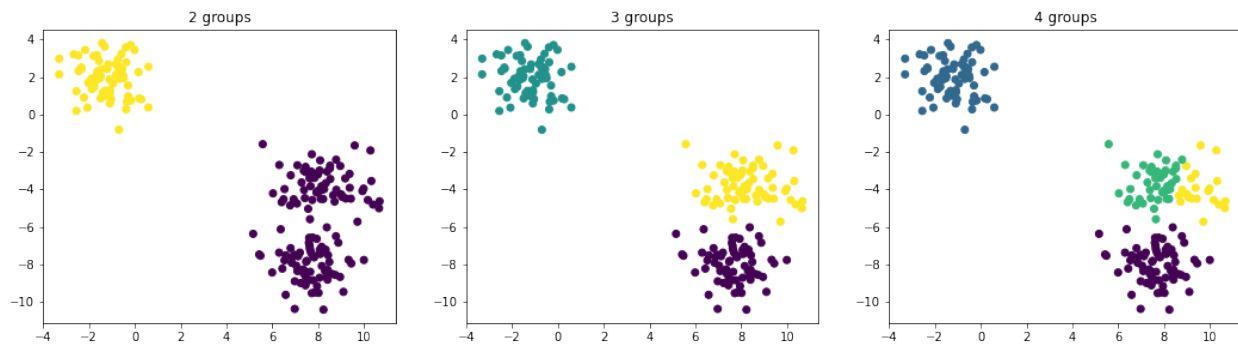
Fig. 6.1: Flow chart for the k-means clustering iterations.

6.2.1 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.

```
N=3
fig, ax = plt.subplots(1,N,figsize=(18,4.5))

for i in range(N) :
    km = KMeans(n_clusters=i+2, random_state=2018); n_grp = km.fit_predict(test_pts)
    ax[i].scatter(test_pts.x, test_pts.y, c=n_grp)
    ax[i].set_title('{0} groups'.format(i+2))
```



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 how ever 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.

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)

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, : :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]))
```

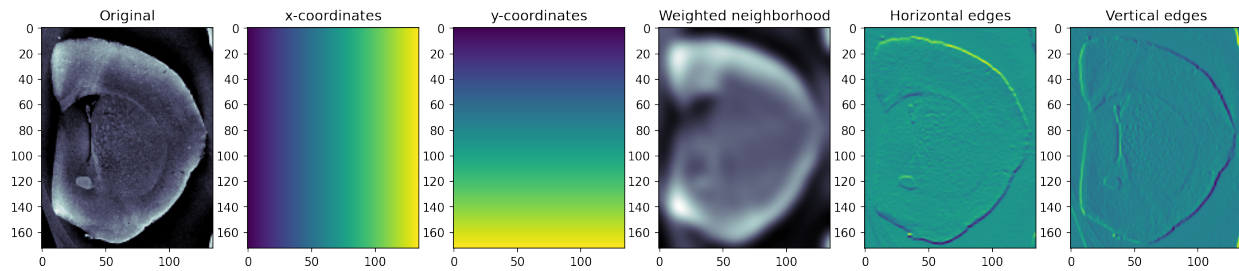
(continues on next page)

(continued from previous page)

```

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');

```



6.4.1 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

```

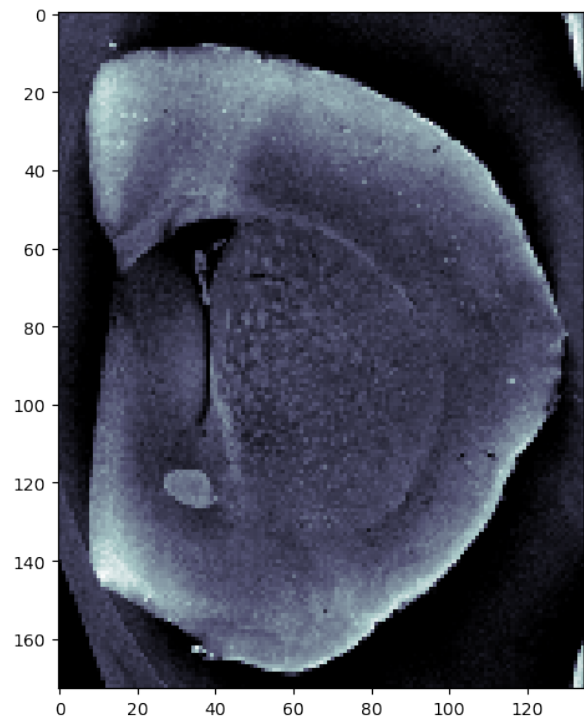
```

cortex_img = imread("figures/cortex.png")[:,::3, ::3]/1000.0
np.random.seed(2018)
fig, (ax0,ax1) = plt.subplots(1, 2,
                             figsize=(12, 8), dpi=100)
ax1.imshow(cortex_img, cmap='bone')
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()))
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');

```

	x	y	intensity
18621	126.0	137.0	6.015
8830	55.0	65.0	18.961
23217	132.0	171.0	57.508
7605	45.0	56.0	5.81
13164	69.0	97.0	18.338
13064	104.0	96.0	16.929
8397	27.0	62.0	0.388
3556	46.0	26.0	35.116
7075	55.0	52.0	21.122
8819	44.0	65.0	15.514



6.4.2 First segmentation attempt with k-means on brain image

We use $N_{clusters} = 4$

```

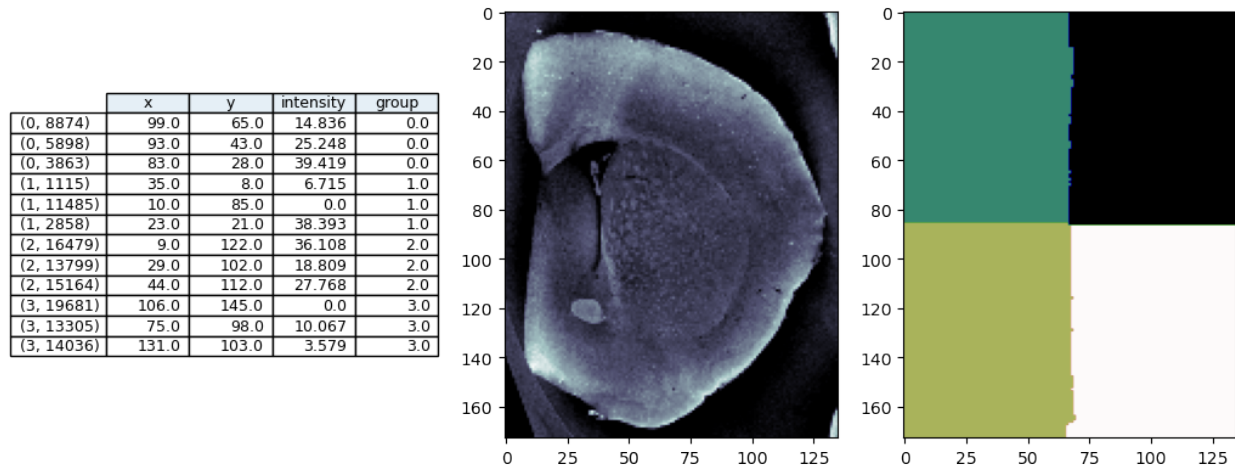
from sklearn.cluster import KMeans
km = KMeans(n_clusters=4, random_state=2018)
cortex_df['group'] = km.fit_predict(cortex_df[['x', 'y', 'intensity']].values)

fig, (ax0, ax1, ax2) = plt.subplots(1, 3,
                                   figsize=(12, 8), dpi=100)
ax1.imshow(cortex_img, cmap='bone')
ax2.imshow(cortex_df['group'].values.reshape(
    cortex_img.shape), 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');

```



6.4.3 Why is the image segmented like this?

6.4.4 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 $\vec{v} = \{\frac{x}{10}, \frac{y}{10}, \text{Intensity}\}$

$N_{clusters}=4$

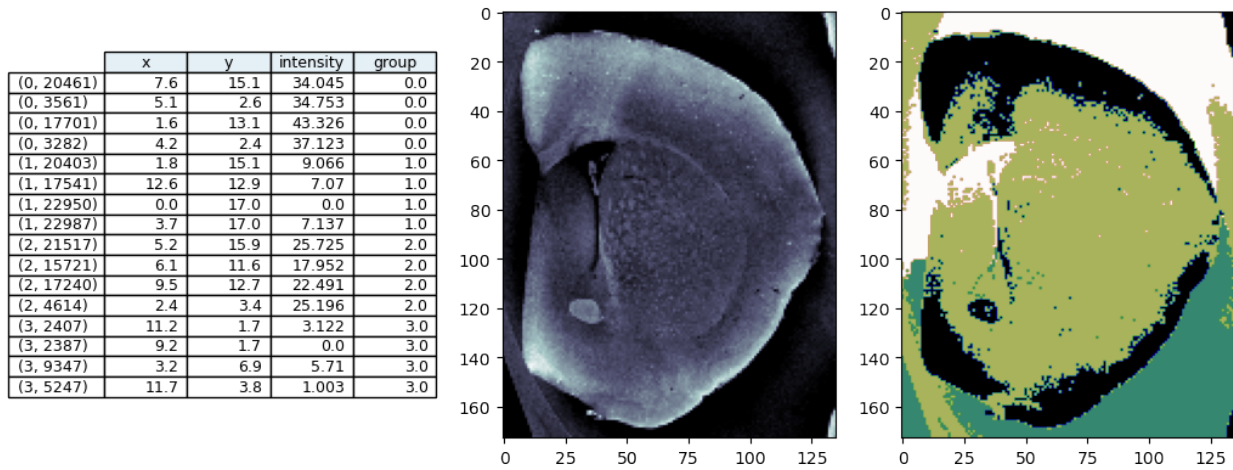
```

km = KMeans(n_clusters=4, random_state=2018)
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)

fig, (ax0,ax1, ax2) = plt.subplots(1, 3,
    figsize=(12, 8), dpi=100)
ax1.imshow(cortex_img, cmap='bone')
ax2.imshow(scale_cortex_df['group'].values.reshape(cortex_img.shape),
    cmap='gist_earth')

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

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

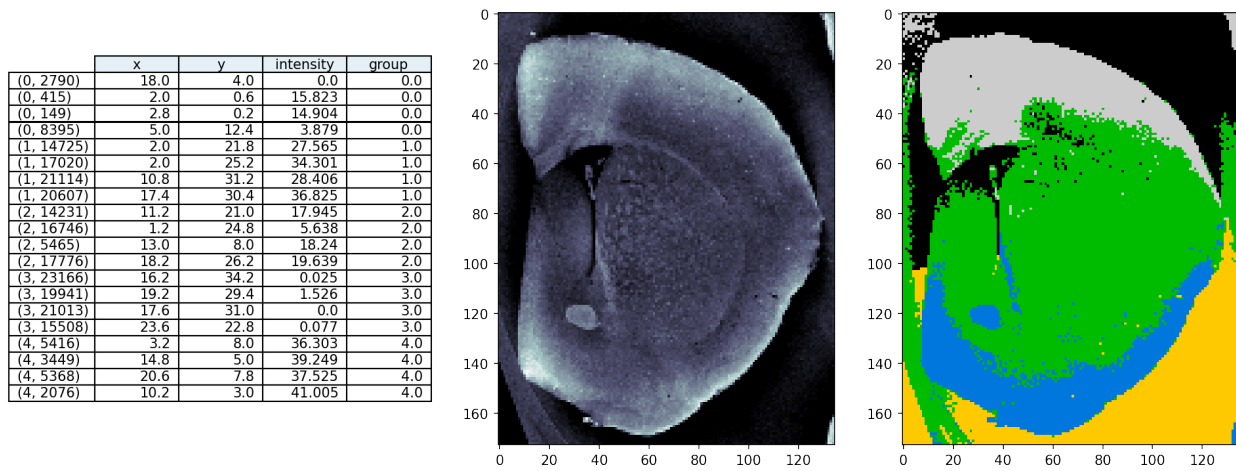
$N_{clusters}=5$

```

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');
    
```

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 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 a great combination as these data are

- Hyper-spectral data (materials science example)

Many materials have a characteristic response 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.

QUAD TREES

Split the image in subregions until a criterion is fulfilled:

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.

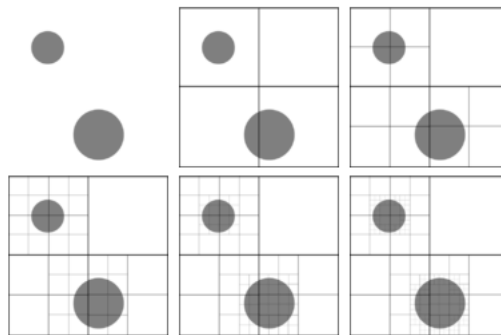


Fig. 7.1: A quad tree decomposition.

7.2 Quad tree example

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

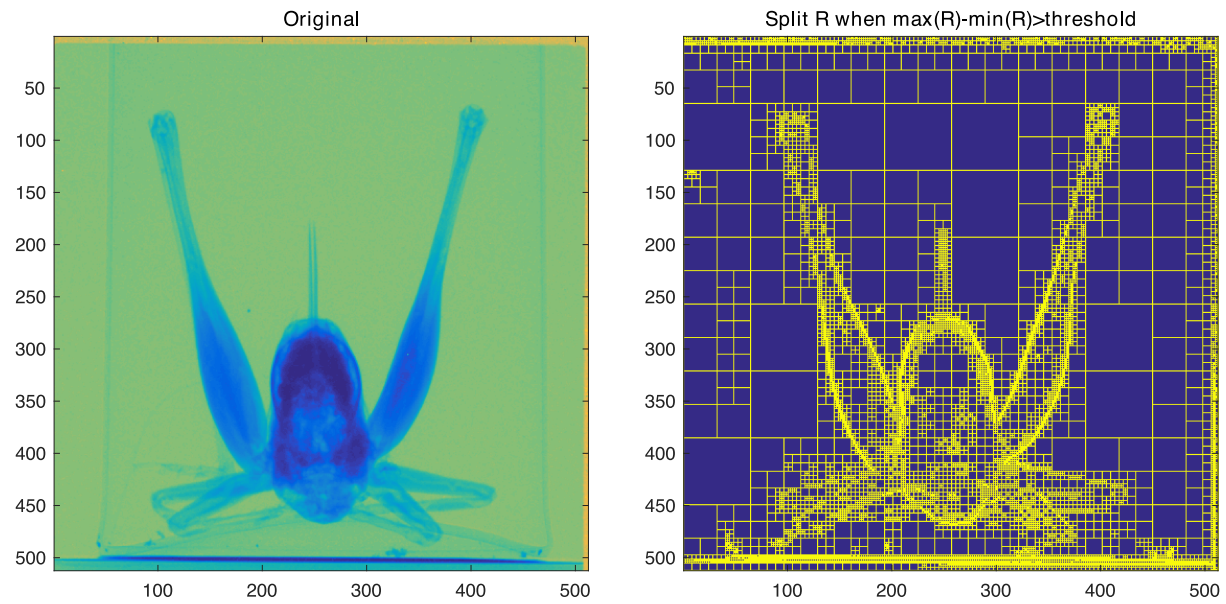


Fig. 7.2: A neutron radiograph of a grasshopper decomposed using a quad tree.

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

8.1 Why use superpixels

Super pixels

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

8.1.1 A super-pixel example

We start with an example of shale with multiple phases

- rock
- clay
- pore

8.1.2 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()])
```

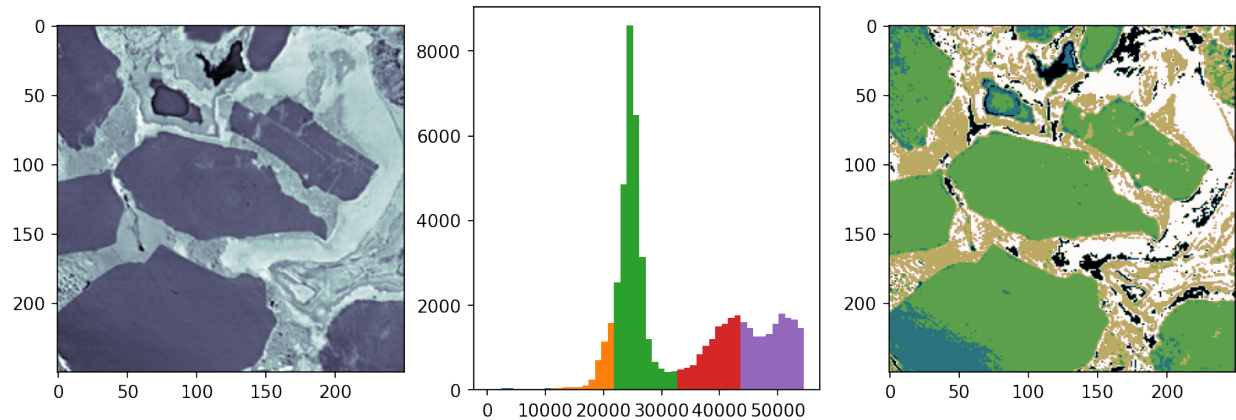
(continues on next page)

(continued from previous page)

```

out_img[thresh_reg] = i
ax3.imshow(out_img, cmap='gist_earth');

```



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

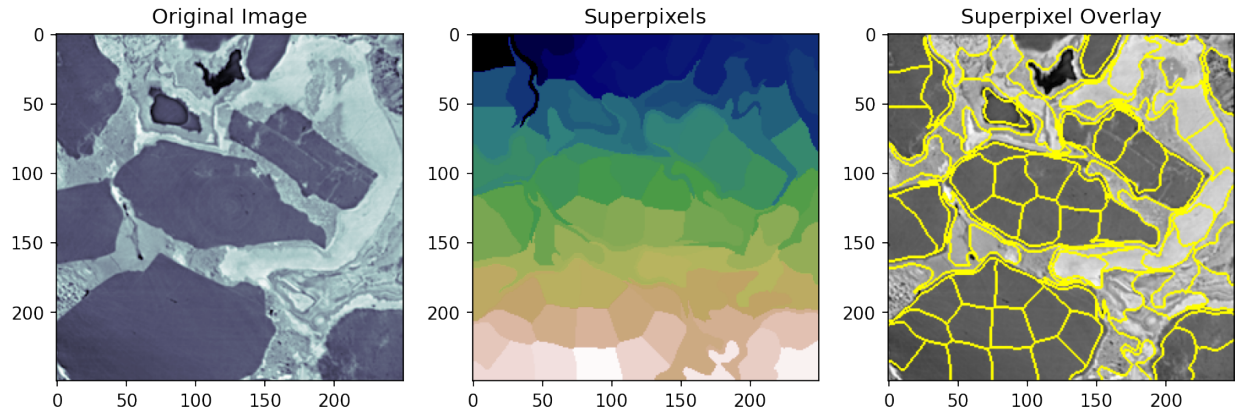
```

from skimage.segmentation import slic, mark_boundaries

shale_segs = slic(shale_img,
                  n_segments=100,
                  compactness=5e-2,
                  sigma=3.0,
                  start_label=1)

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');

```

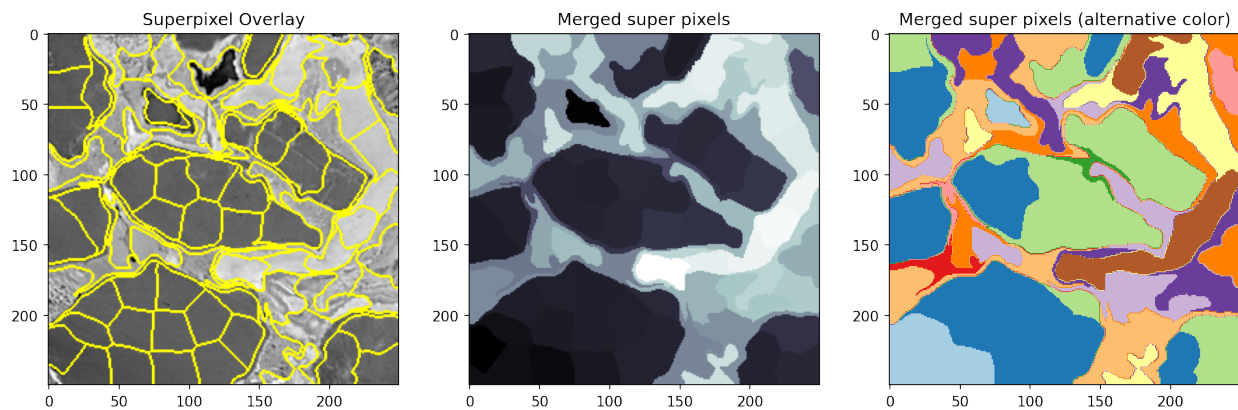


8.3 Merging super pixels

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])

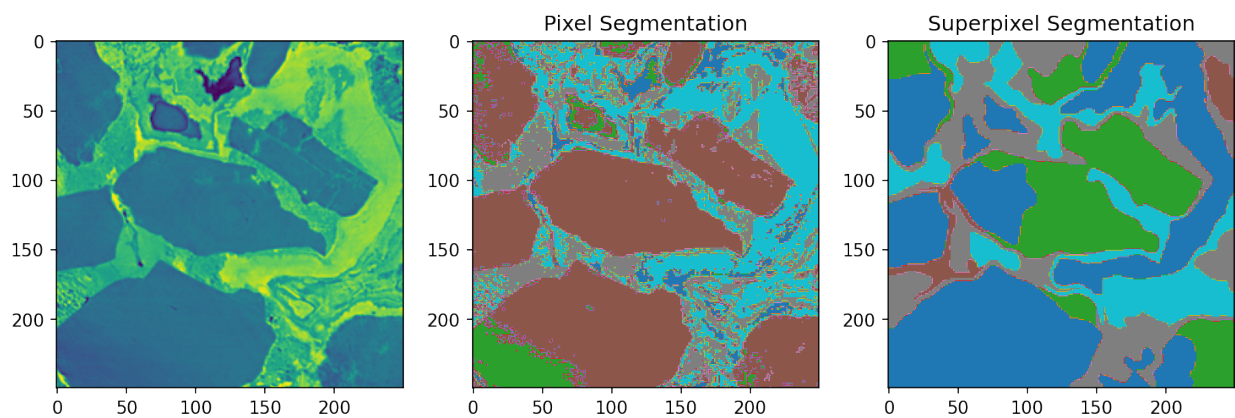
fig, (ax1,ax2,ax3) = plt.subplots(1, 3, figsize=(15, 5), dpi=150)
ax1.imshow(mark_boundaries(shale_img, shale_segs))
ax1.set_title('Superpixel Overlay')
ax2.imshow(flat_shale_img, cmap='bone'); ax2.set_title('Merged super pixels')
ax3.imshow(flat_shale_img, cmap='Paired'); ax3.set_title('Merged super pixels -
↪(alternative color)');
```



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')
ax2.set_title('Pixel Segmentation')
ax3.imshow(sp_out_img, cmap='tab10')
ax3.set_title('Superpixel Segmentation');
```



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 α and β

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

- $P(\{\vec{x}, f(\vec{x})\}|\alpha)$ the probability a given pixel is in phase α given we know it's 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)

SUMMARY

10.1 Histogram based thresholding

- Otsu and others
- Hysteresis thresholding

10.2 Clustering - K-means

- Add position information

10.3 Similar region segmentations

- Quad trees
- Super pixels

SUPERVISED SEGMENTATION APPROACHES

11.1 Overview

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

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

11.2.1 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 DecisionTreeClassifier
import graphviz
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
```

(continues on next page)

(continued from previous page)

```
from sklearn.cluster import KMeans

%matplotlib inline

plt.rcParams["figure.figsize"] = (8, 8)
plt.rcParams["figure.dpi"] = 150
plt.rcParams["font.size"] = 14
plt.rcParams['font.family'] = ['sans-serif']
plt.rcParams['font.sans-serif'] = ['DejaVu Sans']
plt.style.use('ggplot')
sns.set_style("whitegrid", {'axes.grid': False})
```

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.

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.

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.

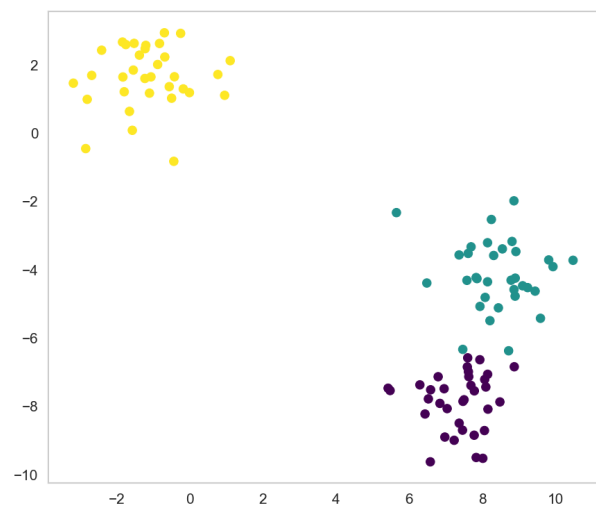
CLASSIFICATION

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
27	-2.81	0.99	2.0
15	8.1	-7.43	0.0
14	7.86	-4.26	1.0
62	7.47	-7.85	0.0
93	-1.84	2.67	2.0
52	0.95	1.11	2.0
5	-1.65	0.64	2.0
9	8.15	-4.35	1.0
66	8.02	-9.52	0.0
71	-0.25	2.93	2.0



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.

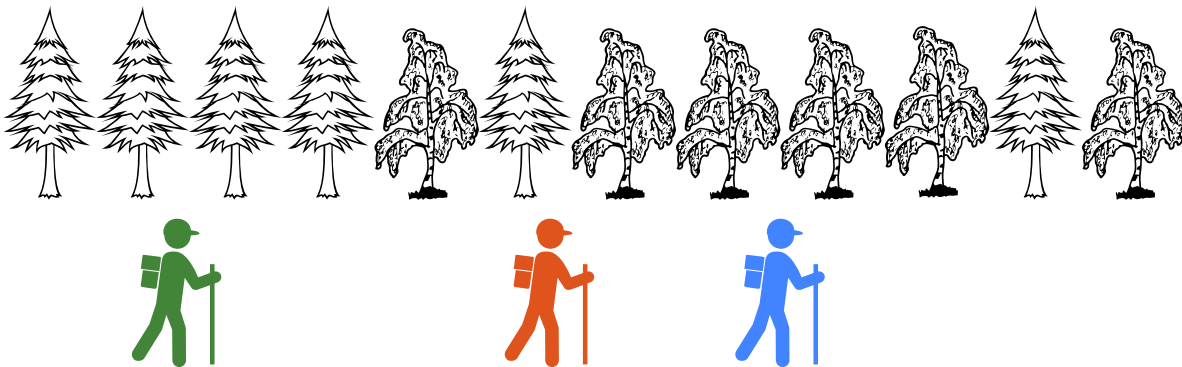
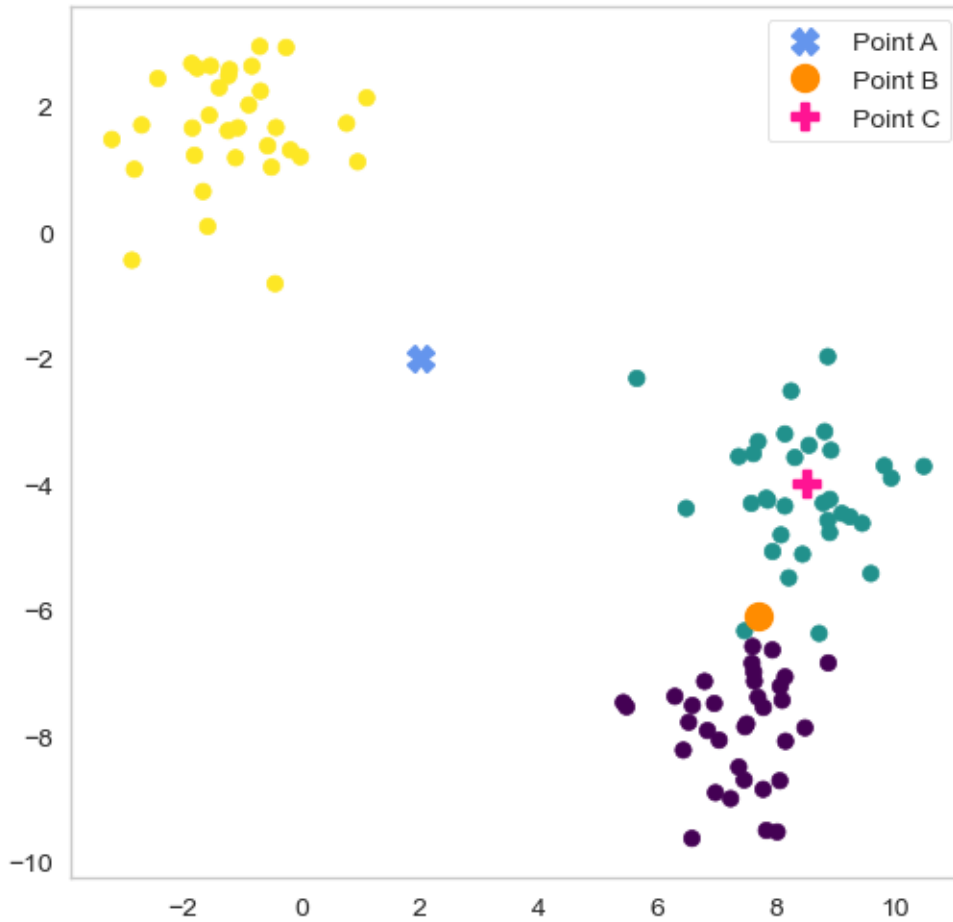


Fig. 13.1: Depending on the where the hiker is standing he makes the the conclusion that is either in a birch or spruce forest.

```
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();
```



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)
```

13.4 Nearest neighbor predictions

13.4.1 Basic test

Same values as training

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

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

13.4.2 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']
```

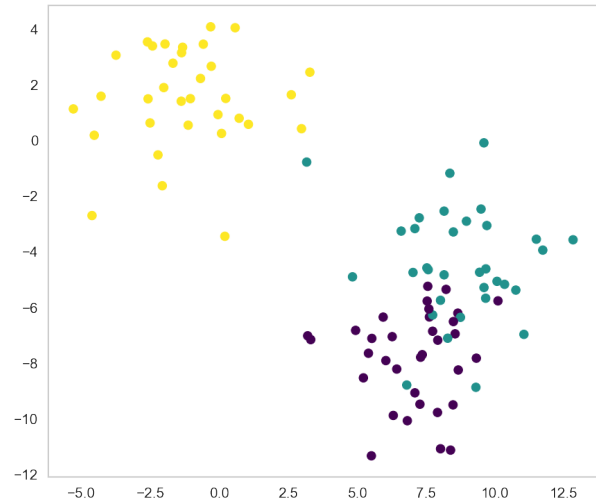
13.5 Let's come back to the blob data

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

```
blob_data, blob_labels = make_blobs(n_samples=100,
                                   cluster_std=2.0,
                                   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
83	8.49	-3.28	1.0
46	9.6	-0.08	1.0
48	-4.62	-2.7	2.0
6	-0.05	0.93	2.0
52	2.98	0.42	2.0
63	-0.68	2.23	2.0
71	0.57	4.05	2.0
86	4.83	-4.9	1.0
88	9.32	-7.82	0.0
44	0.08	0.25	2.0



13.6 Training the model using one neighbor

```
# Define classifier
k_class = KNeighborsClassifier(1)

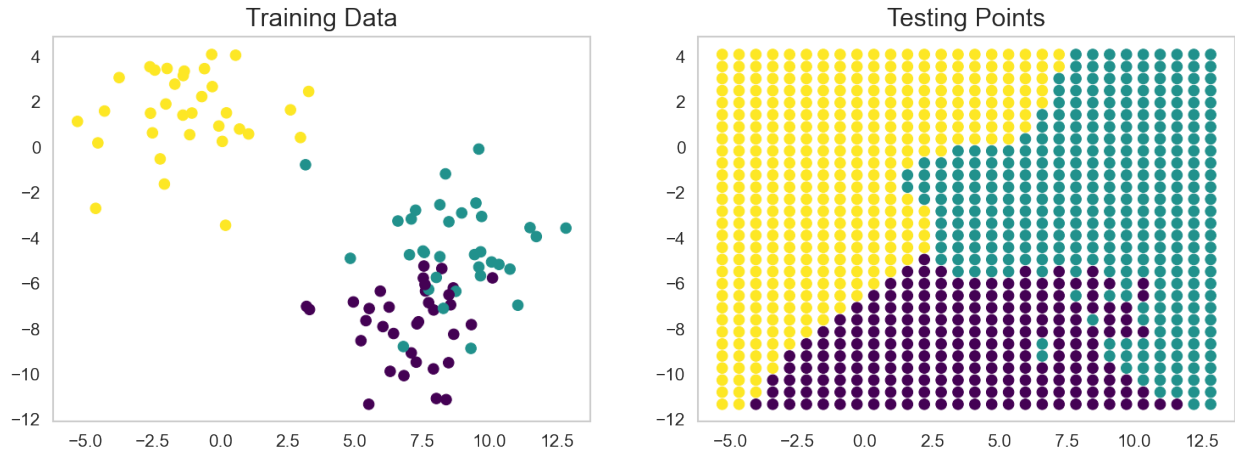
# Train the classifier model with data
k_class.fit(test_pts[['x', 'y']], test_pts['group_id'])
```

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

13.6.1 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');
```



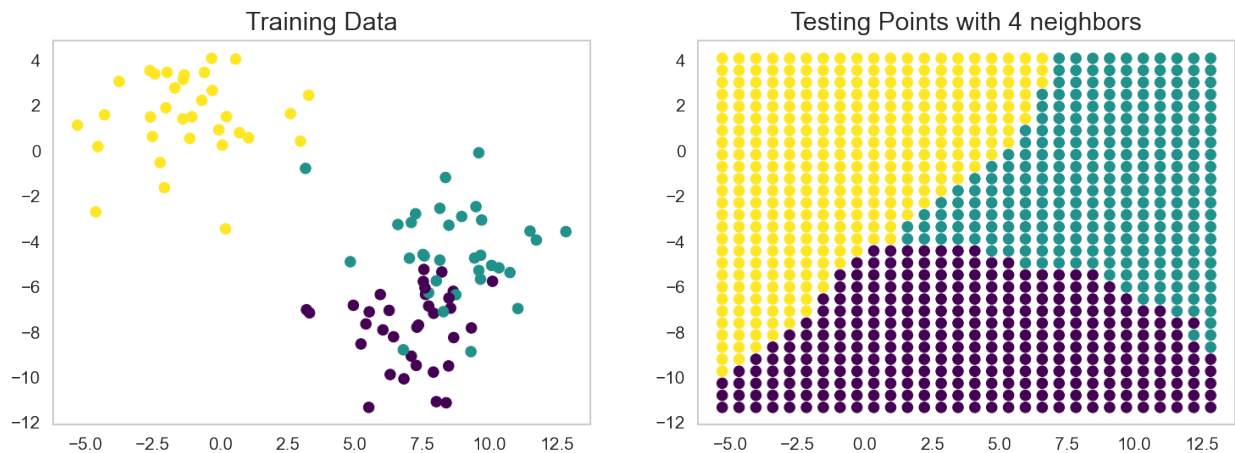
13.7 Stabilizing Results - Increase number of neighbors

- 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');
    
```



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

14.1 We need data to fit...

```
from sklearn.linear_model import LinearRegression

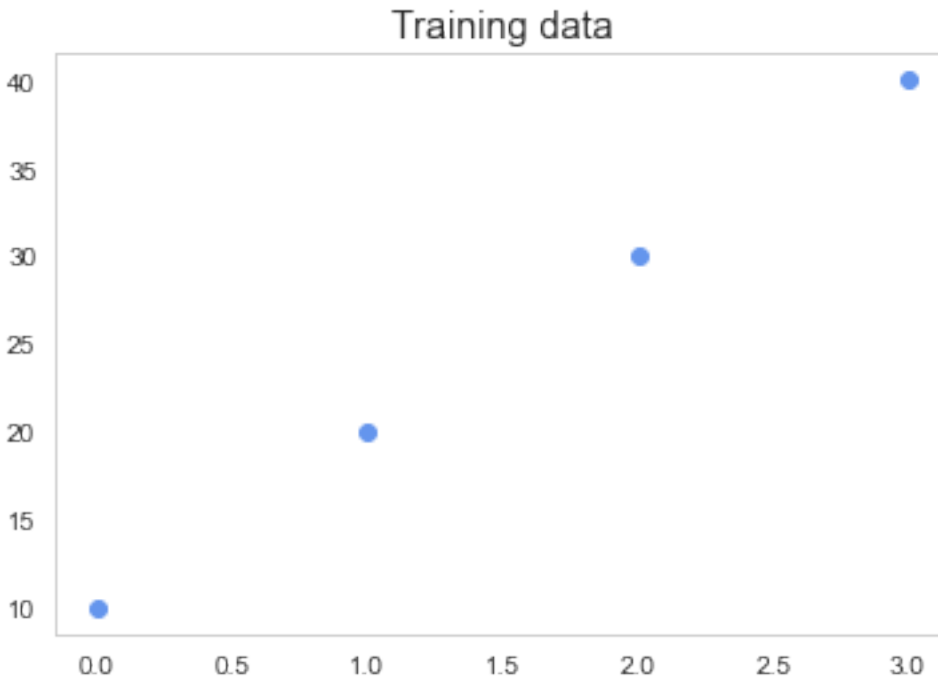
x = np.array([0, 1, 2, 3])
y = np.array([10, 20, 30, 40] )

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
```



14.1.1 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.]
```

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')

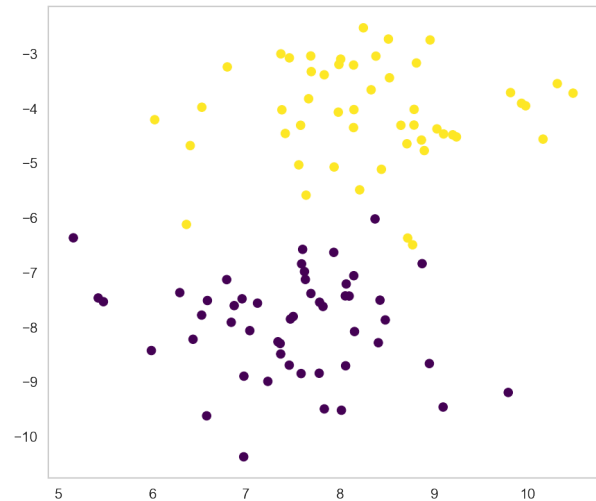
ccolors = plt.cm.BuPu(np.full(3, 0.1))
```

(continues on next page)

(continued from previous page)

```
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
51	9.94	-3.9	1.0
26	7.23	-8.99	0.0
31	6.8	-7.13	0.0
62	9.1	-9.46	0.0
21	7.59	-6.84	0.0
67	7.34	-8.27	0.0
58	7.94	-5.07	1.0
82	10.49	-3.72	1.0
20	8.49	-7.87	0.0
93	8.15	-4.35	1.0



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 1.3469297085944607
```

14.3.1 Evaluate the regression model

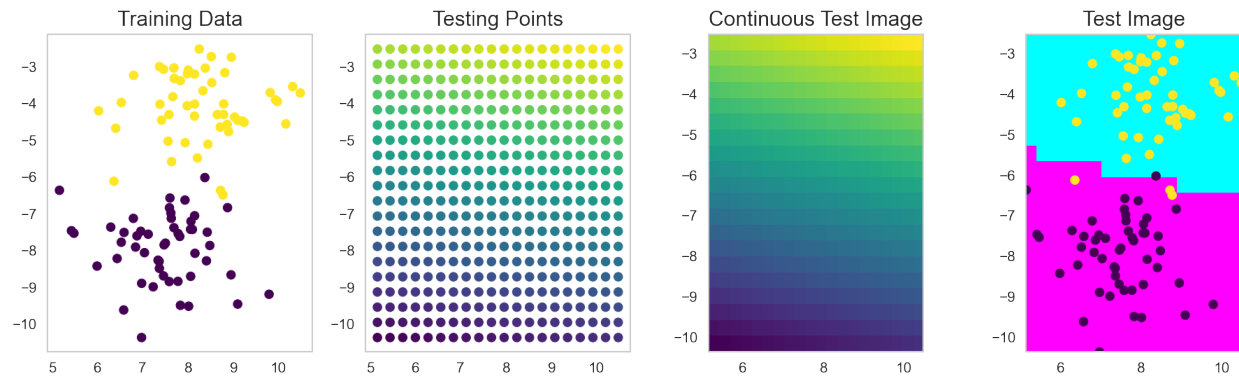
```
xx, yy = np.meshgrid(np.linspace(test_pts.x.min(), test_pts.x.max(), 20),
    ↪np.linspace(test_pts.y.min(), test_pts.y.max(), 20),
    ↪indexing='ij'
    ↪)
grid_pts = pd.DataFrame(dict(x=xx.ravel(), y=yy.ravel()))
grid_pts['predicted_id'] = l_reg.predict(grid_pts[['x', 'y']])
```

```
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(
    ↪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');
```

(continues on next page)

(continued from previous page)

```
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');
```



DECISION TREES

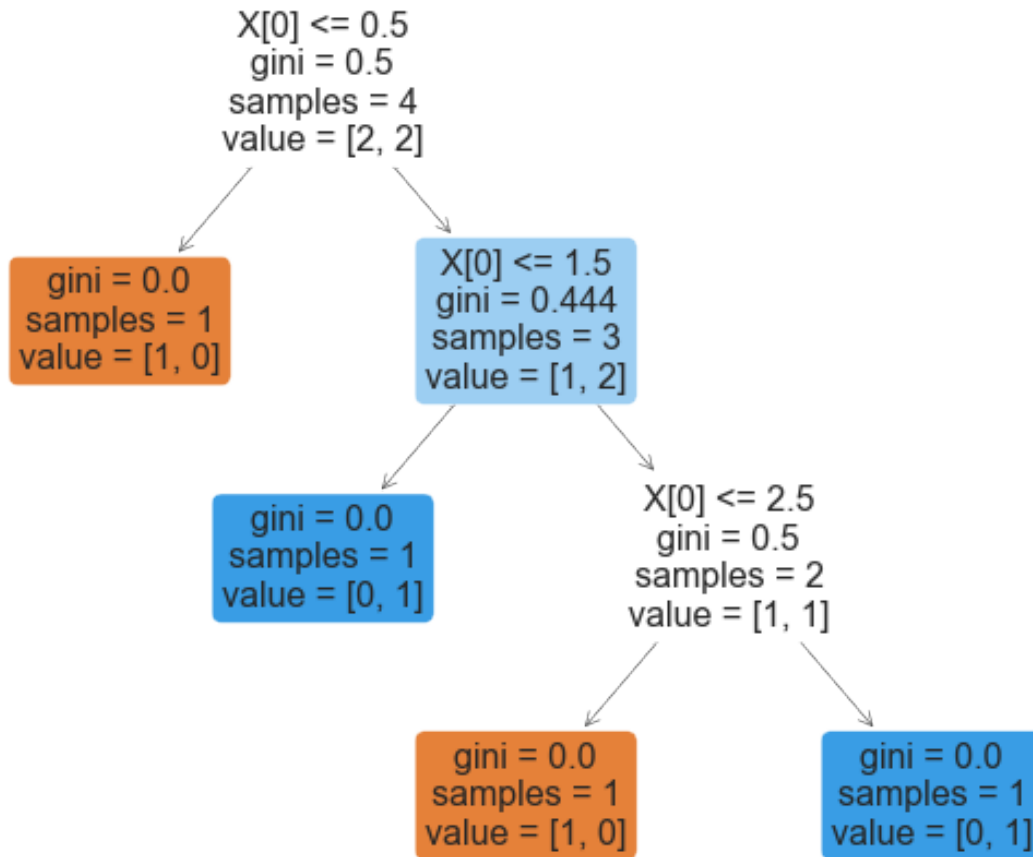
- SciKit Learn documentation on trees
- SciKit Learn trees explained
- Hastie et al., Elements Of Statistical Learning, 2009 Section 9.2 Trees.

```
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))
```

15.1 Create a decision tree classifier

```
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);
```



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

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

```

```

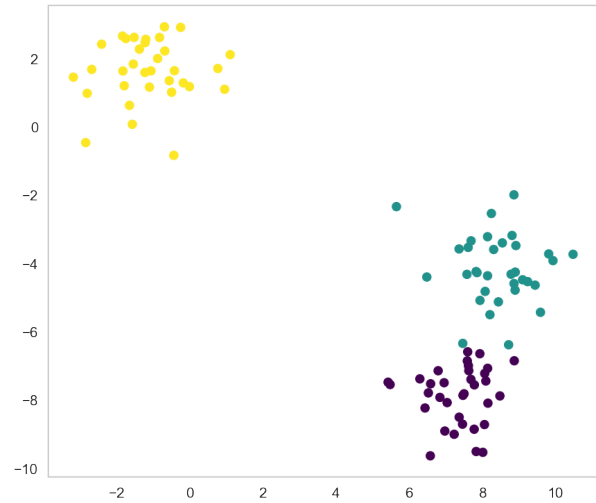
blob_data, blob_labels = make_blobs(n_samples=100, random_state=2018)
test_pts = pd.DataFrame(blob_data, columns=['x', 'y'])
test_pts['group_id'] = blob_labels

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');

```

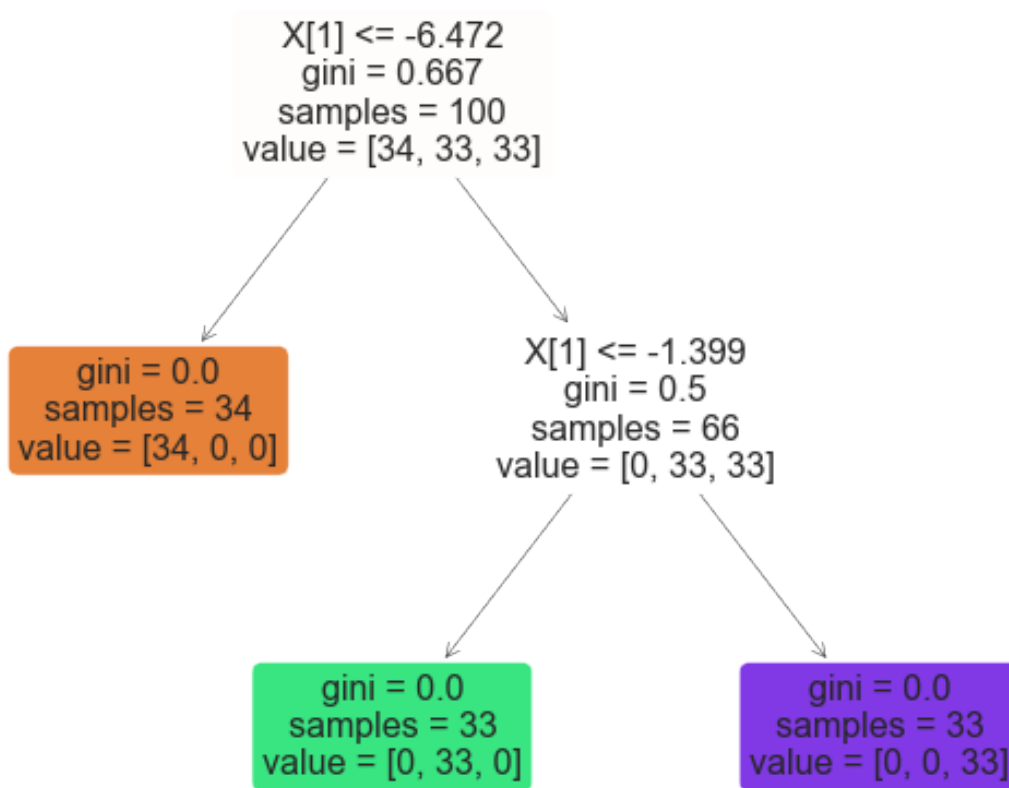
	x	y	group_id
44	-0.5	1.03	2.0
53	8.82	-3.17	1.0
57	7.37	-8.49	0.0
42	-1.79	1.22	2.0
52	0.95	1.11	2.0
98	7.62	-6.98	0.0
24	-0.83	2.63	2.0
46	8.87	-1.98	1.0
0	8.55	-3.39	1.0
13	6.96	-7.48	0.0



15.2.1 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);
```



15.2.2 Let's look at the decision map

```

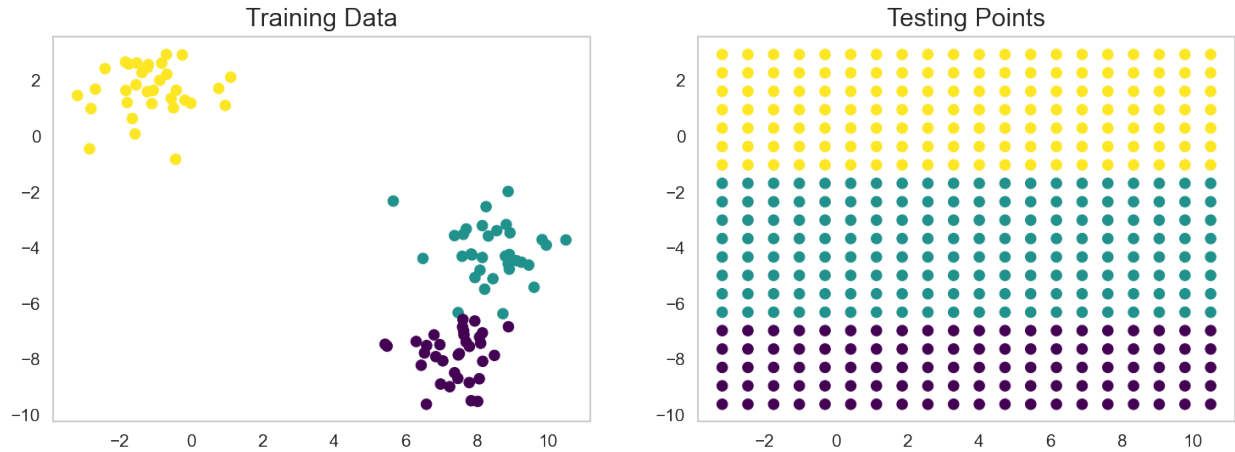
xx, yy = np.meshgrid(np.linspace(test_pts.x.min(), test_pts.x.max(), 20),
                    np.linspace(test_pts.y.min(), test_pts.y.max(), 20),
                    indexing='ij')
)
grid_pts = pd.DataFrame(dict(x=xx.ravel(), y=yy.ravel()))
grid_pts['predicted_id'] = d_tree.predict(grid_pts[['x', 'y']])

```

```

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');

```



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](#)

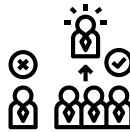
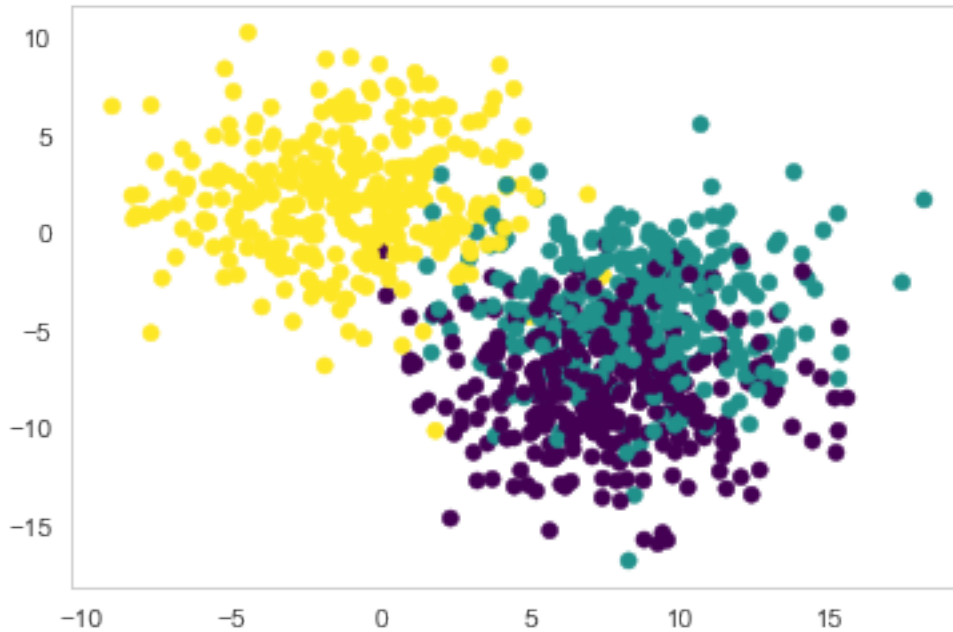


Fig. 15.1: Random forests train trees on different fractions of the data.

```
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');
```



15.3.1 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'])

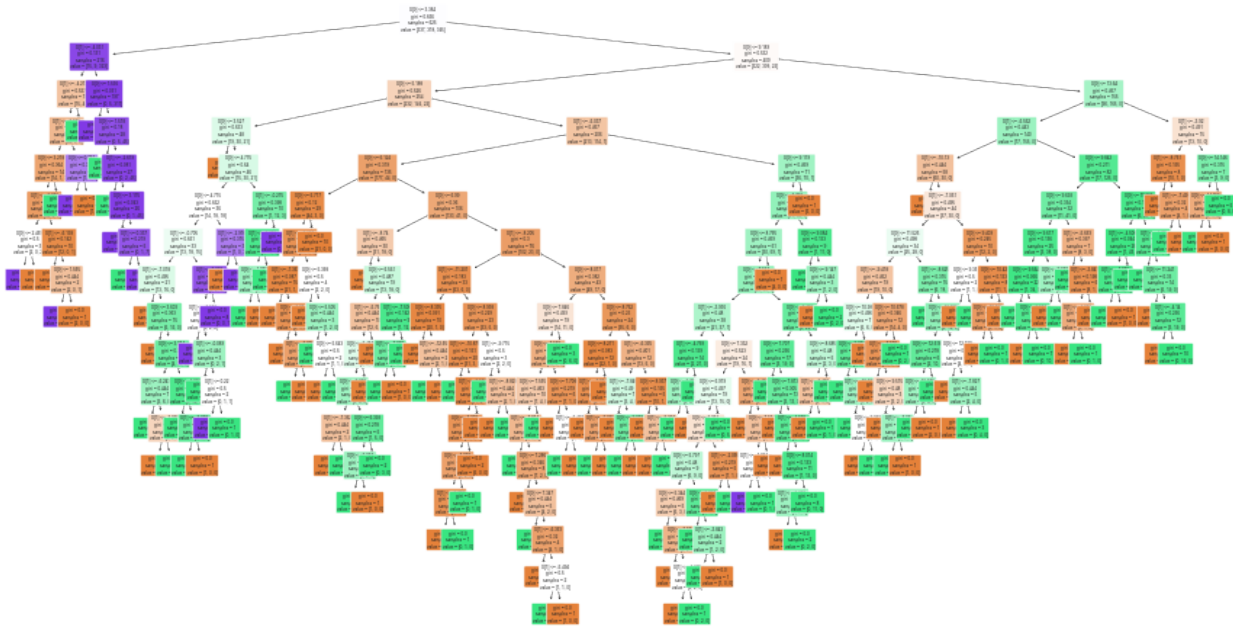
print('Build ', len(rf_class.estimators_), 'decision trees')
```

```
Build 5 decision trees
```

15.3.2 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')

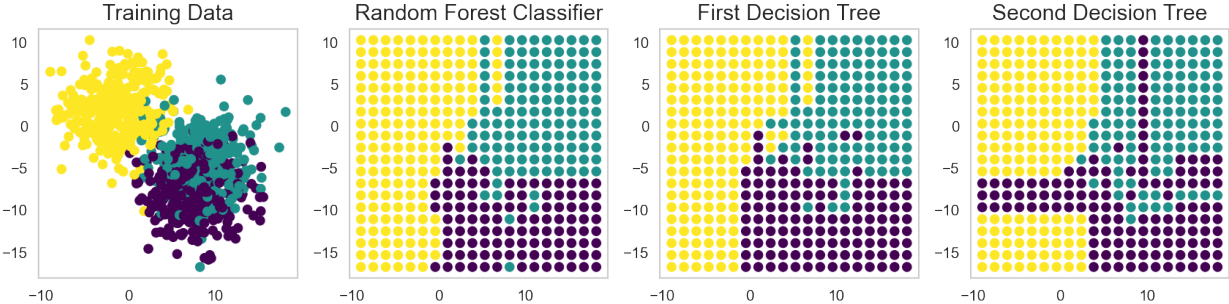
grid_pts = pd.DataFrame(dict(x=xx.ravel(), y=yy.ravel()))
rnd_forest_classes = rf_class.predict(grid_pts[['x', 'y']]) # Prediction of the forest

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');
```



PIPELINES

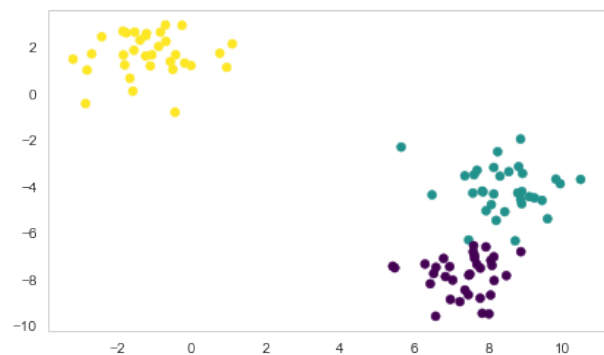
We will use the idea of pipelines generically here to refer to the combination of steps that need to be performed to solve a problem.

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

16.1 Let's the 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');
```

x	y	group_id
-1.8289749226285306	1.6482165325889067	2.0
8.869061857524883	-1.978022202334015	1.0
7.692387397193145	-7.382755171300531	0.0
7.633552458215117	-7.1264027670381855	0.0
8.146828793711666	-4.348565130067033	1.0

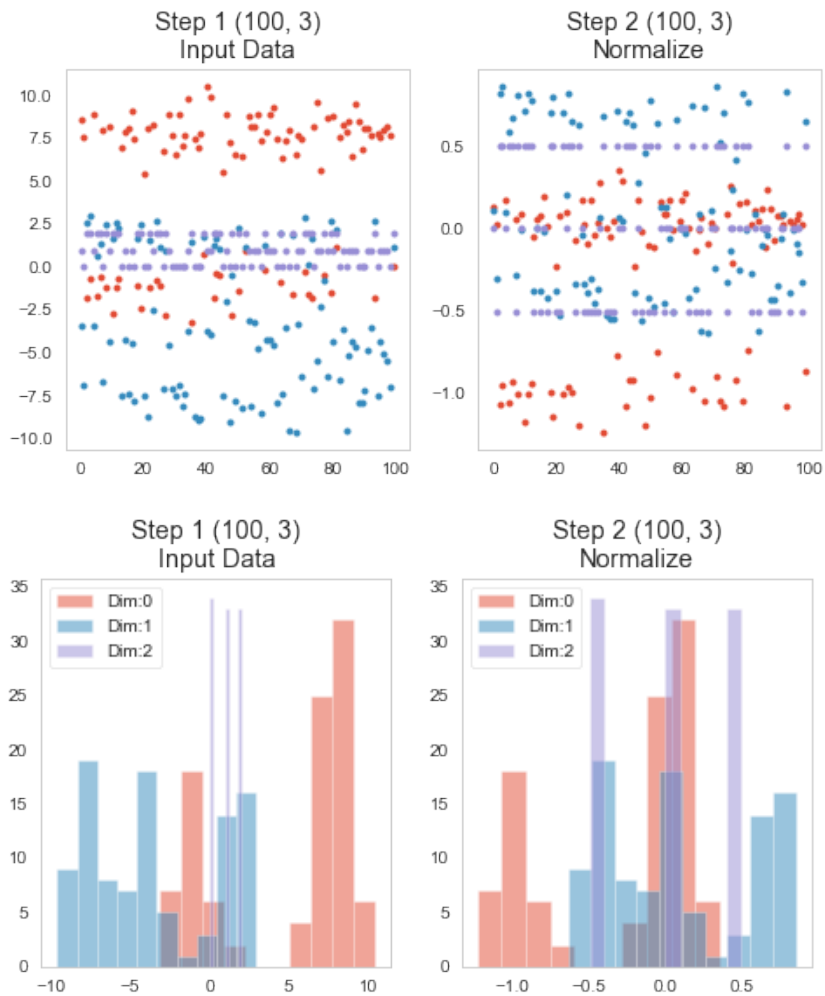


16.2 A basic pipeline

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)
```

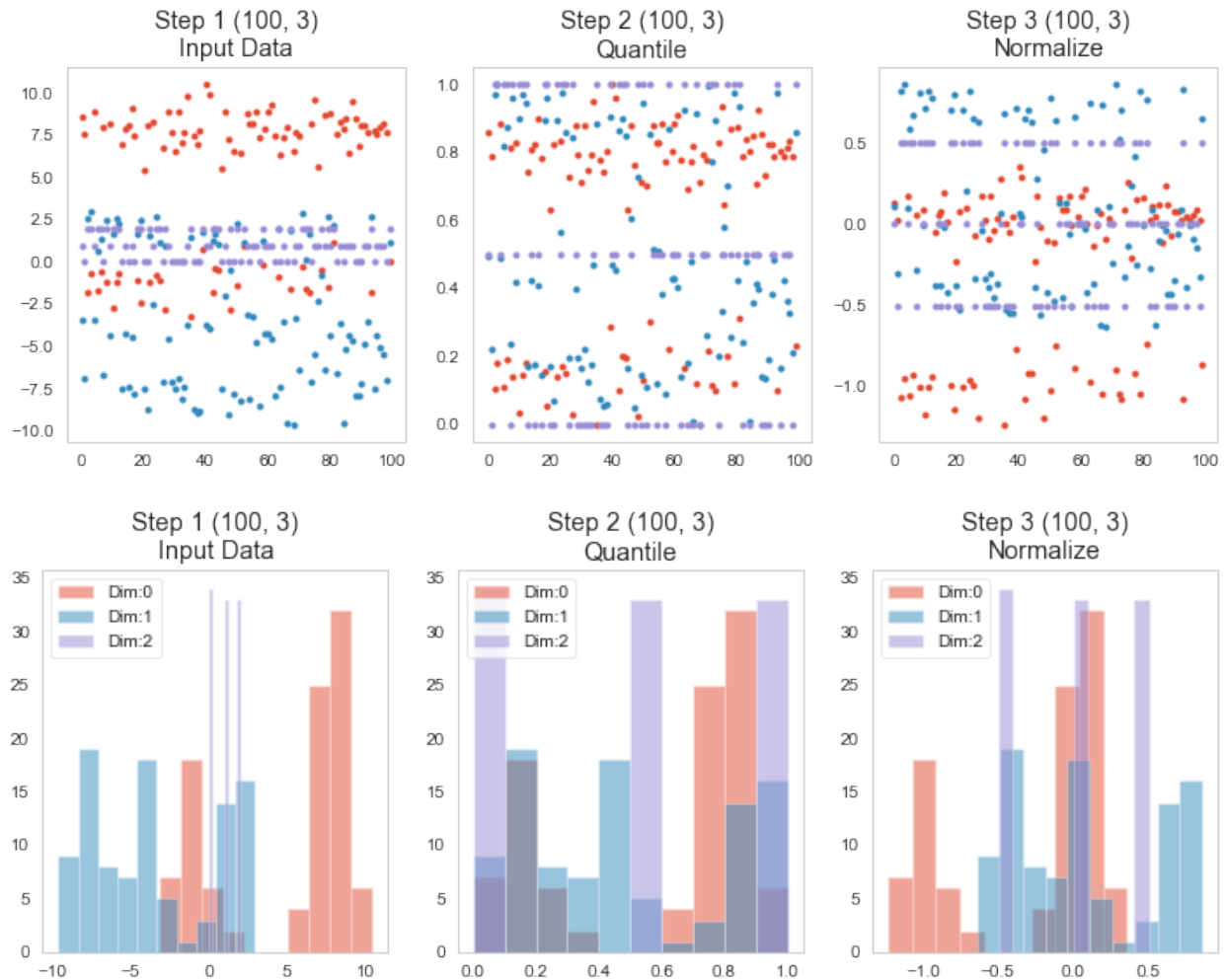


16.2.1 Adding tasks to the pipeline

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)
```

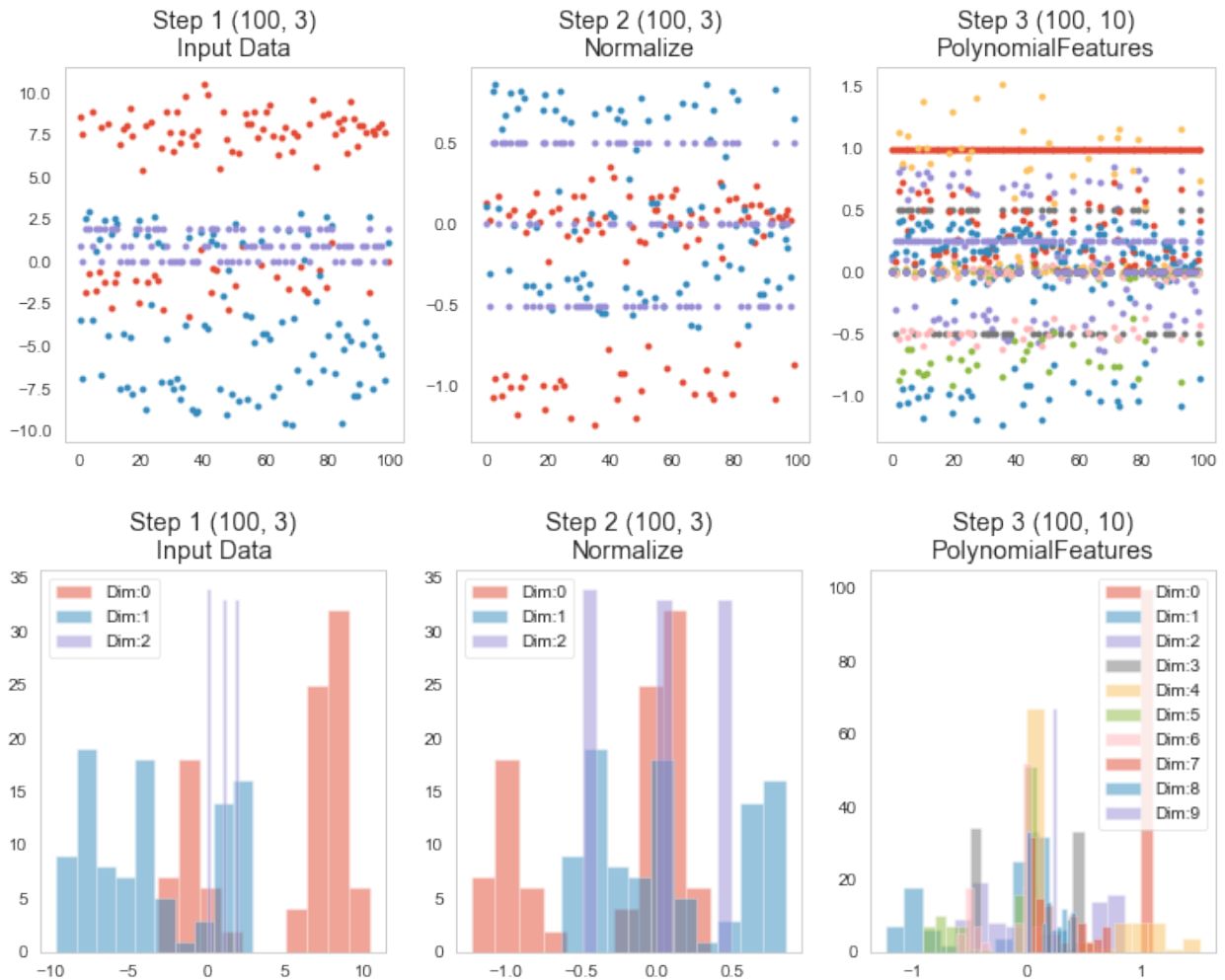


16.2.2 Adding polynomial features to the pipeline

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)
```



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.

17.1 A first example with pipeline processing

17.1.1 Let's load some images

- Images of numbers 0 to 9 (MNIST)
- 8×8 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)
```

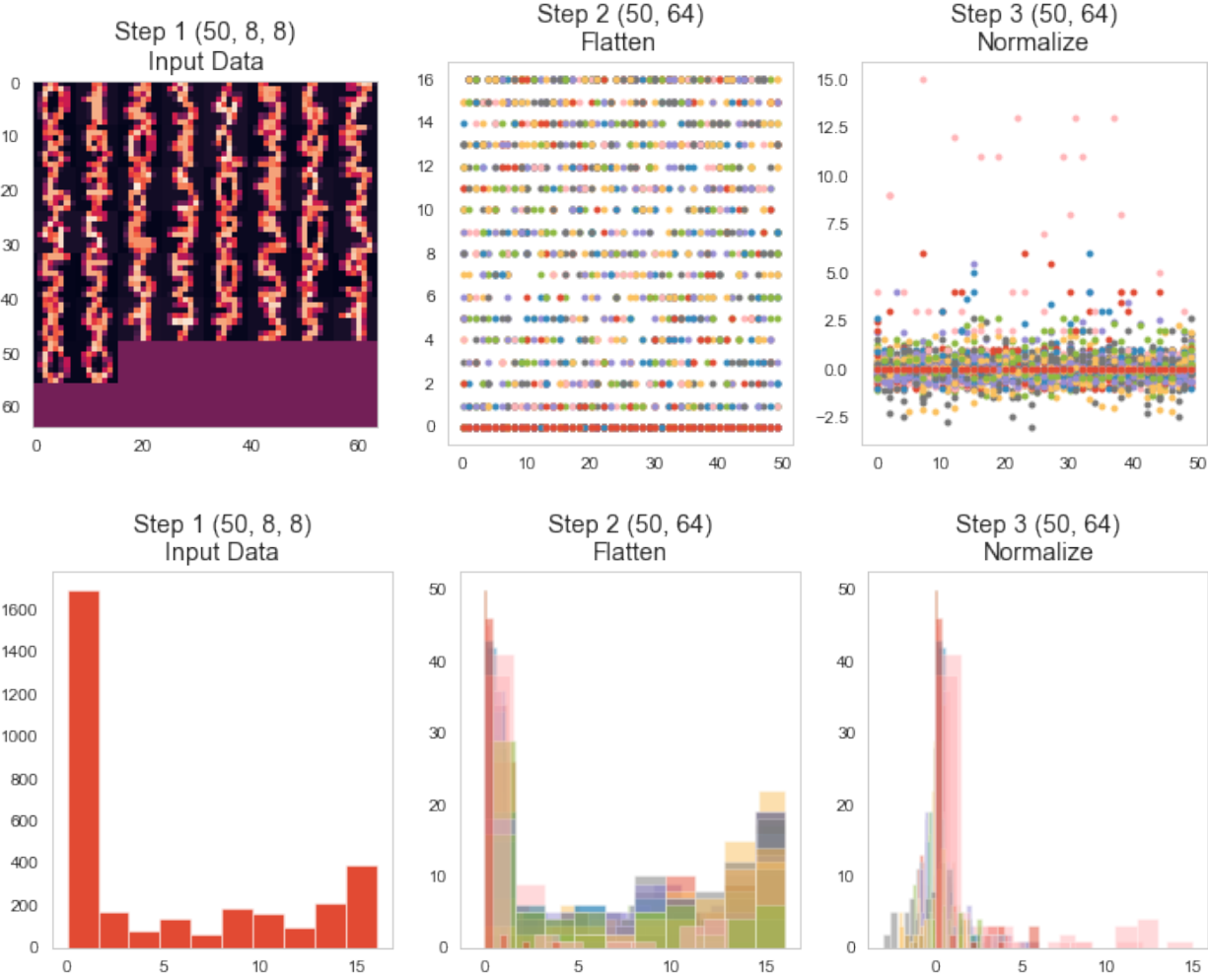
17.1.2 Run a preprocessing pipeline

- 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)
    
```



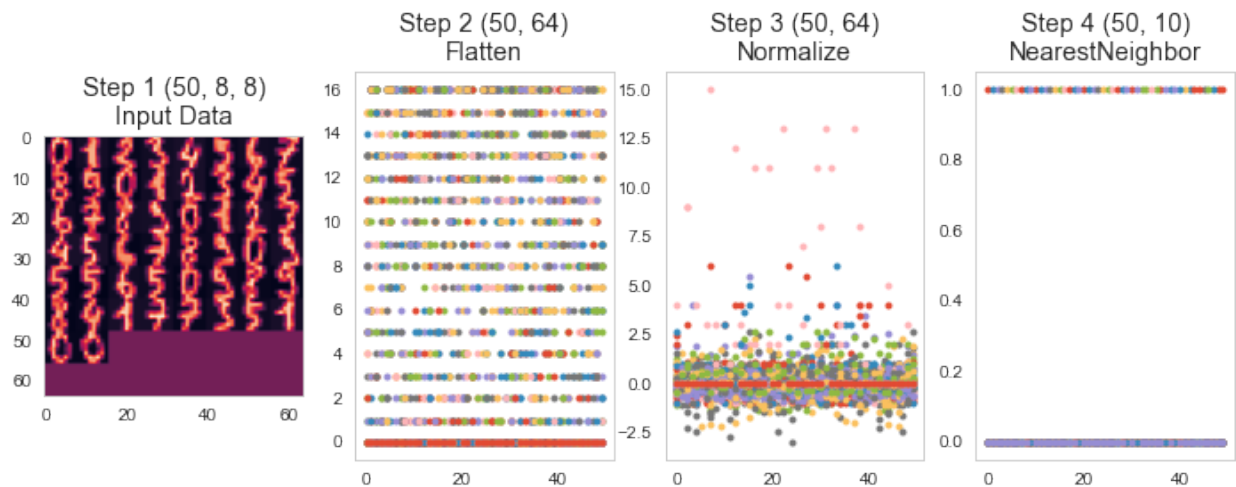
17.1.3 Add a classifier

- 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)
```



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
```

17.2.1 How about the confusion matrix?

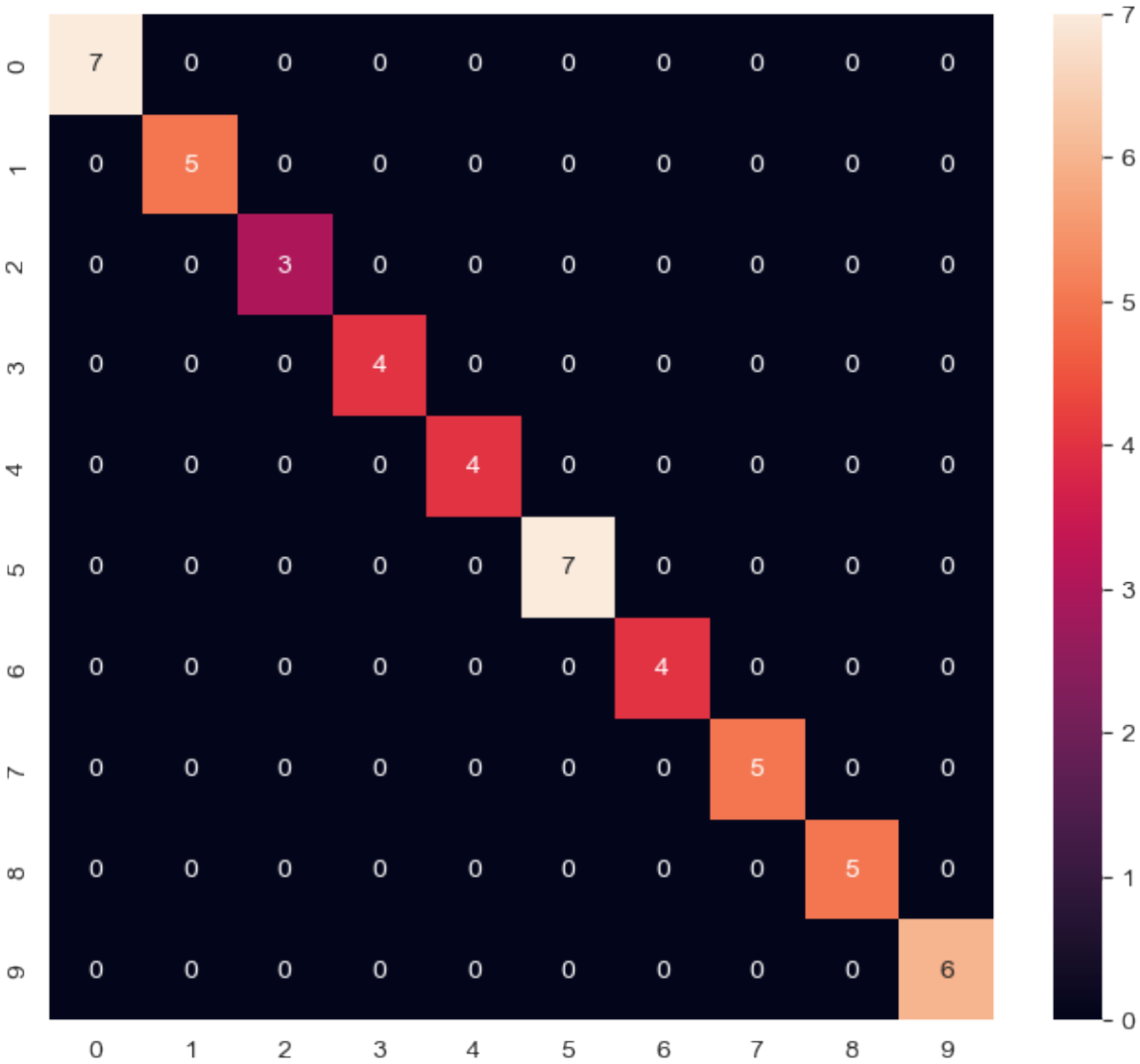
Let's look at the predictions in the confusion matrix.

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

(continues on next page)

(continued from previous page)

```
annot=True,
fmt='d',
ax=ax1);
```



17.2.2 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

(continues on next page)

(continued from previous page)

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
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

17.3 Wow! We've built an amazing algorithm!

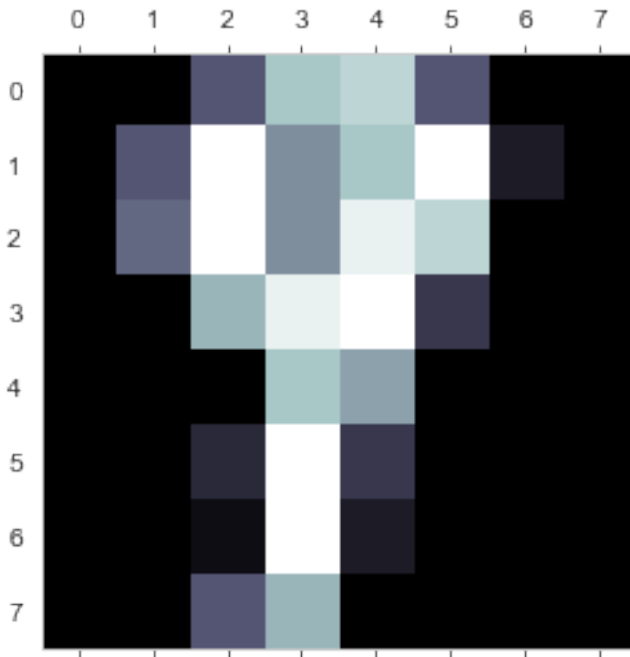
Let's patent it! Call Google!

17.3.1 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
```



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>

REGRESSION USING A PIPELINE

For regression, we can see

- it is very similar 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
```

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]
```

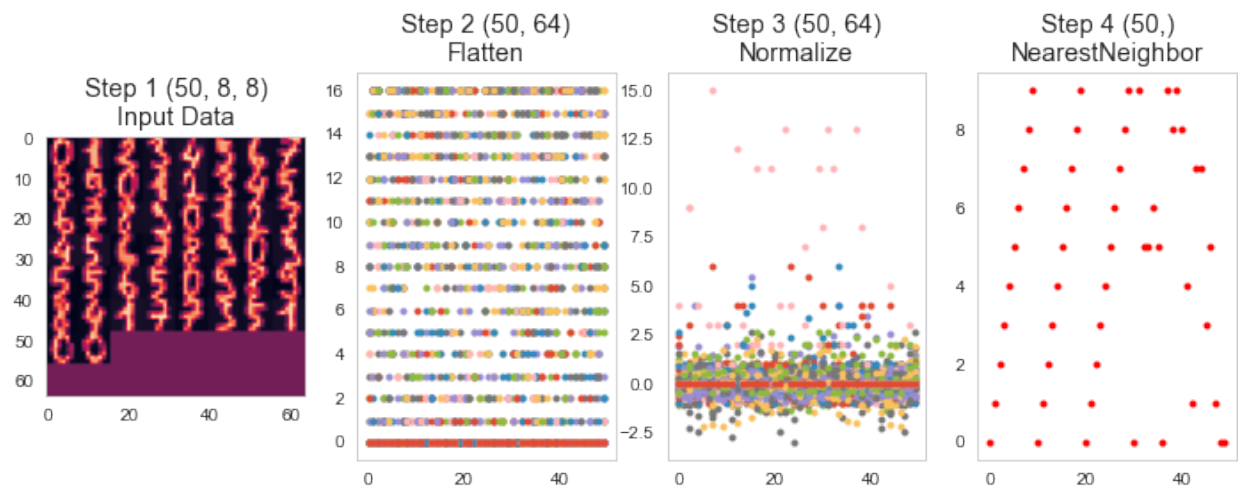
18.2 Run a KNeighbors regression

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)
```



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.

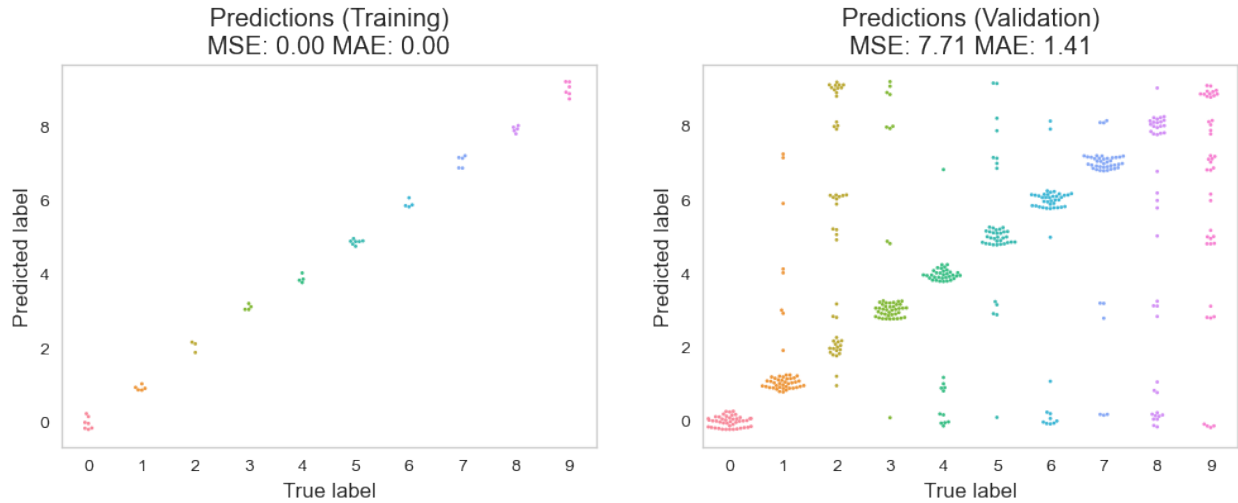
$$MSE = \frac{1}{N} \sum (y_{predicted} - y_{actual})^2$$

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

```
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.
    ↪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)
ax2.set_title('Predictions (Validation)\nmSE: %2.2f MAE: %2.2f' % (np.mean(np.
    ↪square(pred_valid-valid_id)),
                                                                    np.mean(np.
    ↪abs(pred_valid-valid_id))))
ax2.set_xlabel('True label'), ax2.set_ylabel('Predicted label');
```

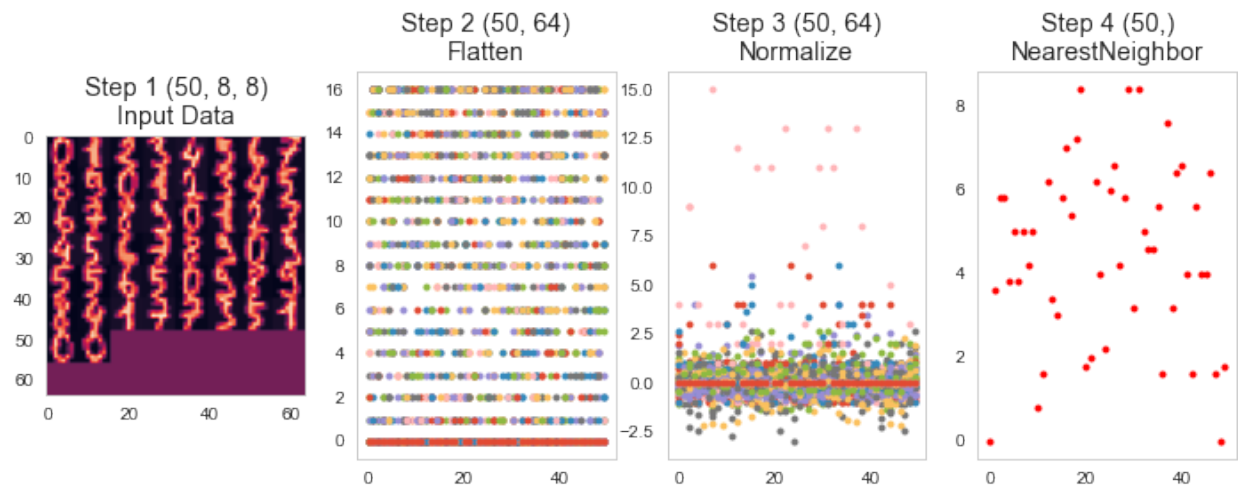
```
/Users/kaestner/miniforge3/lib/python3.9/site-packages/seaborn/categorical.py:1296:
    ↪UserWarning: 9.1% of the points cannot be placed; you may want to decrease the size
    ↪of the markers or use stripplot.
    warnings.warn(msg, UserWarning)
```



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)

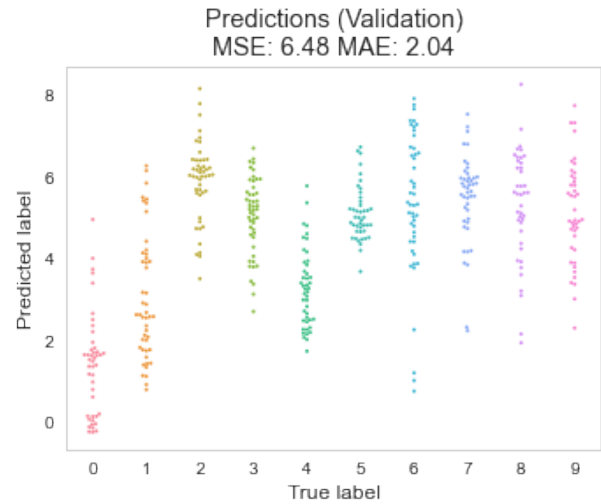
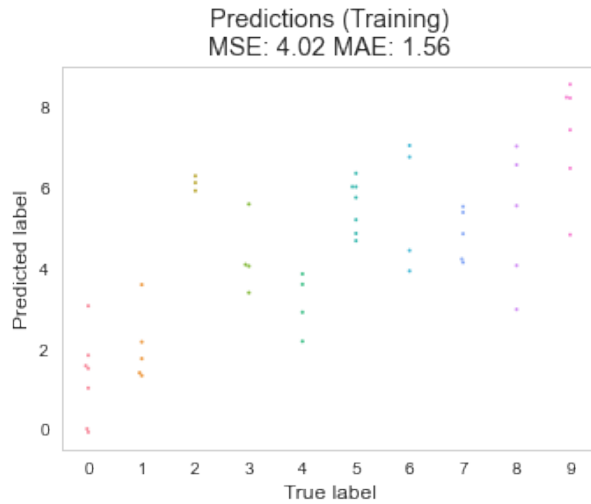
sns.swarmplot(x=digit_id, y=jitter(pred_train), ax=ax1, size=2)
ax1.set_title('Predictions (Training)\nMSE: %2.2f MAE: %2.2f' % (np.mean(np.
↳square(pred_train-digit_id)),
np.mean(np.abs(pred_
↳train-digit_id))))

ax1.set_xlabel('True label'), ax1.set_ylabel('Predicted label')
pred_valid = digit_regress_pipe.predict(valid_data)
```

(continues on next page)

(continued from previous page)

```
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.
    ↳square(pred_valid-valid_id)),
                                                    np.mean(np.
    ↳abs(pred_valid-valid_id)))));
ax2.set_xlabel('True label'), ax2.set_ylabel('Predicted label');
```



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.

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

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

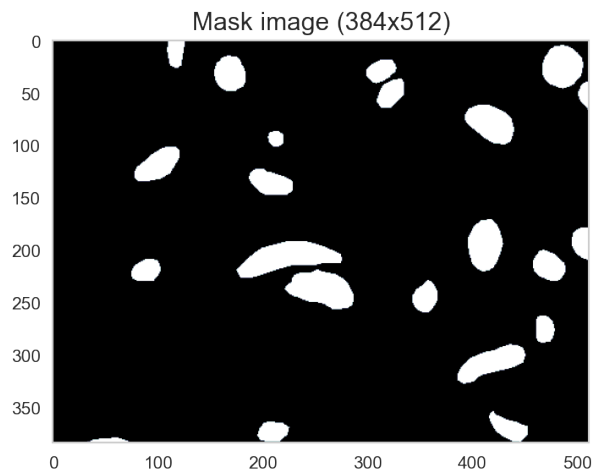
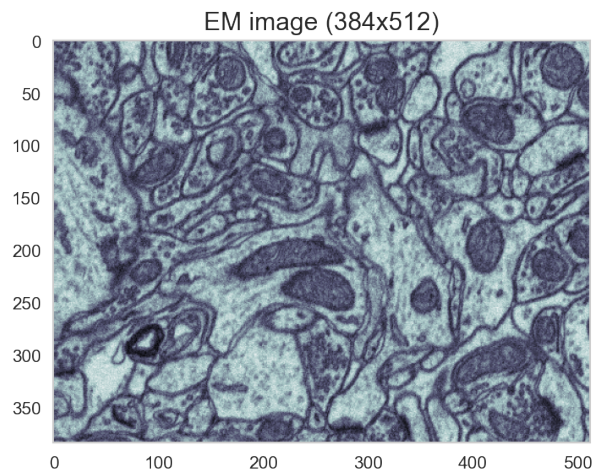
20.3 Preparing the mitochondria image and mask

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]));
```



20.3.1 Training and validation data

We only have one image available...

Solution: Split it into two parts!

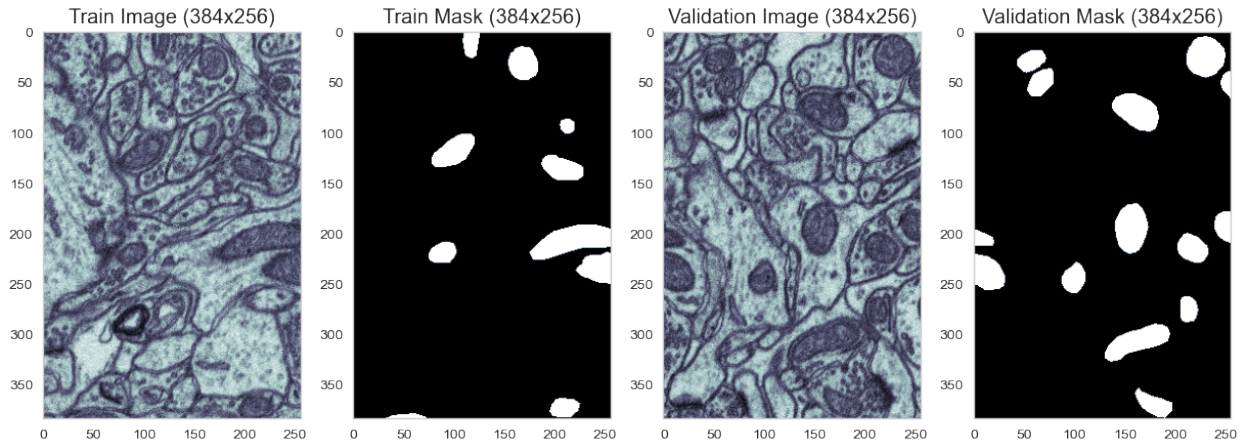
```
train_img, valid_img = cell_img[:, :256], cell_img[:, 256:]
train_mask, valid_mask = cell_seg[:, :256], cell_seg[:, 256:]
```

```
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]))
```

(continues on next page)

(continued from previous page)

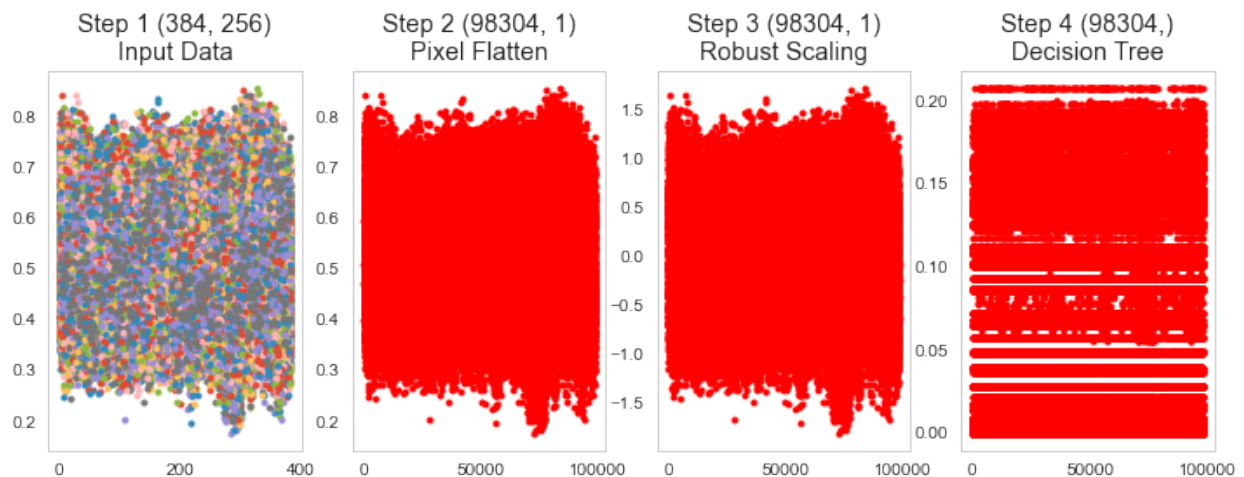
```
ax3.imshow(valid_img, cmap='bone'); ax3.set_title('Validation Image ({0}x{1})'.
    ↳format (valid_img.shape[0],valid_img.shape[1]))
ax4.imshow(valid_mask, cmap='bone'); ax4.set_title('Validation Mask ({0}x{1})'.
    ↳format (valid_mask.shape[0],valid_mask.shape[1]));
```



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]);
```



20.4.1 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')

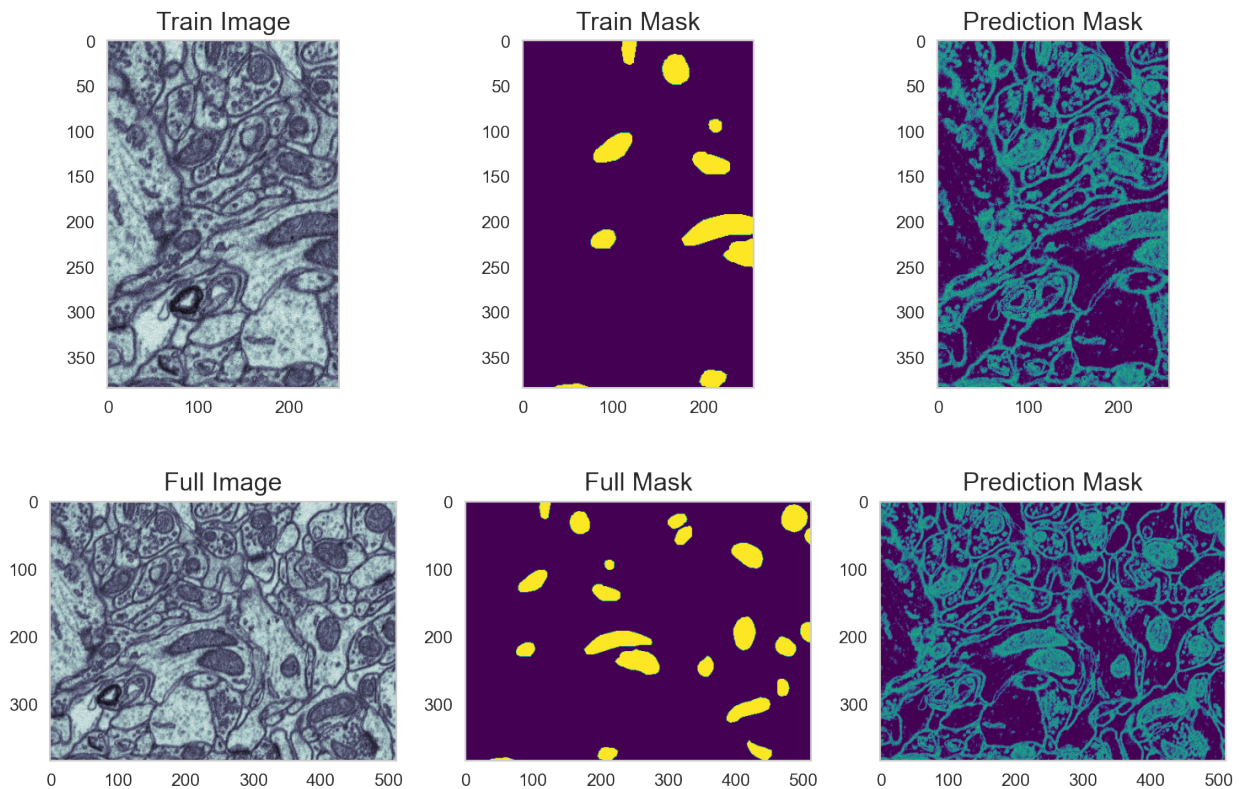
ax5.imshow(train_mask, cmap='viridis'); ax5.set_title('Train Mask')

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');
```



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(
```

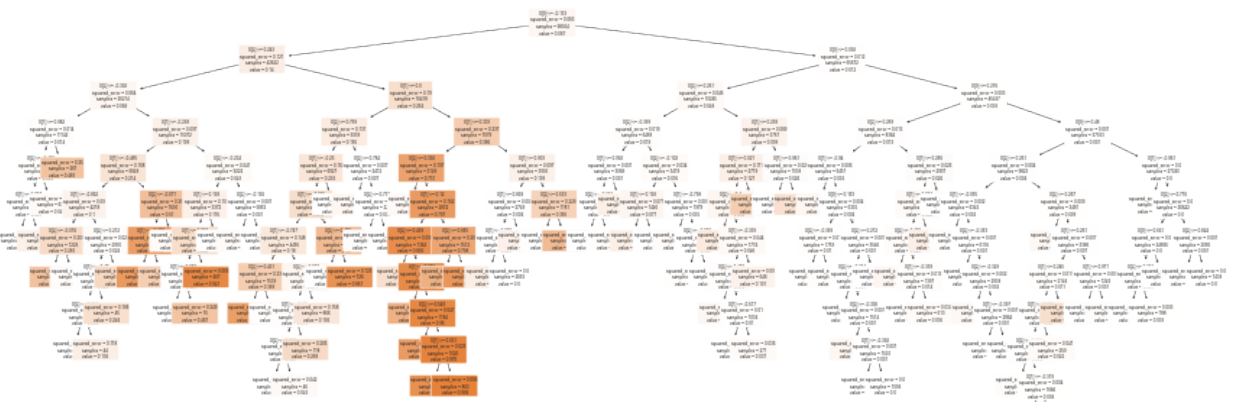
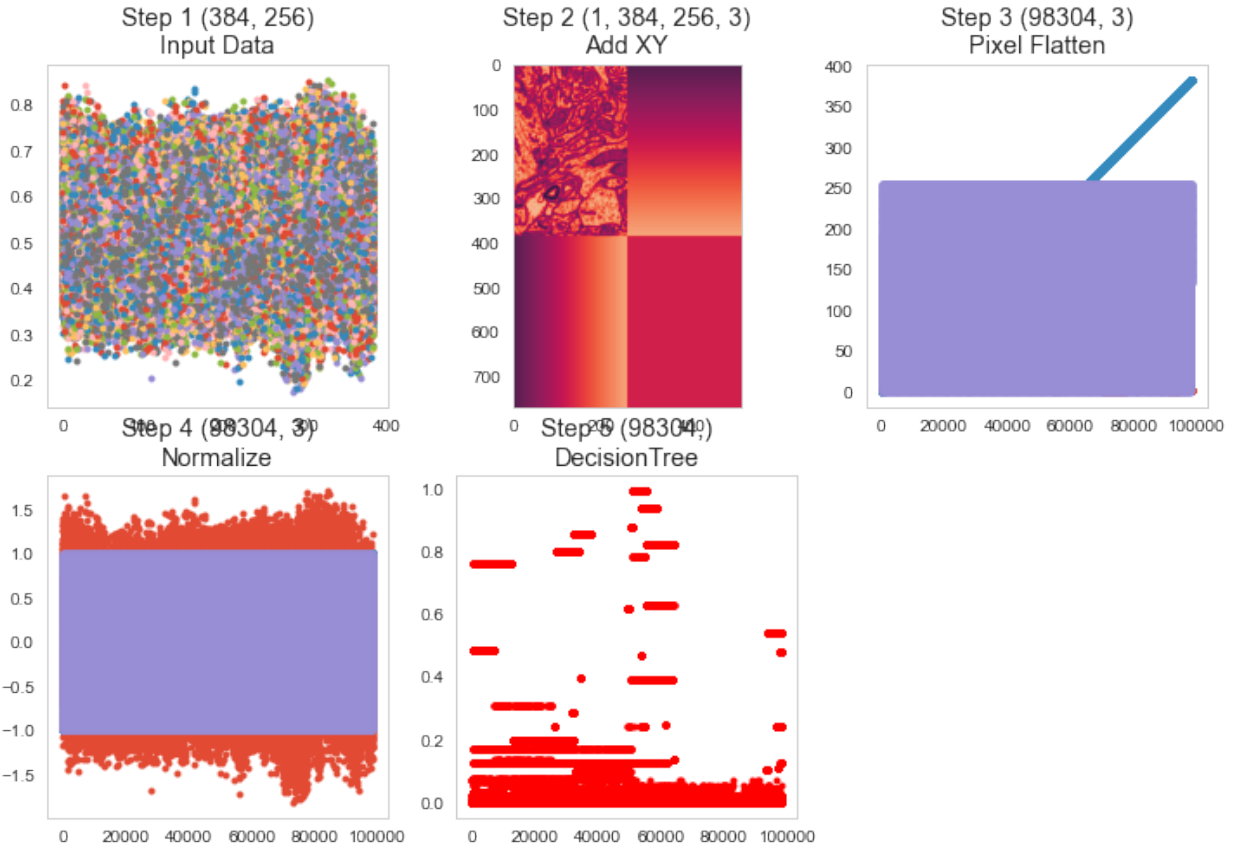
(continues on next page)

(continued from previous page)

```

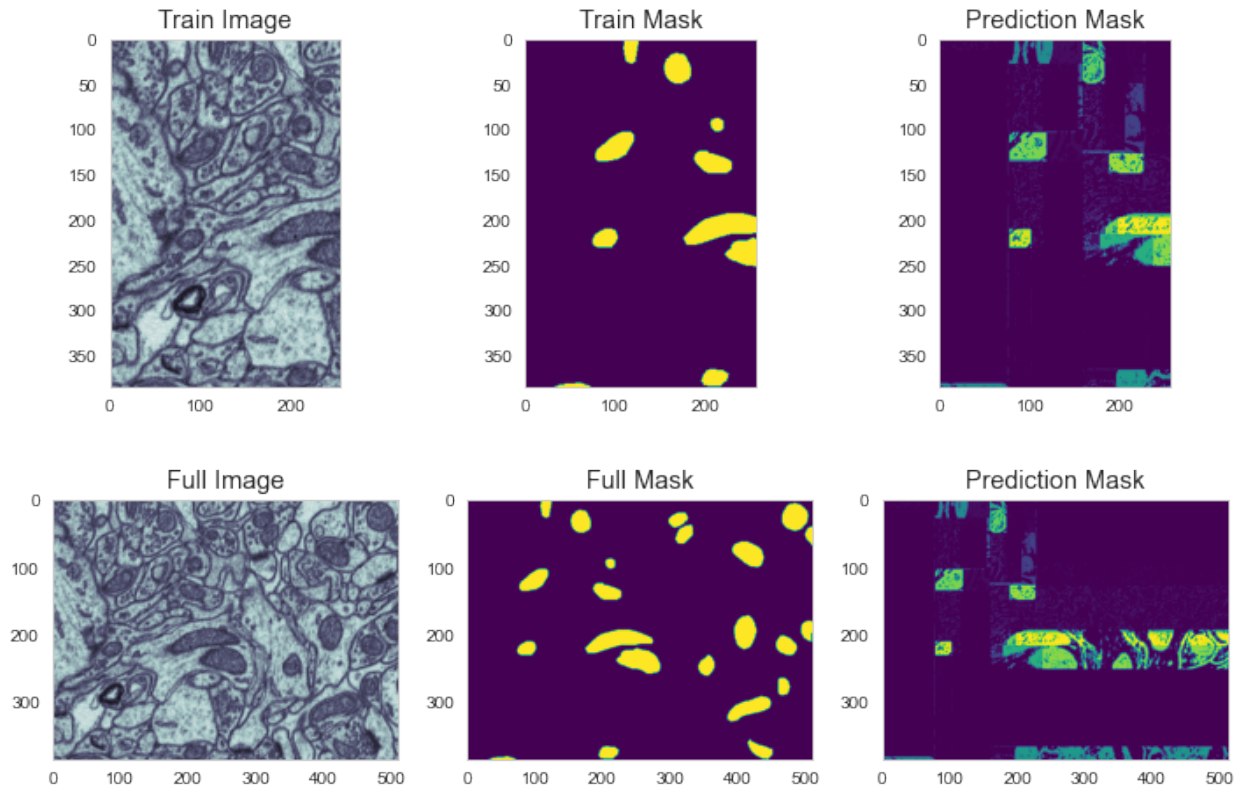
min_samples_split=1000))
])

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);
    
```



20.5.1 Did the segmentation performance improve?

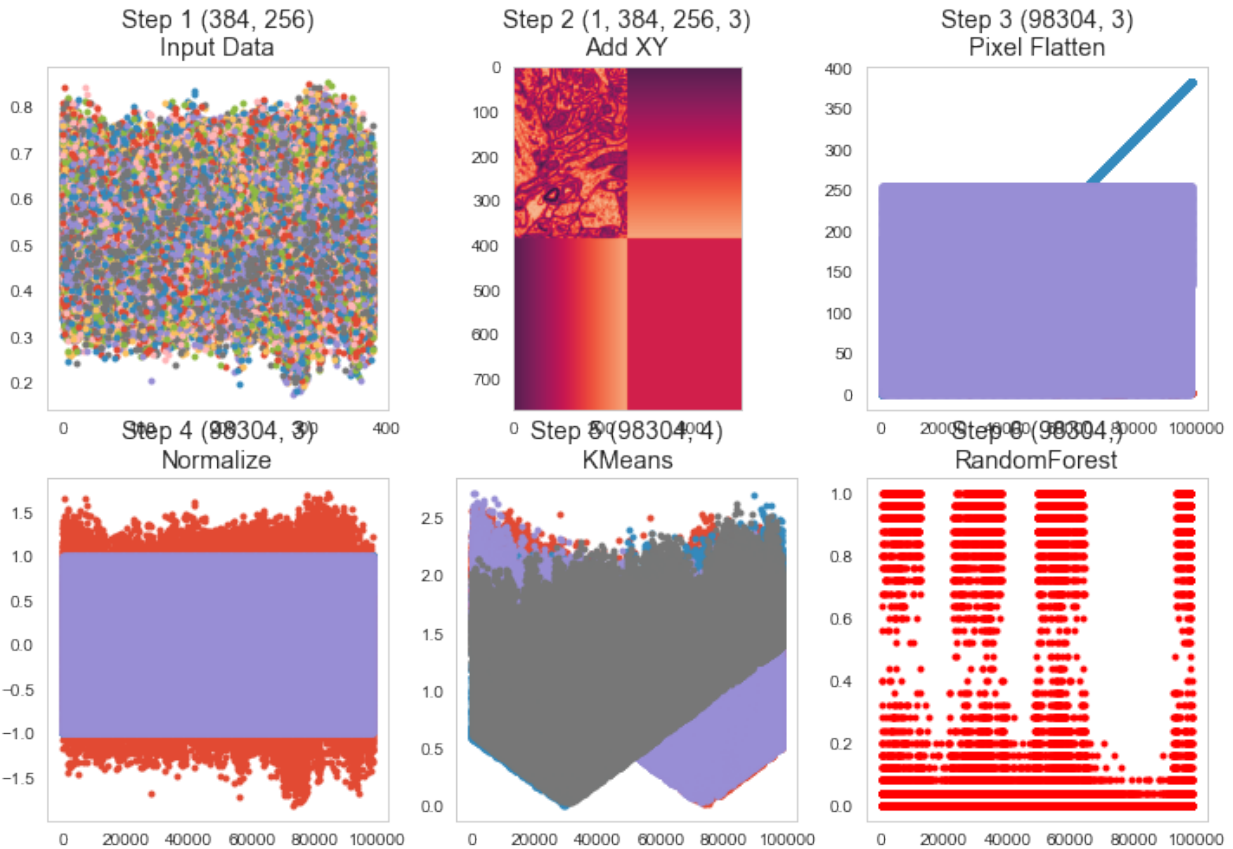
```
fig, ((ax1, ax5, ax2), (ax3, ax6, ax4)) = plt.subplots(2, 3, figsize=(12, 8), dpi=72)
ax1.imshow(train_img, cmap='bone'); ax1.set_title('Train Image')
ax5.imshow(train_mask, cmap='viridis'); ax5.set_title('Train Mask')
ax2.imshow(pred_func(train_img)[: , : , 0], cmap='viridis', vmin=0, vmax=1); ax2.set_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



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)),
                              ('RandomForest', RandomForestRegressor(n_estimators=25))
                              ])

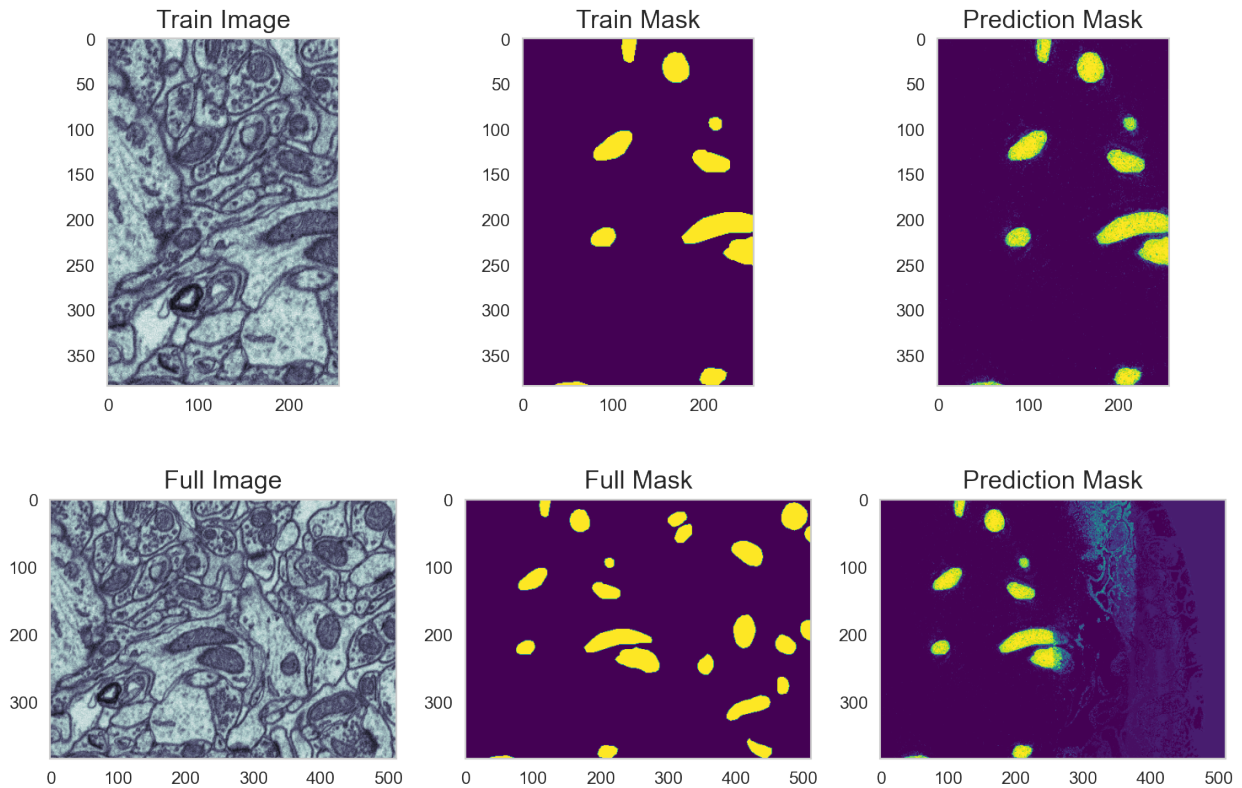
pred_func = fit_img_pipe(rf_xyseg_k_model, train_img, train_mask)
show_pipe(rf_xyseg_k_model, train_img)
```

20.6.1 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');
    
```



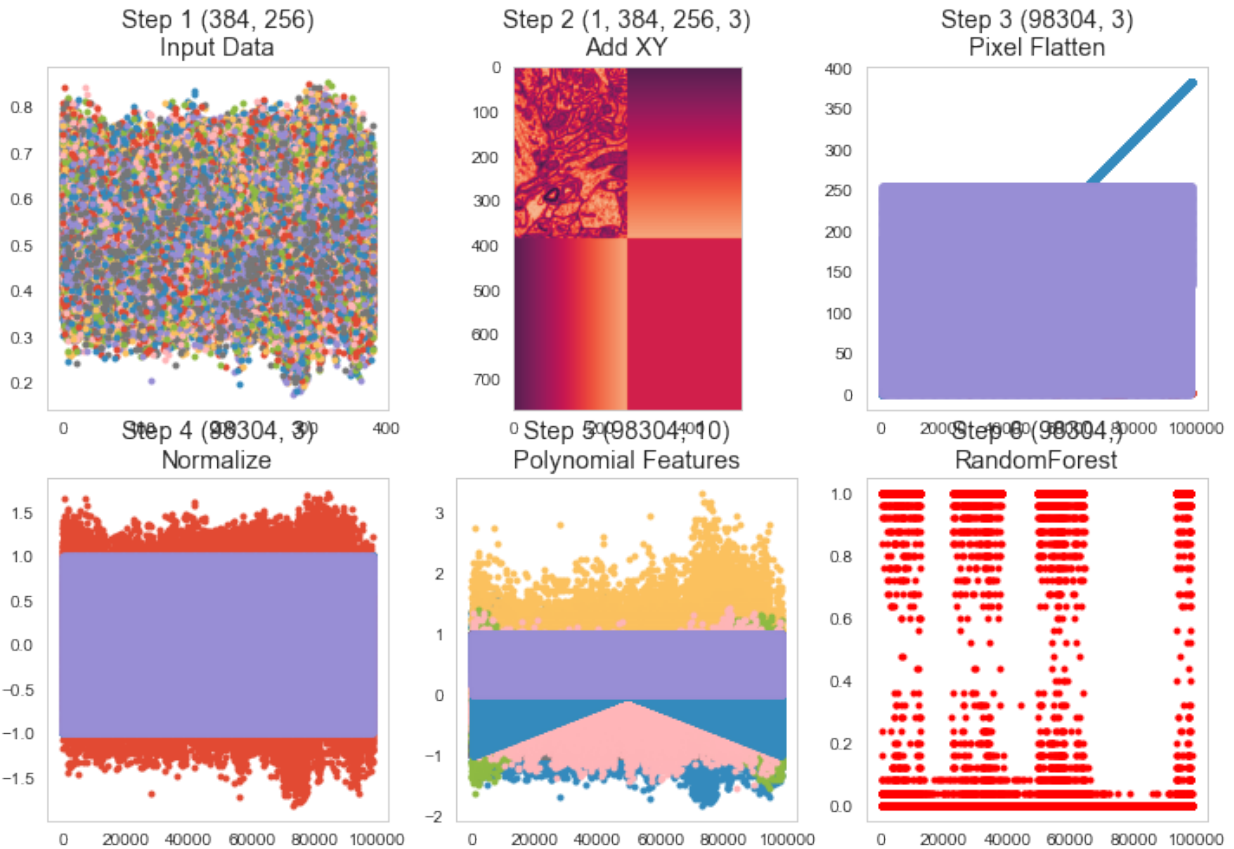
20.7 Trying polynomial features

```

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)

```



20.7.1 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')

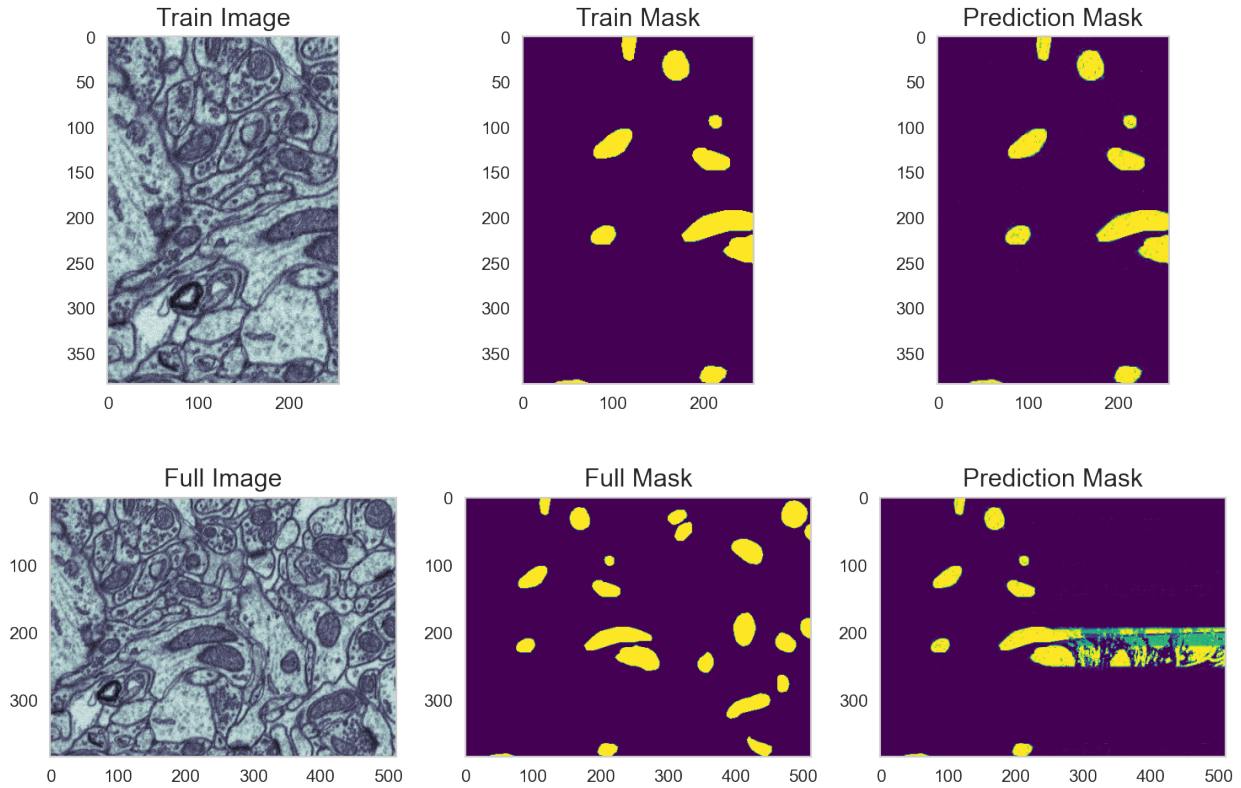
ax5.imshow(train_mask, cmap='viridis'); ax5.set_title('Train Mask')

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')

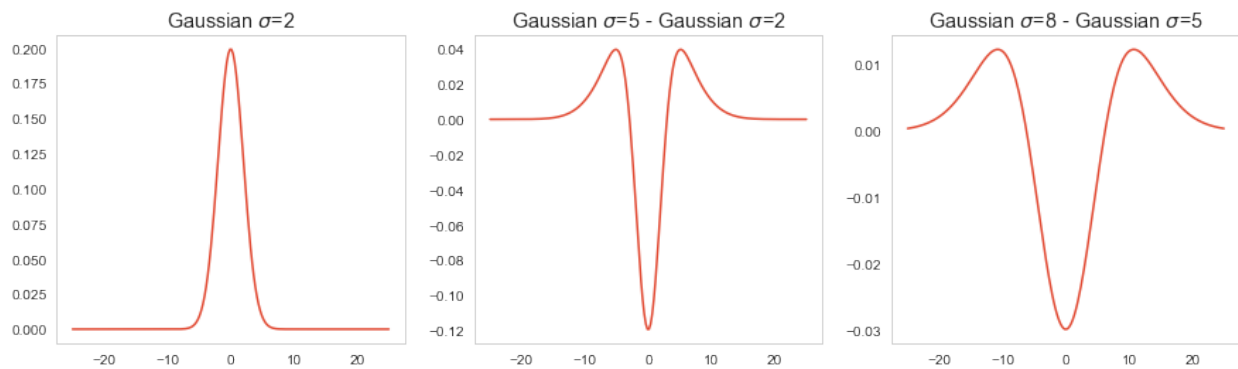
ax4.imshow(pred_func(cell_img)[: , : , 0], cmap='viridis', vmin=0, vmax=1); ax4.set_
title('Prediction Mask');
```



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  $\sigma=5$ ");
```



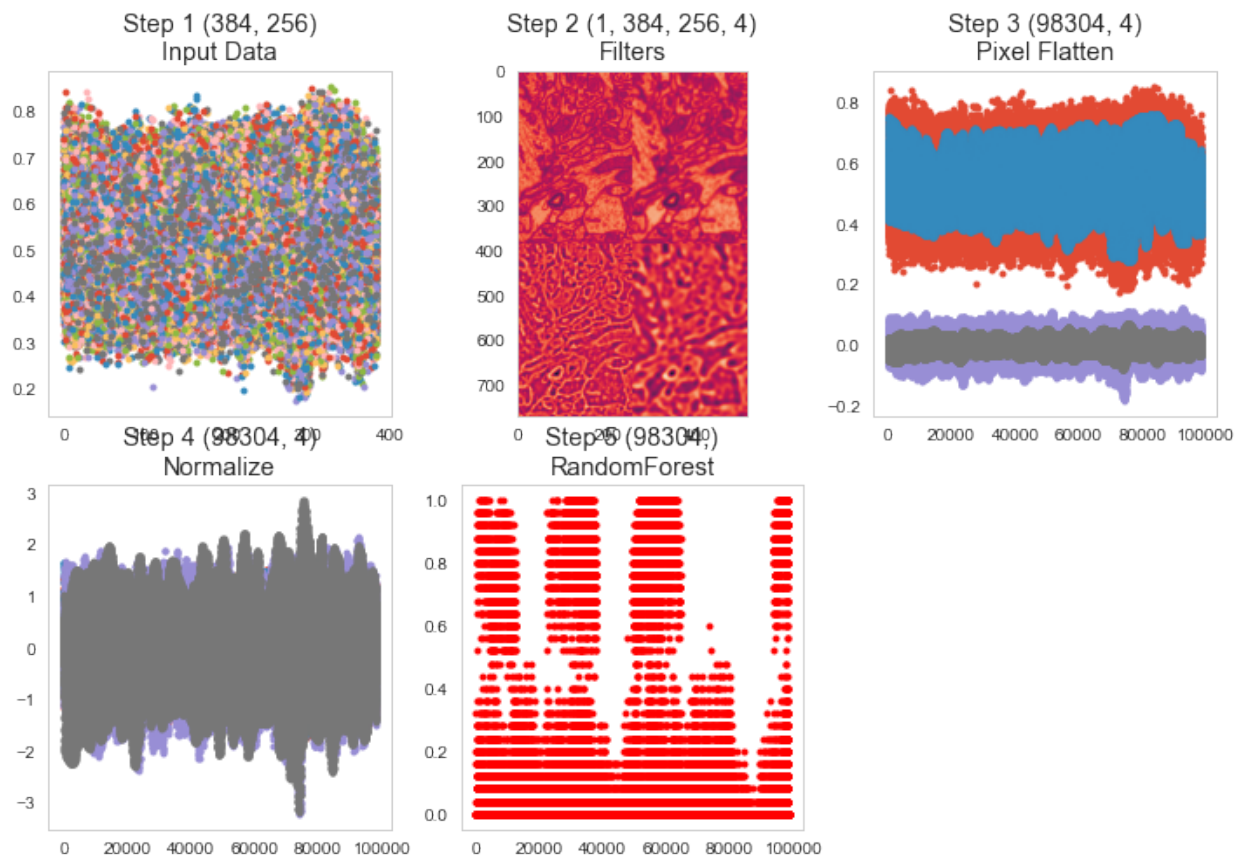
20.8.1 Pipeline with filters

```

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)

```



20.8.2 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')

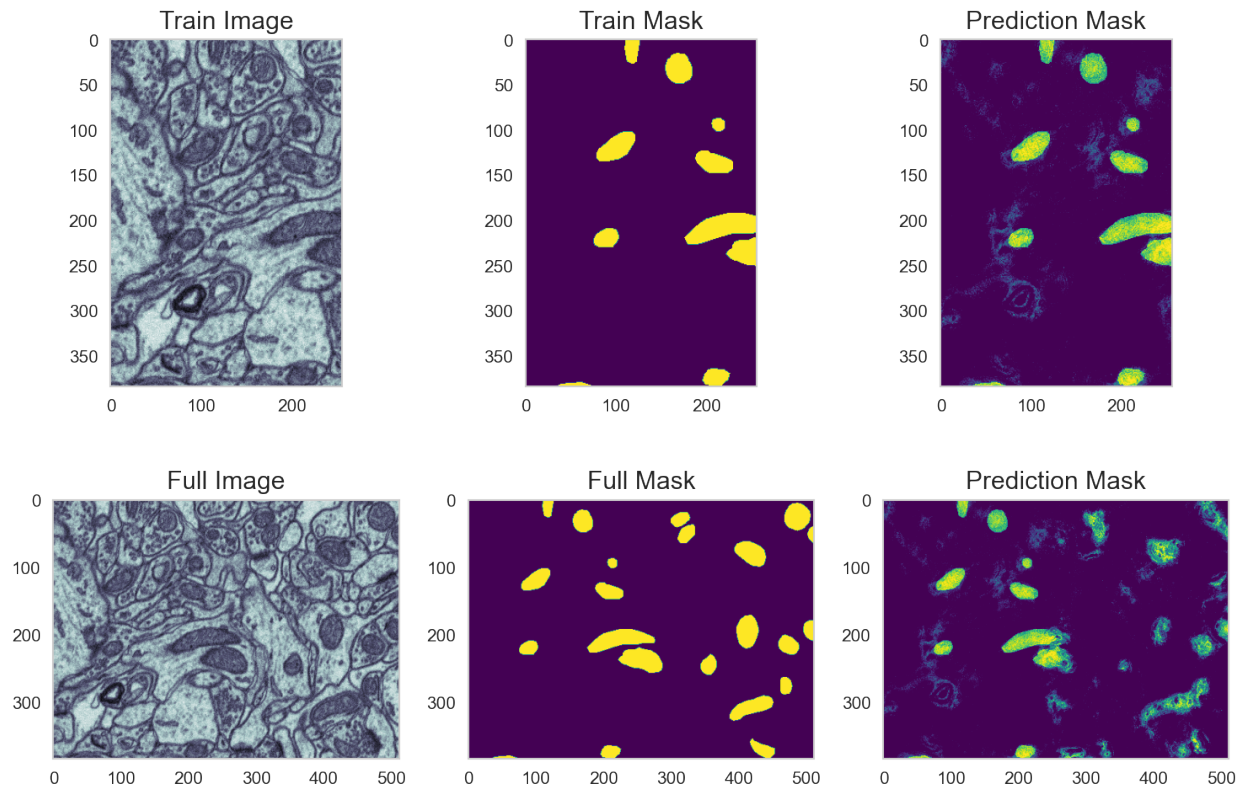
ax3.imshow(cell_img, cmap='bone'); ax3.set_title('Full Image')

```

(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');
```



20.9 Using the Neighborhood

We can also include the whole neighborhood

- shifting the image in x and y by ± 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.

20.9.1 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))

    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)

```

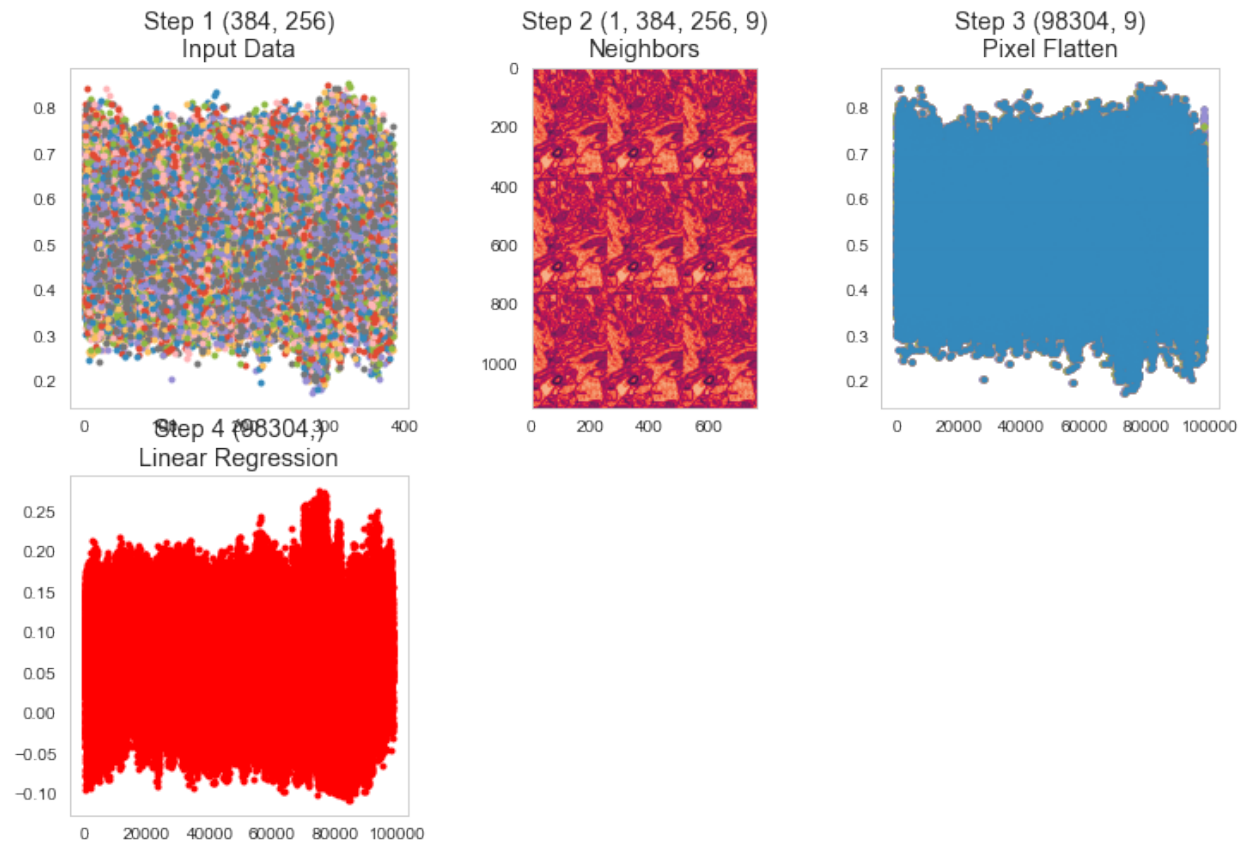
20.9.2 Pipeline with neighborhood features

```

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

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

```



20.9.3 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')

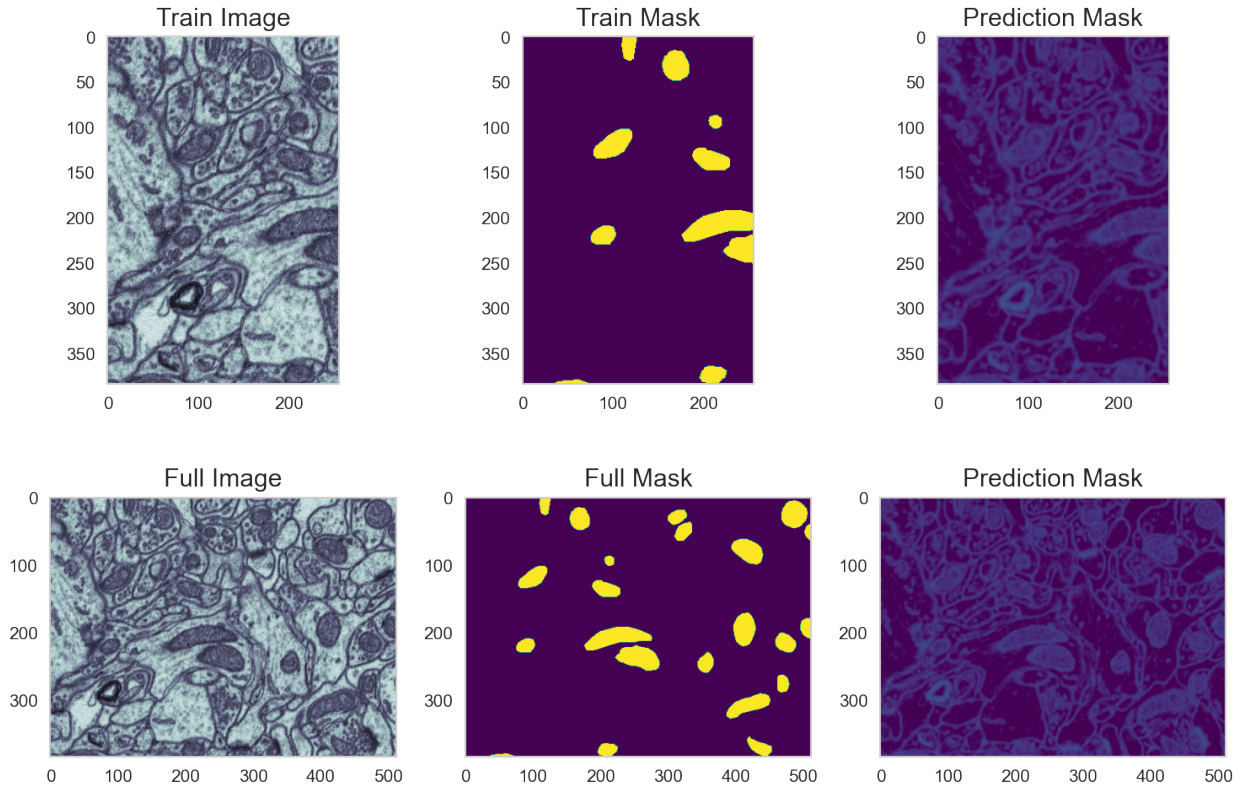
ax5.imshow(train_mask, cmap='viridis'); ax5.set_title('Train Mask')

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');
```

20.9.4 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_neighborseg_model.steps[-1][1].coef_
b = linreg_neighborseg_model.steps[-1][1].intercept_
print('M: [{:0.4}, {:0.4}, {:0.4}, {:0.4}, \033[0m, {:0.4}, {:0.4}, {:0.4},
↪4}, {:0.4}]'.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.
↪1819]
b: 0.4213
```

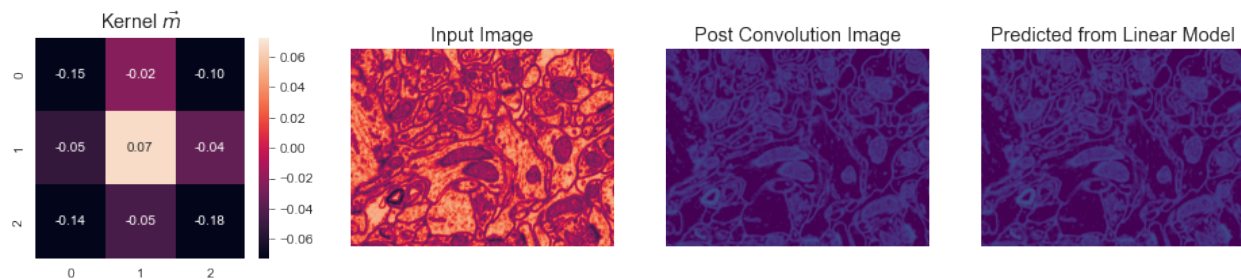
20.9.5 Convolution

The steps we have here make up a convolution and so 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');
```

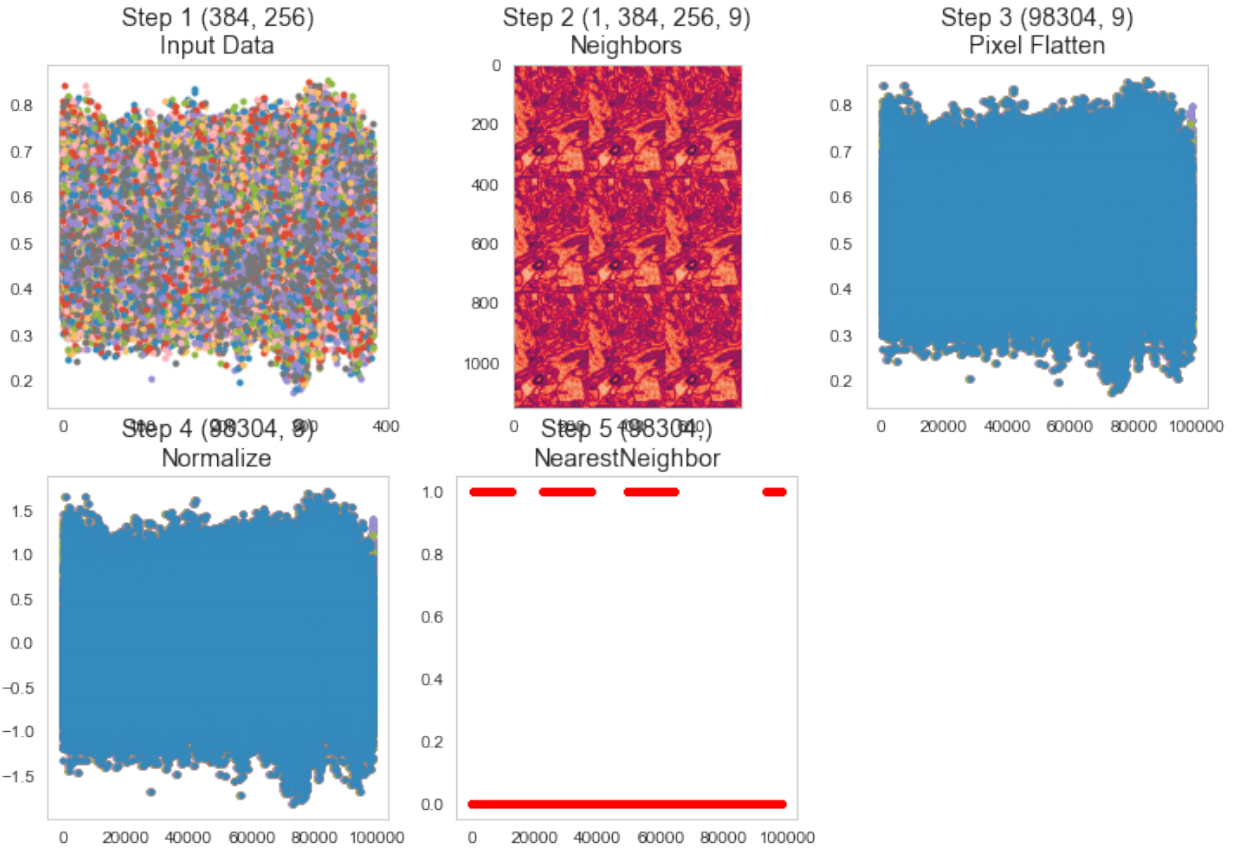


20.10 Nearest Neighbor

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_neighborseg_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_neighborseg_model, train_img, train_mask)
show_pipe(nn_neighborseg_model, train_img)
```



```

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')

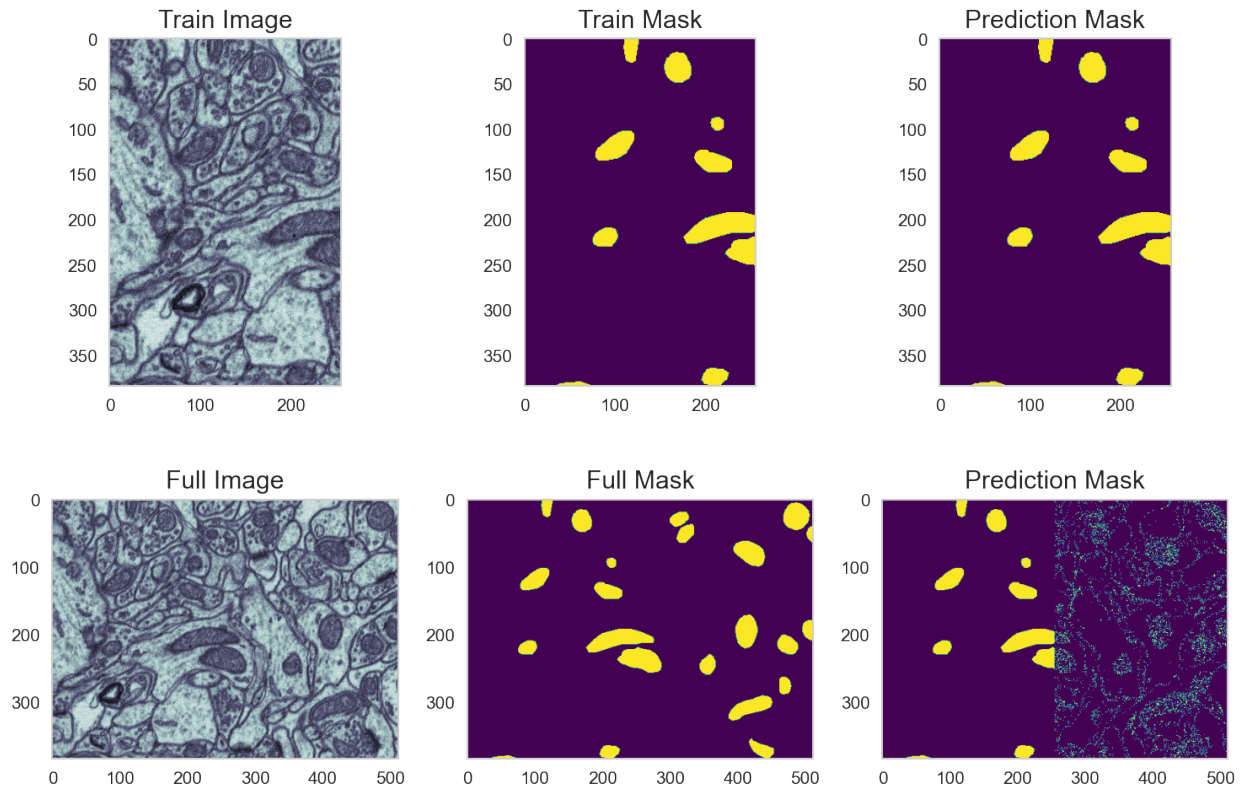
ax5.imshow(train_mask, cmap='viridis'); ax5.set_title('Train Mask')

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');
    
```



20.11 Summarizing the pipeline segmentations

- We have seen that a pipeline can be efficiently used for segmentation tasks.
- 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!

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

21.1 Let's build a small U-Net

```
import numpy as np
import skimage.io as io
import matplotlib.pyplot as plt
from tensorflow.keras.models import Model
from tensorflow.keras.layers import Input, Conv2D, MaxPool2D, UpSampling2D, concatenate
import tensorflow.keras.losses as losses
from IPython.display import SVG
from keras.utils.vis_utils import model_to_dot
base_depth = 32
in_img = Input((None, None, 1), name='Image_Input')
lay_1 = Conv2D(base_depth, kernel_size=(3, 3), padding='same', activation='relu')(in_img)
lay_2 = Conv2D(base_depth, kernel_size=(3, 3), padding='same', activation='relu')(lay_1)
lay_3 = MaxPool2D((2, 2))(lay_2)
lay_4 = Conv2D(base_depth*2, kernel_size=(3, 3), padding='same', activation='relu')(lay_3)
lay_5 = Conv2D(base_depth*2, kernel_size=(3, 3), padding='same', activation='relu')(lay_4)
lay_6 = MaxPool2D((2, 2))(lay_5)
lay_7 = Conv2D(base_depth*4, kernel_size=(3, 3), padding='same', activation='relu')(lay_6)
lay_8 = Conv2D(base_depth*4, kernel_size=(3, 3), padding='same', activation='relu')(lay_7)
lay_9 = UpSampling2D((2, 2))(lay_8)
lay_10 = concatenate([lay_5, lay_9])
lay_11 = Conv2D(base_depth*2, kernel_size=(3, 3), padding='same', activation='relu')(lay_10)
lay_12 = Conv2D(base_depth*2, kernel_size=(3, 3), padding='same', activation='relu')(lay_11)
lay_13 = UpSampling2D((2, 2))(lay_12)
lay_14 = concatenate([lay_2, lay_13])
lay_15 = Conv2D(base_depth, kernel_size=(3, 3), padding='same', activation='relu')(lay_14)
```

(continues on next page)

(continued from previous page)

```
lay_16 = Conv2D(base_depth, kernel_size=(3, 3), padding='same', activation='relu')(lay_
↳15)
lay_17 = Conv2D(1, kernel_size=(1, 1), padding='same',
activation='sigmoid')(lay_16)
t_unet = Model(inputs=[in_img], outputs=[lay_17], name='SmallUNET')
```

A network summary

```
t_unet.summary()
```

Model: "SmallUNET"

```
↳ _____
```

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

=====

Image_Input (InputLayer)	(None, None, None, 0		[]
--------------------------	----------------------	--	----

	1)]		
--	-----	--	--

--	--	--	--

conv2d (Conv2D)	(None, None, None, 320		['Image_Input[0][0]']
-----------------	------------------------	--	-----------------------

	32)		
--	-----	--	--

--	--	--	--

conv2d_1 (Conv2D)	(None, None, None, 9248		['conv2d[0][0]']
-------------------	-------------------------	--	------------------

	32)		
--	-----	--	--

--	--	--	--

max_pooling2d (MaxPooling2D)	(None, None, None, 0		['conv2d_1[0][0]']
------------------------------	----------------------	--	--------------------

	32)		
--	-----	--	--

--	--	--	--

conv2d_2 (Conv2D)	(None, None, None, 18496		['max_pooling2d[0][0]
-------------------	--------------------------	--	-----------------------

↳']			
-----	--	--	--

	64)		
--	-----	--	--

conv2d_3 (Conv2D)	(None, None, None,	36928	['conv2d_2[0][0]']	
		64)		
max_pooling2d_1 (MaxPooling2D)	(None, None, None,	0	['conv2d_3[0][0]']	
		64)		
conv2d_4 (Conv2D)	(None, None, None,	73856	['max_pooling2d_	
		128)	→1[0][0]']	
conv2d_5 (Conv2D)	(None, None, None,	147584	['conv2d_4[0][0]']	
		128)		
up_sampling2d (UpSampling2D)	(None, None, None,	0	['conv2d_5[0][0]']	
		128)		
concatenate (Concatenate)	(None, None, None,	0	['conv2d_3[0][0]',	
		192)	'up_sampling2d[0][0]	
		→']		
conv2d_6 (Conv2D)	(None, None, None,	110656	['concatenate[0][0]']	
		64)		
conv2d_7 (Conv2D)	(None, None, None,	36928	['conv2d_6[0][0]']	

Quantitative Big Imaging - Advanced segmentation

64)
up_sampling2d_1 (UpSampling2D) (None, None, None, 0 ['conv2d_7[0][0]'])
64)
concatenate_1 (Concatenate) (None, None, None, 0 ['conv2d_1[0][0]',
↔1[0][0]'] 96) 'up_sampling2d_
conv2d_8 (Conv2D) (None, None, None, 27680 ['concatenate_1[0][0]
↔'] 32)
conv2d_9 (Conv2D) (None, None, None, 9248 ['conv2d_8[0][0]'] 32)
conv2d_10 (Conv2D) (None, None, None, 33 ['conv2d_9[0][0]'] 1)
=====
Total params: 470,977
Trainable params: 470,977
Non-trainable params: 0
↔ _____

21.1.1 A schematic rendering of the network

```
dot_mod = model_to_dot(t_unet, show_shapes=True, show_layer_names=False, dpi=40)
dot_mod.set_rankdir('UD')
SVG(dot_mod.create_svg())
```

```
<IPython.core.display.SVG object>
```

21.2 New training data

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

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

(continues on next page)

(continued from previous page)

```

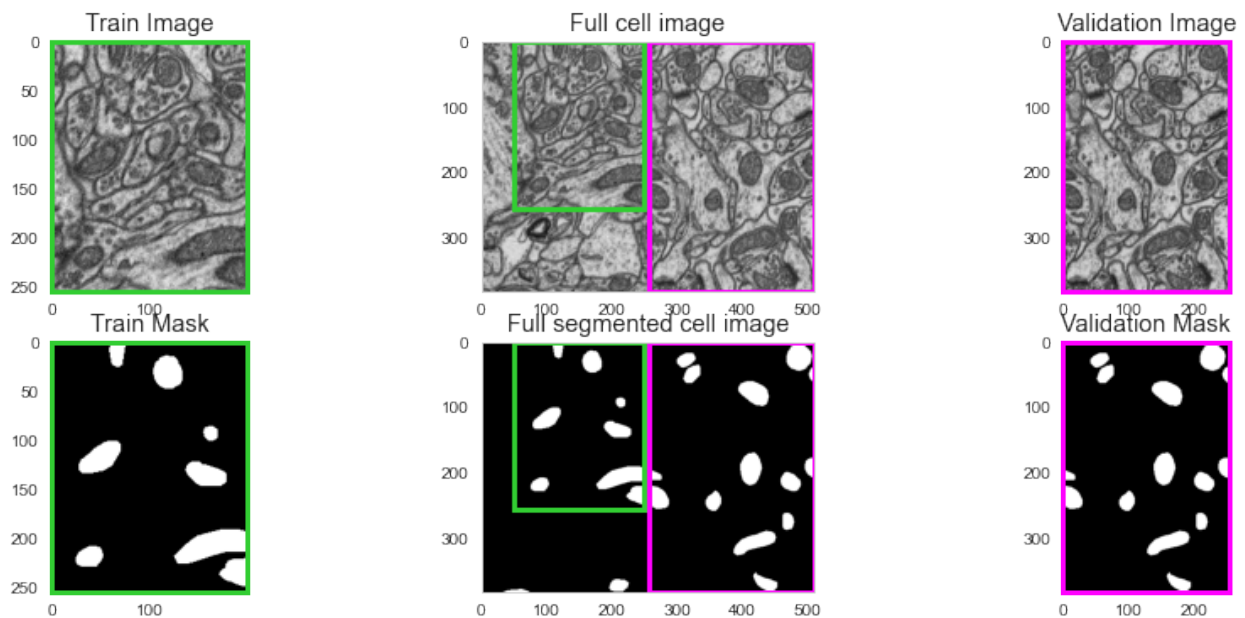
coloraxes(ax[3], 'limegreen')

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')

```



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(prepare_img(cell_img))[0, :, :, 0];
```

```

2022-03-31 18:02:15.075241: W tensorflow/core/platform/profile_utils/cpu_utils.
  →cc:128] Failed to get CPU frequency: 0 Hz

```

```

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')

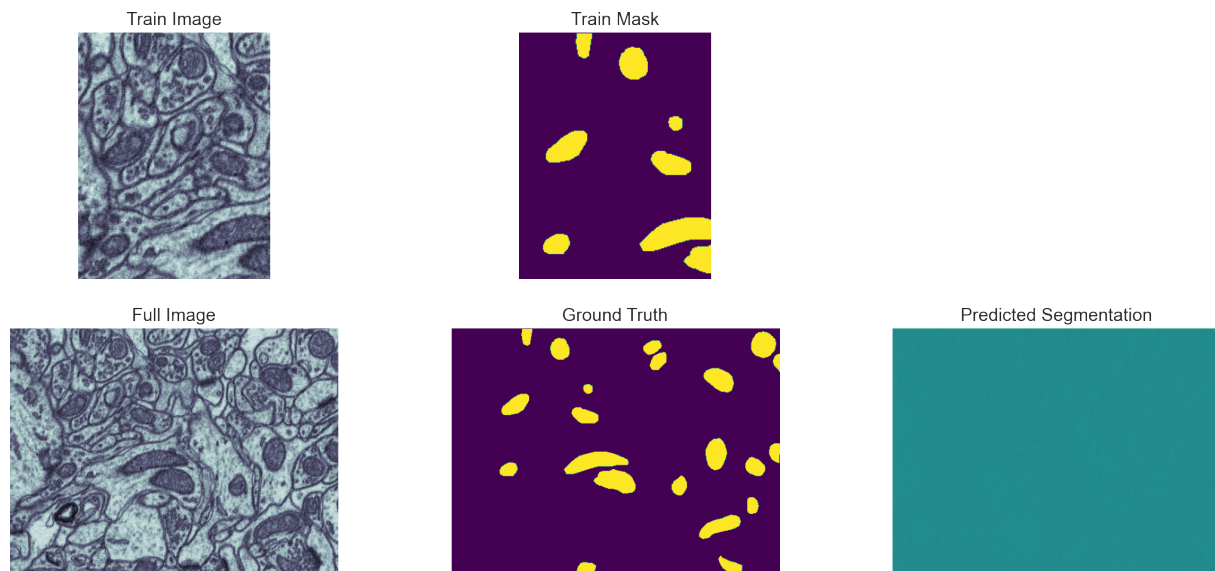
ax3.imshow(cell_img, cmap='bone')
ax3.set_title('Full Image')

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

ax5.imshow(unet_pred,
            cmap='viridis', vmin=0, vmax=1)
ax5.set_title('Predicted Segmentation')

```

```
Text(0.5, 1.0, 'Predicted Segmentation')
```



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

21.4.1 Training conditions

- Loss function - MAE
- Optimizer - Stochastic Gradient Decent
- 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 $DSC = \frac{2|X \cap Y|}{|X| + |Y|} = \frac{2TP}{2TP + FP + FN}$

21.4.2 Let's train the model

```
# from keras.optimizers import adam_v2
from keras.optimizers import sgd_experimental

#opt = adam_v2.Adam(learning_rate=0.02)
opt = sgd_experimental.SGD(learning_rate=0.05)

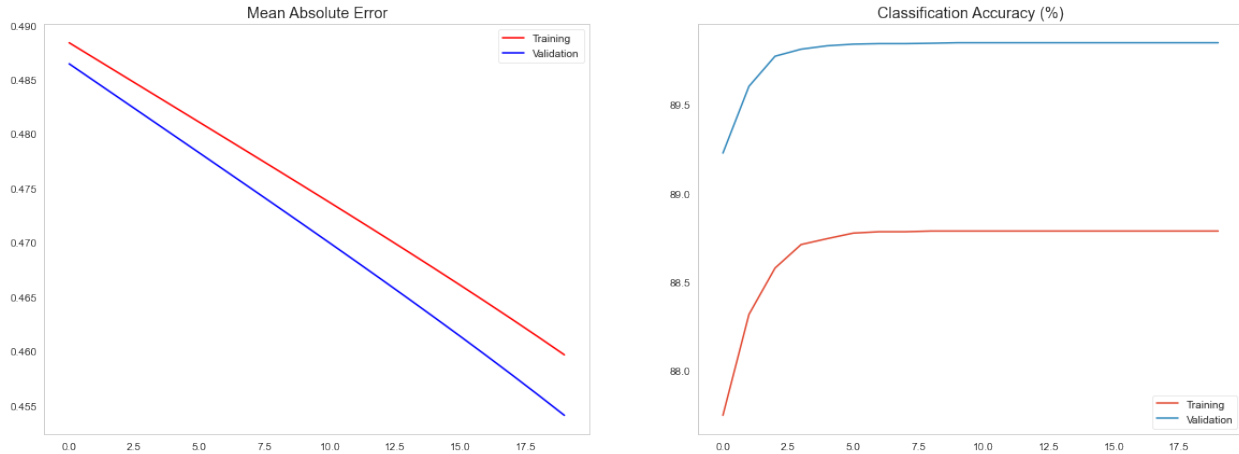
t_unet.compile(
    # we use stochastic gradient descent to optimize
    optimizer='sgd',
    # we use a simple loss metric of mean-squared error to optimize
    loss=losses.MeanSquaredError(),
    # we keep track of the number of pixels correctly classified and the mean_
    ↪absolute error as well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

loss_history = t_unet.fit(prepare_img(train_img, n=5),
                          prepare_mask(train_mask, n=5),
                          validation_data=(prepare_img(valid_img.astype(float)),
                                           prepare_mask(valid_mask.astype(float))),
                          epochs=20, verbose=0)
```

```
fig, (ax1, ax2) = plt.subplots(1, 2,
                               figsize=(20, 7))

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.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_title('Classification Accuracy (%)')
ax2.legend();
```



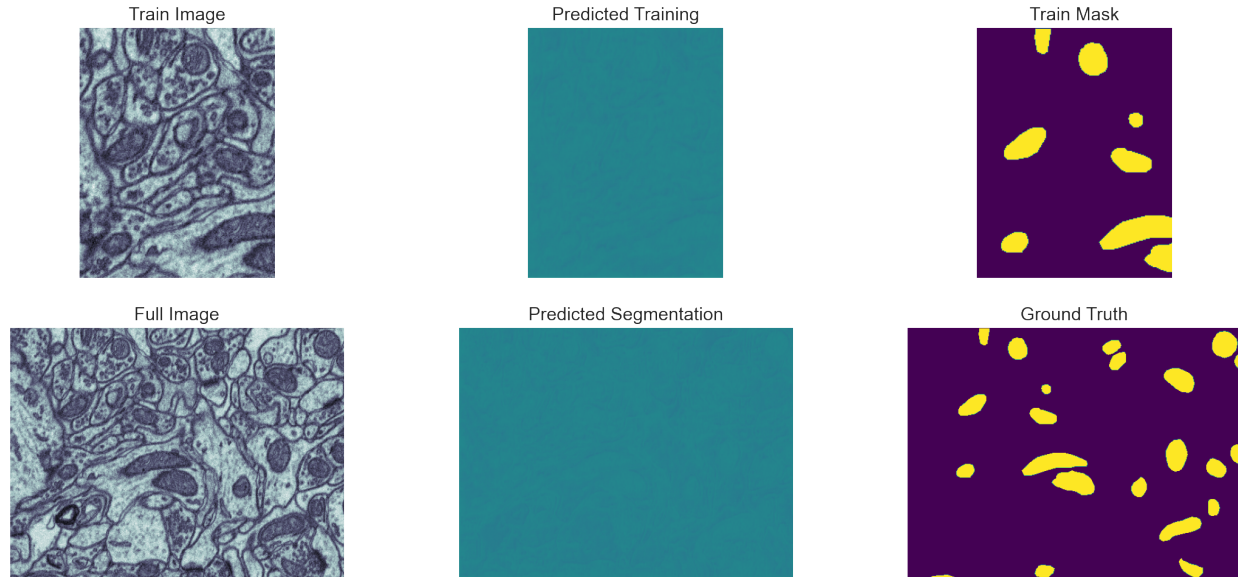
```

fig, m_axs = plt.subplots(2, 3,
                          figsize=(18, 8), dpi=150)
for c_ax in m_axs.flatten():
    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(t_unet.predict(prepare_img(train_img))[0, :, :, 0],
            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')
unet_pred = t_unet.predict(prepare_img(cell_img))[0, :, :, 0]
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');

```



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',  
    # we use stochastic gradient descent to optimize  
    optimizer=sgd_experimental.SGD(learning_rate=0.5),  
    # 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(prepare_img(train_img),  
                          prepare_mask(train_mask),  
                          validation_data=(prepare_img(valid_img),  
                                          prepare_mask(valid_mask)),  
                          epochs=20, verbose=0)
```

```
WARNING:tensorflow:The dtype of the source tensor must be floating (e.g. tf.float32)_  
↪when calling GradientTape.gradient, got tf.int64
```

```
WARNING:tensorflow:Gradients do not exist for variables ['iteration:0', 'learning_  
↪rate:0'] when minimizing the loss. If you're using `model.compile()`, did you_  
↪forget to provide a `loss` argument?
```

```
WARNING:tensorflow:Gradients do not exist for variables ['iteration:0', 'learning_
↳rate:0'] when minimizing the loss. If you're using `model.compile()`, did you
↳forget to provide a `loss` argument?
```

```
WARNING:tensorflow:The dtype of the source tensor must be floating (e.g. tf.float32)
↳when calling GradientTape.gradient, got tf.int64
```

```
WARNING:tensorflow:Gradients do not exist for variables ['iteration:0', 'learning_
↳rate:0'] when minimizing the loss. If you're using `model.compile()`, did you
↳forget to provide a `loss` argument?
```

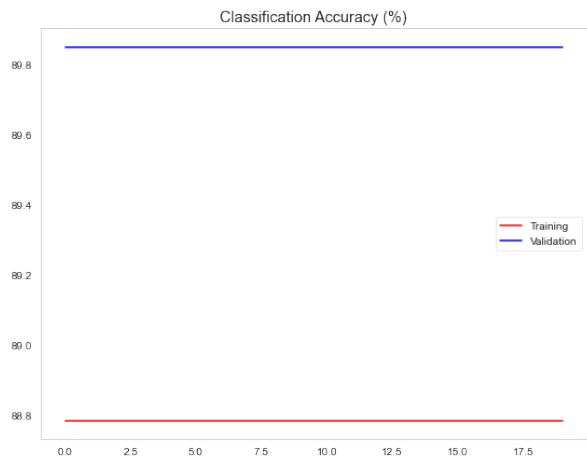
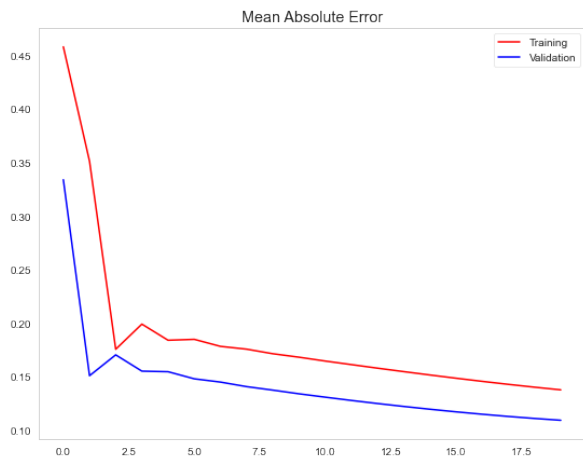
```
WARNING:tensorflow:Gradients do not exist for variables ['iteration:0', 'learning_
↳rate:0'] when minimizing the loss. If you're using `model.compile()`, did you
↳forget to provide a `loss` argument?
```

```
2022-03-31 18:02:28.079800: W tensorflow/python/util/util.cc:368] Sets are not
↳currently considered sequences, but this may change in the future, so consider
↳avoiding using them.
```

```
fig, (ax1, ax2) = plt.subplots(1, 2,
                               figsize=(20, 7))
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.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_title('Classification Accuracy (%)')
ax2.legend()
```

```
<matplotlib.legend.Legend at 0x2b861b700>
```



```

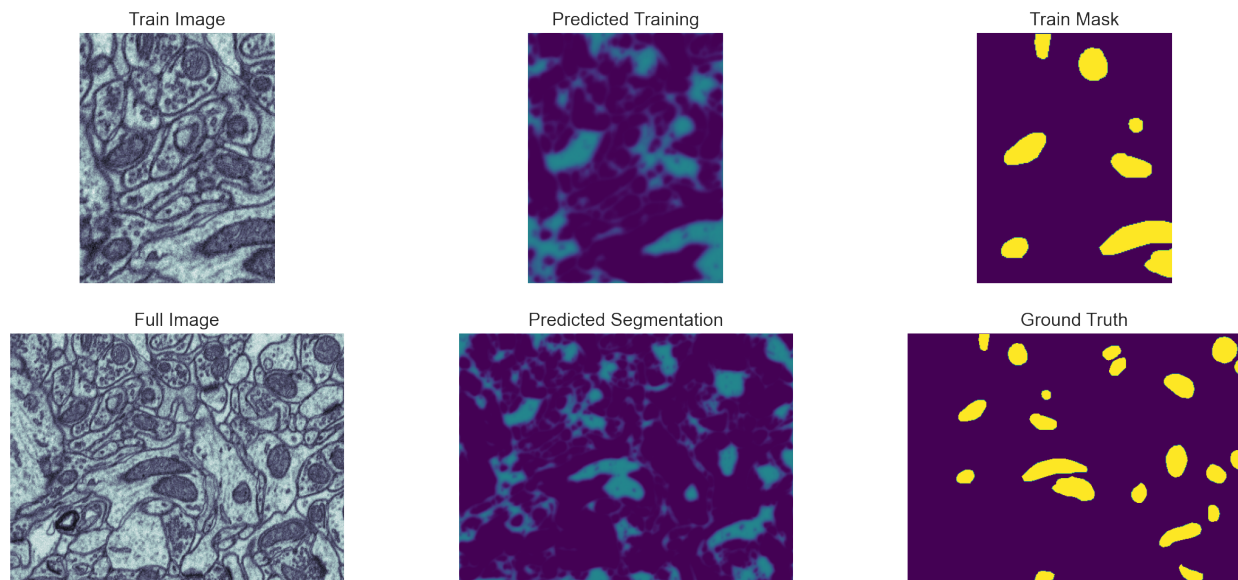
fig, m_axs = plt.subplots(2, 3,
                          figsize=(18, 8), dpi=150)
for c_ax in m_axs.flatten():
    c_ax.axis('off')
((ax1, ax15, ax2), (ax3, ax4, ax5)) = m_axs
ax1.imshow(train_img, cmap='bone')
ax1.set_title('Train Image')
ax15.imshow(t_unet.predict(prepare_img(train_img))[0, :, :, 0],
            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')
unet_pred = t_unet.predict(prepare_img(cell_img))[0, :, :, 0]
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');

```

WARNING:tensorflow:5 out of the last 5 calls to <function Model.make_predict_function.
↳<locals>.predict_function at 0x2c84065e0> triggered tf.function retracing. Tracing
↳is expensive and the excessive number of tracings could be due to (1) creating @tf.
↳function repeatedly in a loop, (2) passing tensors with different shapes, (3)
↳passing Python objects instead of tensors. For (1), please define your @tf.function
↳outside of the loop. For (2), @tf.function has experimental_relax_shapes=True
↳option that relaxes argument shapes that can avoid unnecessary retracing. For (3),
↳please refer to https://www.tensorflow.org/guide/function#controlling_retracing and
↳https://www.tensorflow.org/api_docs/python/tf/function for more details.



SUMMARY

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