

1. Business Problem

Steel is one of the most important building materials of modern times. Steel buildings are resistant to natural and man-made wear which has made the material ubiquitous around the world. The production process of flat sheet steel is especially delicate. From heating and rolling, to drying and cutting, several machines touch flat steel by the time it's ready to ship.

Today, Severstal is leading the charge in efficient steel mining and production. Severstal is now looking for machine learning to identify defects in steel which will help make production of steel more efficient. This competition will help engineers improve the defect detection algorithm by localizing and classifying surface defects on a steel sheet.

2. Source of Data

It is a Kaggle competition held by Severstal.

<https://www.severstal.com/eng/>

Data is available at <https://www.kaggle.com/c/severstal-steel-defect-detection>

3. Data Overview

- train_images/-folder with 12568 .jpg training images.
- test_images/-folder with 5516 .jpg test images (we are segmenting and classifying these images).
- train.csv -training annotations which provide segments for defects with total 4 defect classes (ClassId=[1,2,3,4]).
- sample_submission.csv -a sample submission file in the correct format (for each ImageId 4 rows, one for each of the 4 defect classes).
- Each image is of 256x1600 resolution

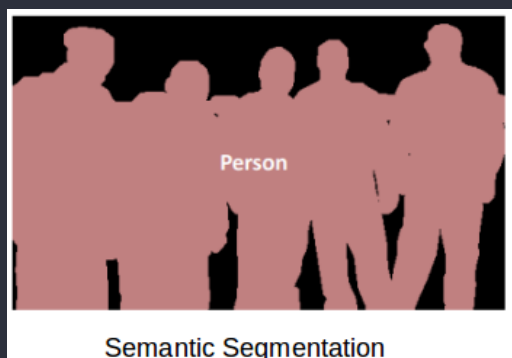
4. Mapping real world problem with Deep Learning problem

4.1. Image Segmentation

Image segmentation is the process of partitioning a digital image into multiple segments. A segmentation model returns much more detailed information about the image. More precisely, image segmentation is the process of assigning a label to every pixel in an image such that pixels with the same label share certain characteristics. Image segmentation has many applications in medical imaging, self-driving cars and satellite imaging to name a few.

The Different Types of Image Segmentation

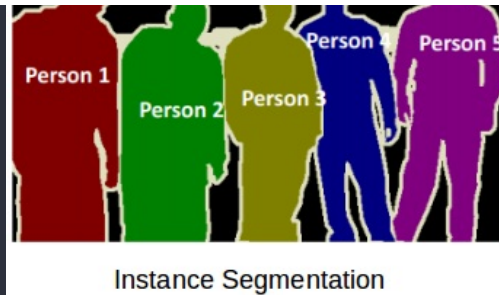
A) Semantic Segmentation



Every pixel belongs to a particular class(either background or person). Also, all the pixels belonging to a particular class are represented by the same color (background as black and person as pink). This is an example of semantic segmentation.

B) Instance Segmentation





Here also assigned a particular class to each pixel of the image. However, different objects of the same class have different colors (Person 1 as red, Person 2 as green, background as black, etc.). This is an example of instance segmentation

In Steel defect detection we will go for Semantic Segmentation.

4.2. Encoded Pixels

In order to reduce the submission file size, metric uses run-length encoding on the pixel values. Instead of submitting an exhaustive list of indices for segmentation, I will submit pairs of values that contain a start position and a run length. E.g. '1 3' implies starting at pixel 1 and running a total of 3 pixels (1,2,3). The competition format requires a space delimited list of pairs. For example, '1 3 10 5' implies pixels 1,2,3,10,11,12,13,14 are to be included in the mask. The metric checks that the pairs are sorted, positive, and the decoded pixel values are not duplicated. The pixels are numbered from top to bottom, then left to right: 1 is pixel (1,1), 2 is pixel (2,1), etc.

5. Performance Metrics

The Dice coefficient can be used to compare the pixel-wise agreement between a predicted segmentation and its corresponding ground truth.

The formula is given by:

$$2 \times \frac{A \cap B}{A \cup B}$$

where A is the predicted set of pixels and B is the ground truth. The Dice coefficient is defined to be 1 when both A and B are empty.

6. Objective

Each image may have no defects, a defect of a single class, or defects of multiple classes. For each image one must segment defects of each class (ClassId = [1,2,3,4]).

```
In [ ]:
import pandas as pd
import numpy as np
from matplotlib import pyplot as plt
import os
import matplotlib.patches as patches
import re
import random
import pickle
import cv2
import seaborn as sns
from PIL import Image
import warnings
warnings.filterwarnings("ignore")

In [ ]:
from sklearn.model_selection import train_test_split
```

```
In [ ]:
from google.colab import drive
drive.mount('/content/drive')
```

Mounted at /content/drive

7. EDA & Data Preparation

```
In [ ]:
train_csv=pd.read_csv('/content/drive//My Drive/Steel_Detection /train.csv')
train_csv.head()
```

	ImageId	ClassId	EncodedPixels
0	0002cc93b.jpg	1	29102 12 29346 24 29602 24 29858 24 30114 24 3...
1	0007a71bf.jpg	3	18661 28 18863 82 19091 110 19347 110 19603 11...
2	000a4bcdd.jpg	1	37607 3 37858 8 38108 14 38359 20 38610 25 388...
3	000f6bf48.jpg	4	131973 1 132228 4 132483 6 132738 8 132993 11 ...
4	0014fce06.jpg	3	229501 11 229741 33 229981 55 230221 77 230468...

```
In [ ]:
train_csv.shape
```

(7095, 3)

```
In [ ]:
Image_id=[]
label=[]
train_folder_path='/content/drive//My Drive/Steel_Detection /train_images'
for i in os.listdir(train_folder_path): #https://www.geeksforgeeks.org/python-os-listdir-meth
od/
    for j in range(1,5):
        Image_id.append(i)
        label.append(j)

x={'ImageId':Image_id,'ClassId':label} #https://www.geeksforgeeks.org/creating-a-pandas-data
frame/
train_img=pd.DataFrame(x)
train_img.head(10)
```

	ImageId	ClassId
0	eb5aec756.jpg	1
1	eb5aec756.jpg	2
2	eb5aec756.jpg	3
3	eb5aec756.jpg	4
4	e9b77950e.jpg	1
5	e9b77950e.jpg	2
6	e9b77950e.jpg	3
7	e9b77950e.jpg	4
8	eb7ec1f85.jpg	1
9	eb7ec1f85.jpg	2

```
In [ ]:
#https://www.geeksforgeeks.org/python-get-unique-values-list/
image_size=set()
train_folder_path='/content/drive//My Drive/Steel_Detection /train_images'
for i in os.listdir(train_folder_path): #https://www.geeksforgeeks.org/python-os-listdir-meth
od/
    image_size.add(cv2.imread(train_folder_path+'/'+i).shape)
unique_image=list(image_size)
for x in unique_image:
    print (x)
```

(355, 155, 3)

```
(256, 1600, 3)
```

- Each image is of 256x1600 resolution

```
In [ ]:
```

```
#https://stackoverflow.com/questions/53645882/pandas-merging-10
df=pd.merge(train_img,train_csv,how='outer',on=['ImageId','ClassId'])
df.fillna('',inplace=True)
df.head()
```

	ImageId	ClassId	EncodedPixels
0	eb5aec756.jpg	1	
1	eb5aec756.jpg	2	
2	eb5aec756.jpg	3	
3	eb5aec756.jpg	4	
4	e9b77950e.jpg	1	378485 4 378733 13 378985 18 379241 18 379496 ...

```
In [ ]:
```

```
#https://www.analyticsvidhya.com/blog/2020/03/pivot-table-pandas-python/
train=pd.pivot_table(df,values='EncodedPixels',index='ImageId',columns='ClassId',aggfunc=np.sum).astype(str)
train=train.reset_index()
train.columns=['image_id','rle_1','rle_2','rle_3','rle_4']
train.head()
```

	image_id	rle_1	rle_2	rle_3	rle_4
0	0002cc93b.jpg	29102 12 29346 24 29602 24 29858 24 30114 24 3...			
1	00031f466.jpg				
2	000418bfc.jpg				
3	000789191.jpg				
4	0007a71bf.jpg		18661 28 18863 82 19091 110 19347 110 19603 11...		

```
In [ ]:
```

```
#Stratified corrosion is a type of corrosion that progresses parallel to the metal surface in
such a manner that underlying layers are gradually separated.
#For stratified sampling, we have taken stratified based on minority label priority
#https://economictimes.indiatimes.com/definition/stratified-sampling
defect=[]
stratify=[]
for i in range(len(train)):
    if (train['rle_1'][i] != '' or train['rle_2'][i] != '' or train['rle_3'][i] != '' or train['rle_4'][i] != ''):
        defect.append(1)
    else:
        defect.append(0)

    if train['rle_1'][i] != '':
        stratify.append(1)
    elif train['rle_2'][i] != '':
        stratify.append(2)
    elif train['rle_3'][i] != '':
        stratify.append(3)
    elif train['rle_4'][i] != '':
        stratify.append(4)
    else:
        stratify.append(0)
train['defect']=defect
train['stratify']=stratify
```

```
In [ ]:
```

```
train.head()
```

	image_id	rle_1	rle_2	rle_3	rle_4	defect	stratify
--	----------	-------	-------	-------	-------	--------	----------

	image_id	29102	12	29346	24	29602	24	29858	24	30114	24	30370	24	30626	24	30882	24	31138	24	31394	24	31650	24	31906	24	32162	24	32418	24	32674	24	32930	24	33186	24	33442	24	33698	24	33954	24	34210	24	34466	24	34722	24	34978	24	35234	24	35490	24	35746	24	36002	24	36258	24	36514	24	36770	24	37026	24	37282	24	37538	24	37794	24	38050	24	38306	24	38562	24	38818	24	39074	24	39330	24	39586	24	39842	24	40098	24	40354	24	40610	24	40866	24	41122	24	41378	24	41634	24	41890	24	42146	24	42402	24	42658	24	42914	24	43170	24	43426	24	43682	24	43938	24	44194	24	44450	24	44706	24	44962	24	45218	24	45474	24	45730	24	45986	24	46242	24	46498	24	46754	24	47010	24	47266	24	47522	24	47778	24	48034	24	48290	24	48546	24	48802	24	49058	24	49314	24	49570	24	49826	24	50082	24	50338	24	50594	24	50850	24	51106	24	51362	24	51618	24	51874	24	52130	24	52386	24	52642	24	52898	24	53154	24	53410	24	53666	24	53922	24	54178	24	54434	24	54690	24	54946	24	55202	24	55458	24	55714	24	55970	24	56226	24	56482	24	56738	24	56994	24	57250	24	57506	24	57762	24	58018	24	58274	24	58530	24	58786	24	59042	24	59298	24	59554	24	59810	24	60066	24	60322	24	60578	24	60834	24	61090	24	61346	24	61602	24	61858	24	62114	24	62370	24	62626	24	62882	24	63138	24	63394	24	63650	24	63906	24	64162	24	64418	24	64674	24	64930	24	65186	24	65442	24	65698	24	65954	24	66210	24	66466	24	66722	24	66978	24	67234	24	67490	24	67746	24	68002	24	68258	24	68514	24	68770	24	69026	24	69282	24	69538	24	69794	24	70050	24	70306	24	70562	24	70818	24	71074	24	71330	24	71586	24	71842	24	72098	24	72354	24	72610	24	72866	24	73122	24	73378	24	73634	24	73890	24	74146	24	74402	24	74658	24	74914	24	75170	24	75426	24	75682	24	75938	24	76194	24	76450	24	76706	24	76962	24	77218	24	77474	24	77730	24	77986	24	78242	24	78498	24	78754	24	79010	24	79266	24	79522	24	79778	24	80034	24	80290	24	80546	24	80802	24	81058	24	81314	24	81570	24	81826	24	82082	24	82338	24	82594	24	82850	24	83106	24	83362	24	83618	24	83874	24	84130	24	84386	24	84642	24	84898	24	85154	24	85410	24	85666	24	85922	24	86178	24	86434	24	86690	24	86946	24	87202	24	87458	24	87714	24	87970	24	88226	24	88482	24	88738	24	88994	24	89250	24	89506	24	89762	24	90018	24	90274	24
0	0001	29102	12	29346	24	29602	24	29858	24	30114	24	30370	24	30626	24	30882	24	31138	24	31394	24	31650	24	31906	24	32162	24	32418	24	32674	24	32930	24	33186	24	33442	24	33698	24	33954	24	34210	24	34466	24	34722	24	34978	24	35234	24	35490	24	35746	24	36002	24	36258	24	36514	24	36770	24	37026	24	37282	24	37538	24	37794	24	38050	24	38306	24	38562	24	38818	24	39074	24	39330	24	39586	24	39842	24	40098	24	40354	24	40610	24	40866	24	41122	24	41378	24	41634	24	41890	24	42146	24	42402	24	42658	24	42914	24	43170	24	43426	24	43682	24	43938	24	44194	24	44450	24	44706	24	44962	24	45218	24	45474	24	45730	24	45986	24	46242	24	46498	24	46754	24	47010	24	47266	24	47522	24	47778	24	48034	24	48290	24	48546	24	48802	24	49058	24	49314	24	49570	24	49826	24	50082	24	50338	24	50594	24	50850	24	51106	24	51362	24	51618	24	51874	24	52130	24	52386	24	52642	24	52898	24	53154	24	53410	24	53666	24	53922	24	54178	24	54434	24	54690	24	54946	24	55202	24	55458	24	55714	24	55970	24	56226	24	56482	24	56738	24	56994	24	57250	24	57506	24	57762	24	58018	24	58274	24	58530	24	58786	24	59042	24	59298	24	59554	24	59810	24	60066	24	60322	24	60578	24	60834	24	61090	24	61346	24	61602	24	61858	24	62114	24	62370	24	62626	24	62882	24	63138	24	63394	24	63650	24	63906	24	64162	24	64418	24	64674	24	64930	24	65186	24	65442	24	65698	24	65954	24	66210	24	66466	24	66722	24	66978	24	67234	24	67490	24	67746	24	68002	24	68258	24	68514	24	68770	24	69026	24	69282	24	69538	24	69794	24	70050	24	70306	24	70562	24	70818	24	71074	24	71330	24	71586	24	71842	24	72098	24	72354	24	72610	24	72866	24	73122	24	73378	24	73634	24	73890	24	74146	24	74402	24	74658	24	74914	24	75170	24	75426	24	75682	24	75938	24	76194	24	76450	24	76706	24	76962	24	77218	24	77474	24	77730	24	77986	24	78242	24	78498	24	78754	24	79010	24	79266	24	79522	24	79778	24	80034	24	80290	24	80546	24	80802	24	81058	24	81314	24	81570	24	81826	24	82082	24	82338	24	82594	24	82850	24	83106	24	83362	24	83618	24	83874	24	84130	24	84386	24	84642	24	84898	24	85154	24	85410	24	85666	24	85922	24	86178	24	86434	24	86690	24	86946	24	87202	24	87458	24	87714	24	87970	24	88226	24	88482	24	88738	24	88994	24	89250	24	89506	24	89762	24	90018	24	90274	24
1	00031f466.jpg																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																

	image_id	rle_1	rle_2	rle_3	rle_4	defect	stratify	defect_1	defect_2	defect_3	defect_4	total_defects
0	0002cc93b.jpg	29102 12 29346 24 29602 24 29858 24 30114 24 3...				1	1	1	0	0	0	1
1	00031f466.jpg					0	0	0	0	0	0	0
2	000418bfc.jpg					0	0	0	0	0	0	0
3	000789191.jpg					0	0	0	0	0	0	0
4	0007a71bf.jpg		18661 28 18863 82 19091 110 19347 110 19603 11...			1	3	0	0	1	0	1

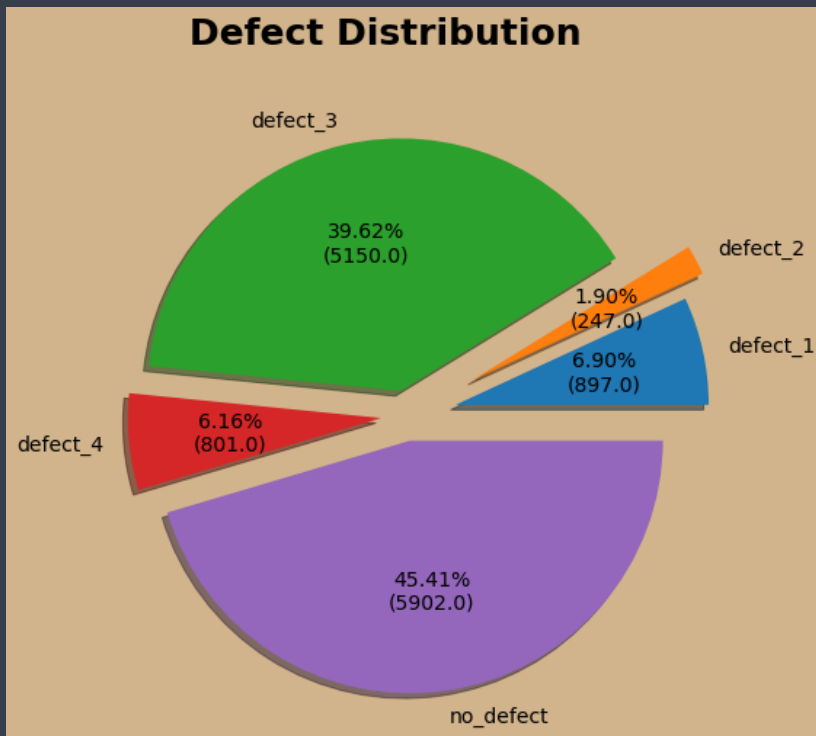
```
In [ ]: train.shape
```

```
In [ ]: len(test_image)
```

```

a = p * 100 / 100
return "{:.2f}%\n({:.1f})".format(p,a)
for i in range(len(train)):
    if train['rle_1'][i] != '':
        defect_1+=1
    if train['rle_2'][i] != '':
        defect_2+=1
    if train['rle_3'][i] != '':
        defect_3+=1
    if train['rle_4'][i] != '':
        defect_4+=1
    if train['defect'][i] == 0:
        no_defect+=1
labels=['defect_1','defect_2','defect_3','defect_4','no_defect']
sizes=[defect_1,defect_2,defect_3,defect_4,no_defect]
explode=(0.2,0.3,0.1,0.1,0.1)
fig,ax=plt.subplots(figsize=(16,8))
ax.pie(sizes,explode=explode,labels=labels,textprops={'fontsize': 14},autopct=lambda p: func(
sizes,p),shadow=True)
fig.suptitle('Defect Distribution',fontsize=25,fontweight='bold')
fig.set_facecolor("tan")
plt.show()

```



- The dataset is very imbalanced.
- Data augmentation and resampling techniques will be required to perform the defect detection.

```

In [ ]:
def patch1(bar,ax):
    #https://stackoverflow.com/questions/52080991/display-percentage-above-bar-chart-in-matplot
    lib
    for p in bar.patches:
        width=p.get_width()
        height=p.get_height()
        x,y=p.get_xy()
        ax.annotate('{:.1f}%'.format(height/100),(x+width/2,y+height*1.02),ha='center',fontSize=14)

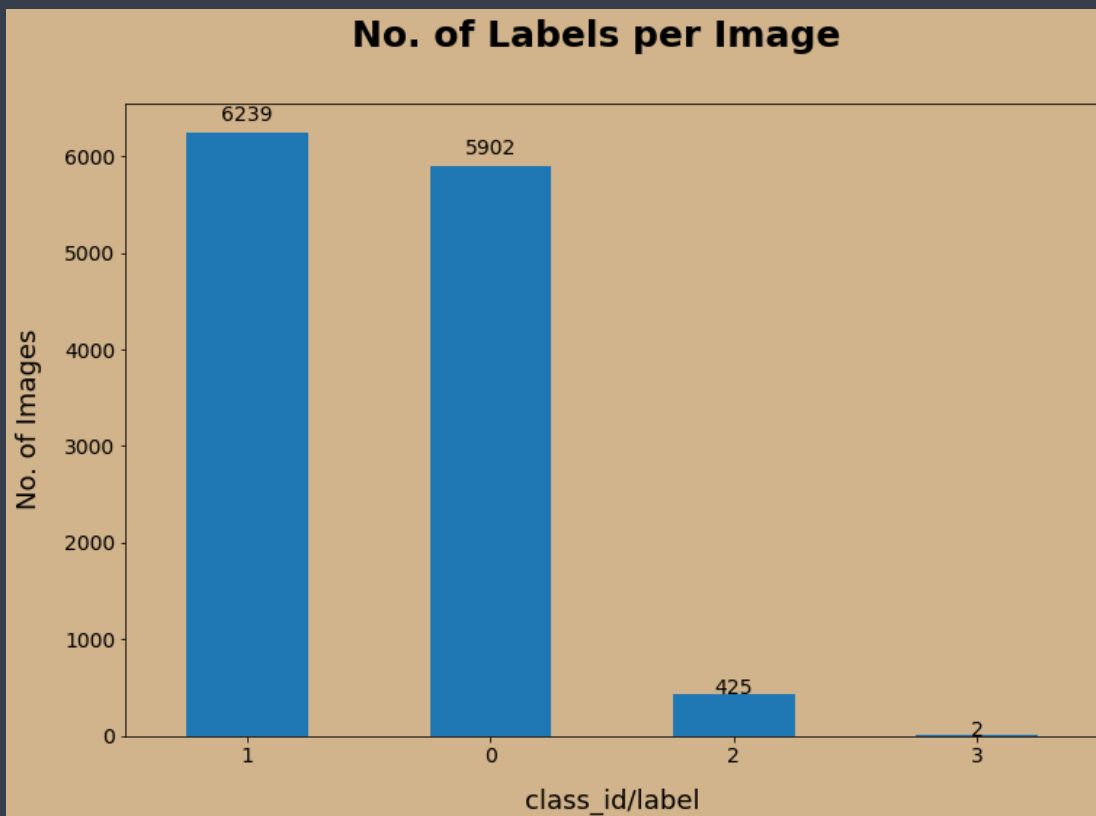
```

```

In [ ]:
fig,ax=plt.subplots(figsize=(12,8))
a=train['total_defects'].value_counts().plot(kind='bar')
patch1(a,ax)
ax.set_xlabel("class_id/label",fontSize=18,labelpad=15)
ax.set_ylabel("No. of Images",fontSize=18,labelpad=15)

```

```
plt.xticks(rotation='horizontal',fontsize=14)
plt.yticks(fontsize=14)
fig.suptitle('No. of Labels per Image',fontsize=25,fontweight='bold')
ax.set_facecolor("tan")
fig.set_facecolor("tan")
plt.show()
```



- There are 5902 images with no labels
- There are 6239 images with 1 label
- There are 425 images with 2 labels
- There are 2 images with 3 labels

- Almost half of images doesn't contain any defects
- Most of images with defects contain the defects of only one type
- In rare cases an image contains the defects of two different types.

```
In [ ]:
#https://www.kaggle.com/paulorzp/rle-functions-run-lenght-encode-decode
def rle_to_mask(rle):
    # CONVERT RLE TO MASK
    if (pd.isnull(rle)) | (rle=='') | (rle=='-1'):
        return np.zeros((256,1600),dtype=np.uint8)

    height= 256
    width = 1600
    mask= np.zeros( width*height ,dtype=np.uint8)

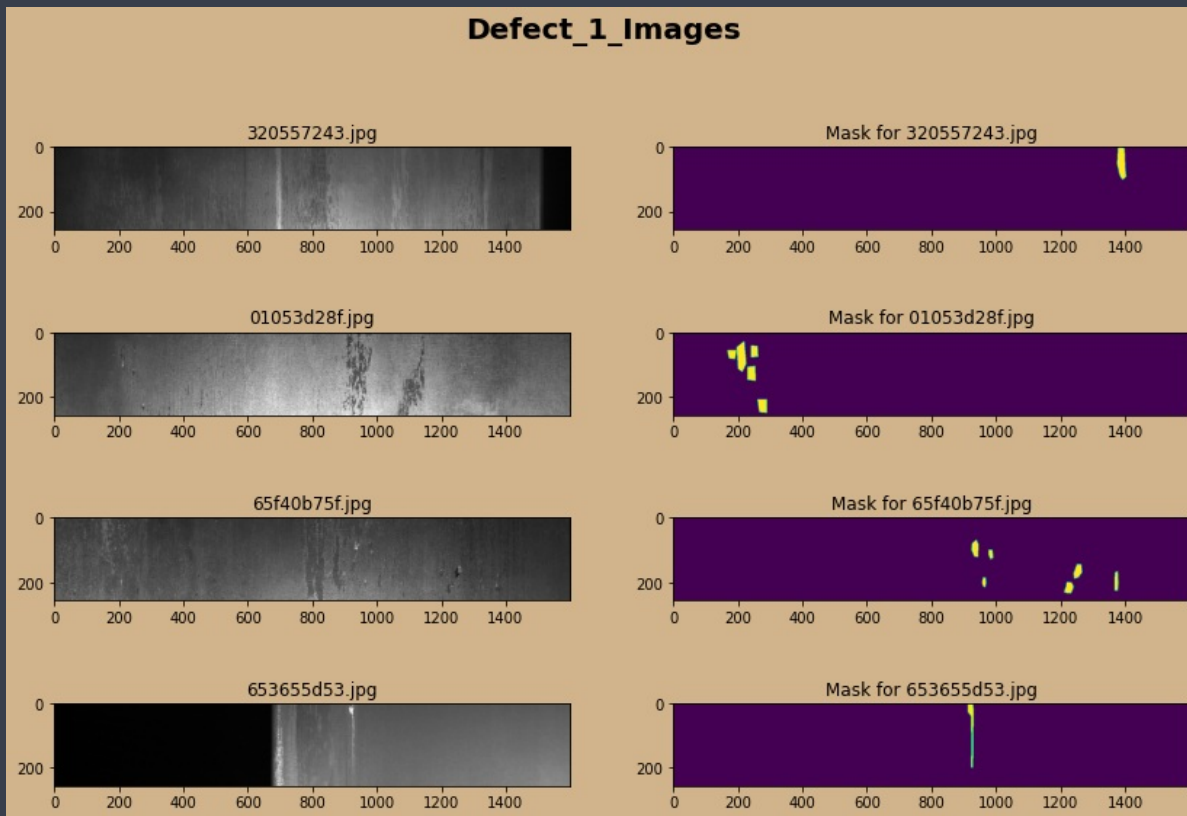
    array = np.asarray([int(x) for x in rle.split()])
    starts = array[0::2]-1
    lengths = array[1::2]
    for index, start in enumerate(starts):
        mask[int(start):int(start+lengths[index])]=1

    return mask.reshape((height,width),order='F')
```

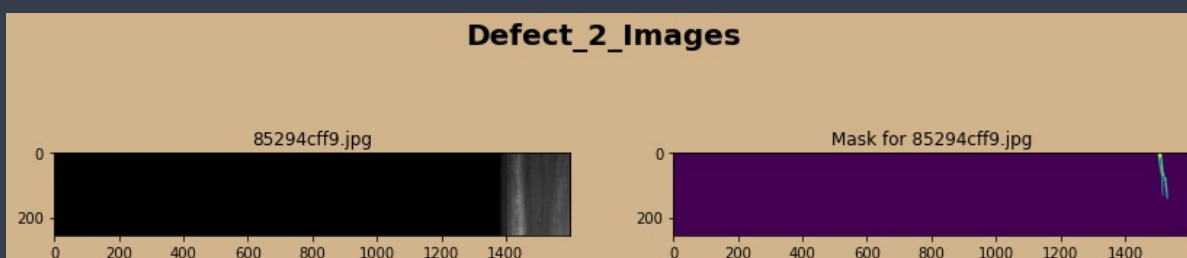
```
In [ ]:
def plot_mask(rle_defect,k):
    x=rle_defect.columns[2]
    train_folder_path='/content/drive//My Drive/Steel_Detection /train_images/'
```

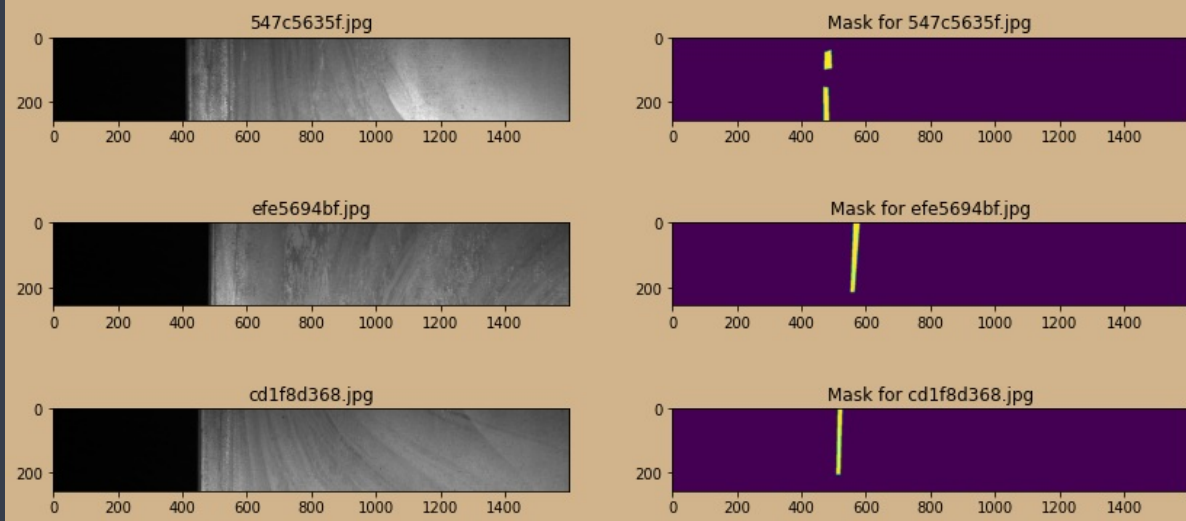
```
# Create figure and axes
fig,ax=plt.subplots(4,2,figsize=(14,9))
fig.suptitle('Defect_'+str(k)+'_Images',fontsize=20,fontweight='bold')
for i in range(4):
    image_id=rle_defect['image_id'][i]
    rle=rle_defect[x][i]
    im=Image.open(train_folder_path+str(image_id))
    ax[i,0].imshow(im)
    ax[i,0].set_title(image_id)
    mask=rle_to_mask(rle)
    ax[i,1].imshow(mask)
    ax[i,1].set_title("Mask for "+str(image_id))
fig.set_facecolor("tan")
plt.show()
```

```
In [ ]:
#https://www.geeksforgeeks.org/how-to-randomly-select-rows-from-pandas-dataframe/
rle_defect=train[train['defect_1']==1]
rle_defect=rle_defect[['image_id','rle_1']]
rle_defect=rle_defect.sample(n=4)
rle_defect=rle_defect.reset_index()
plot_mask(rle_defect,1)
```

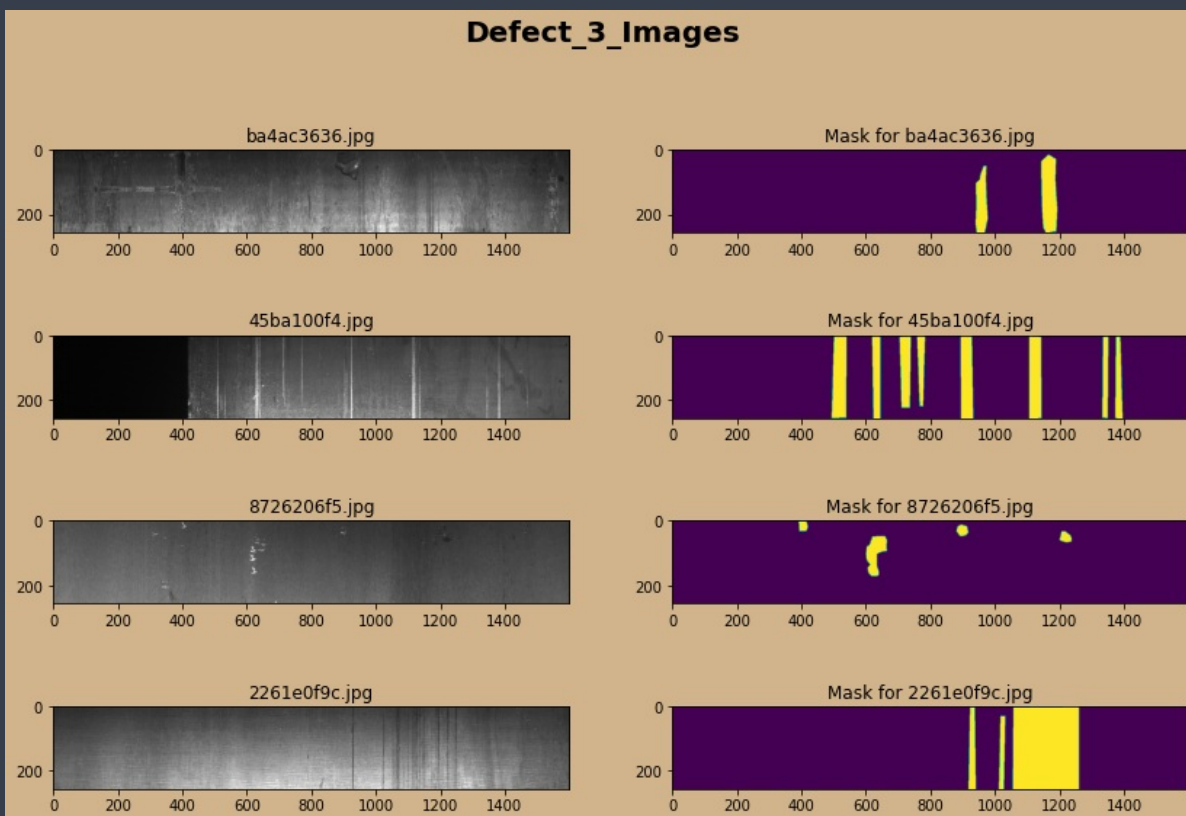


```
In [ ]:
#https://www.geeksforgeeks.org/how-to-randomly-select-rows-from-pandas-dataframe/
rle_defect=train[train['defect_2']==1]
rle_defect=rle_defect[['image_id','rle_2']]
rle_defect=rle_defect.sample(n=4)
rle_defect=rle_defect.reset_index()
plot_mask(rle_defect,2)
```

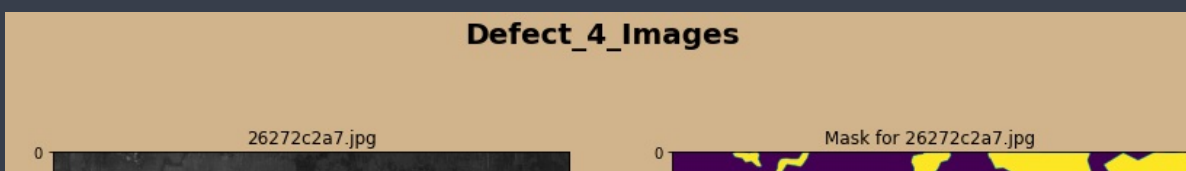


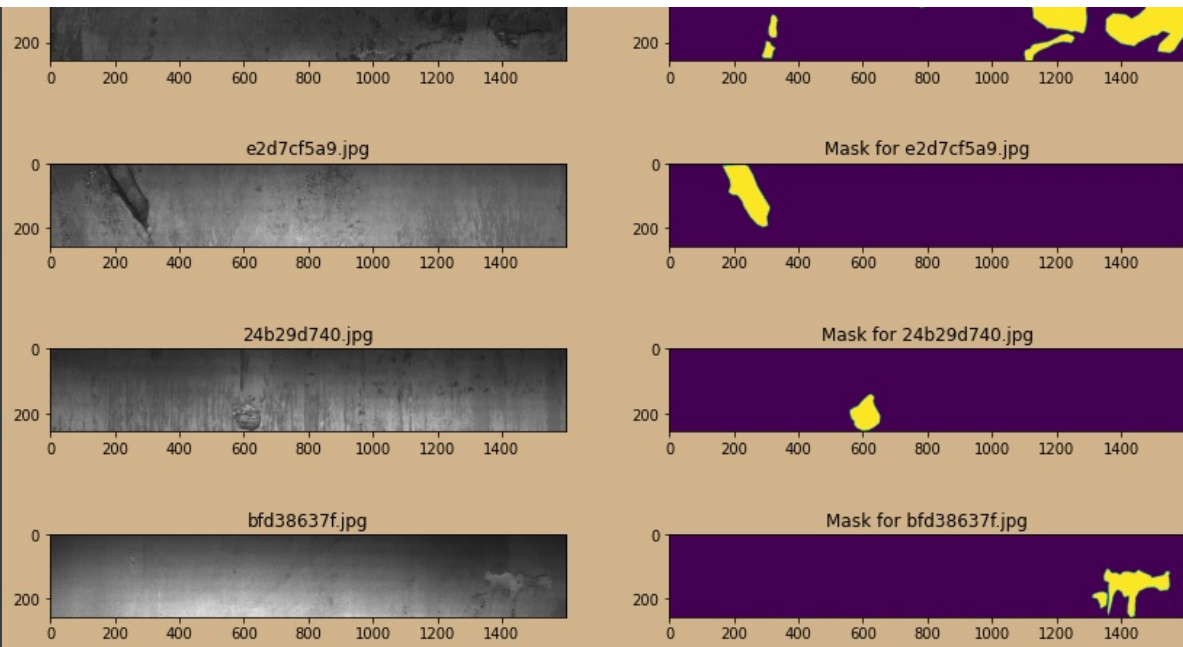


```
In [ ]:
#https://www.geeksforgeeks.org/how-to-randomly-select-rows-from-pandas-dataframe/
rle_defect=train[train['defect_3']==1]
rle_defect=rle_defect[['image_id','rle_3']]
rle_defect=rle_defect.sample(n=4)
rle_defect=rle_defect.reset_index()
plot_mask(rle_defect,3)
```



```
In [ ]:
#https://www.geeksforgeeks.org/how-to-randomly-select-rows-from-pandas-dataframe/
rle_defect=train[train['defect_4']==1]
rle_defect=rle_defect[['image_id','rle_4']]
rle_defect=rle_defect.sample(n=4)
rle_defect=rle_defect.reset_index()
plot_mask(rle_defect,4)
```





```
In [ ]:
def mask_areas(rle_defect):
    area=[]
    for i in rle_defect:
        mask=np.sum(rle_to_mask(i))
        area.append(np.sum(rle_to_mask(i)))
    return area
```

```
In [ ]:
rle_defect=train[train['defect_1']==1]
rle_defect=rle_defect['rle_1']
rle_1_area=mask_areas(rle_defect)

rle_defect=train[train['defect_2']==1]
rle_defect=rle_defect['rle_2']
rle_2_area=mask_areas(rle_defect)

rle_defect=train[train['defect_3']==1]
rle_defect=rle_defect['rle_3']
rle_3_area=mask_areas(rle_defect)

rle_defect=train[train['defect_4']==1]
rle_defect=rle_defect['rle_4']
rle_4_area=mask_areas(rle_defect)
```

```
In [ ]:
fig,ax=plt.subplots(2,2,figsize=(15,7))

ax[0,0].hist(x=rle_1_area)
ax[0,0].set_xlabel("Mask Area : Defect_1",fontsize=13)
ax[0,0].set_ylabel("No. of Images",fontsize=13)
ax[0,0].set_facecolor("tan")

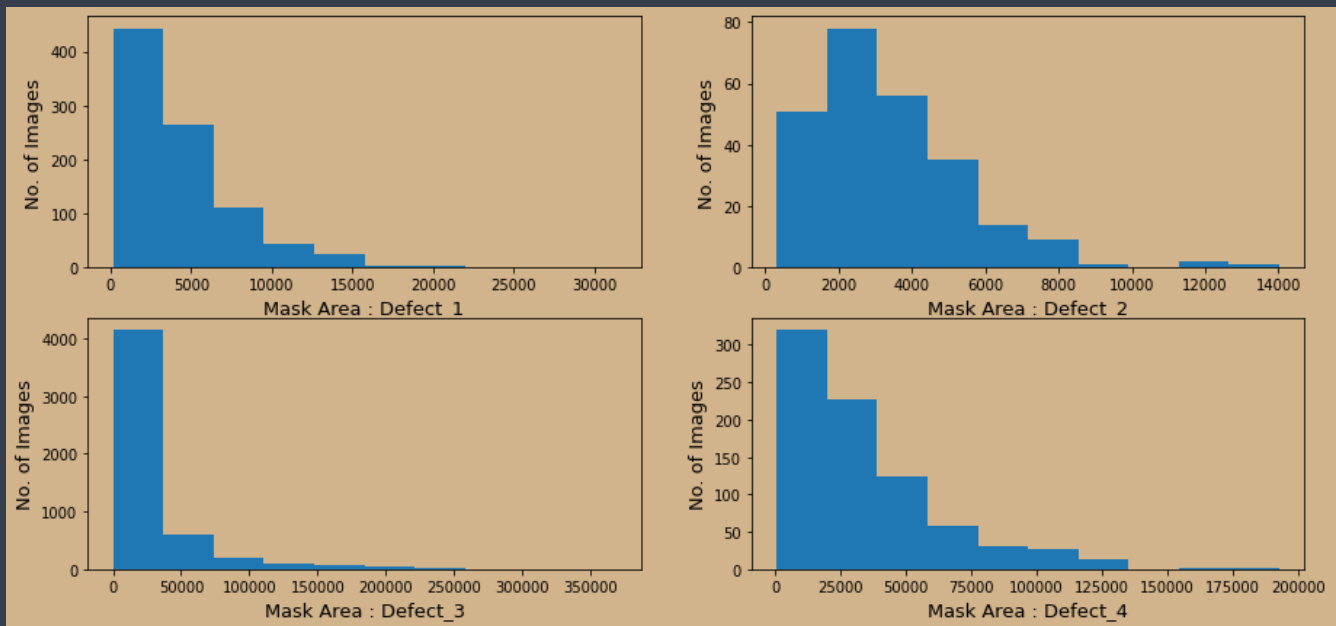
ax[0,1].hist(x=rle_2_area)
ax[0,1].set_xlabel("Mask Area : Defect_2",fontsize=13)
ax[0,1].set_ylabel("No. of Images",fontsize=13)
ax[0,1].set_facecolor("tan")

ax[1,0].hist(x=rle_3_area)
ax[1,0].set_xlabel("Mask Area : Defect_3",fontsize=13)
ax[1,0].set_ylabel("No. of Images",fontsize=13)
ax[1,0].set_facecolor("tan")

ax[1,1].hist(x=rle_4_area)
ax[1,1].set_xlabel("Mask Area : Defect_4",fontsize=13)
```

```
ax[1,1].set_ylabel("No. of Images",fontsize=13)
ax[1,1].set_facecolor("tan")

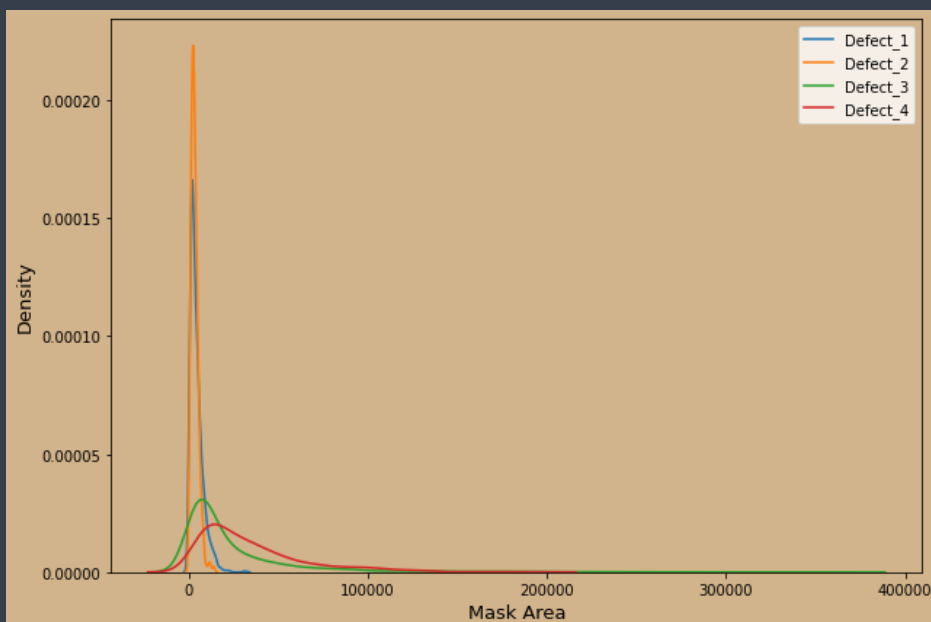
fig.set_facecolor("tan")
plt.show()
```



- Mask area for each defect will help to decide area thresholds during segment prediction (later at the time of modelling).

```
In [ ]:
fig,ax=plt.subplots(figsize=(10,7))

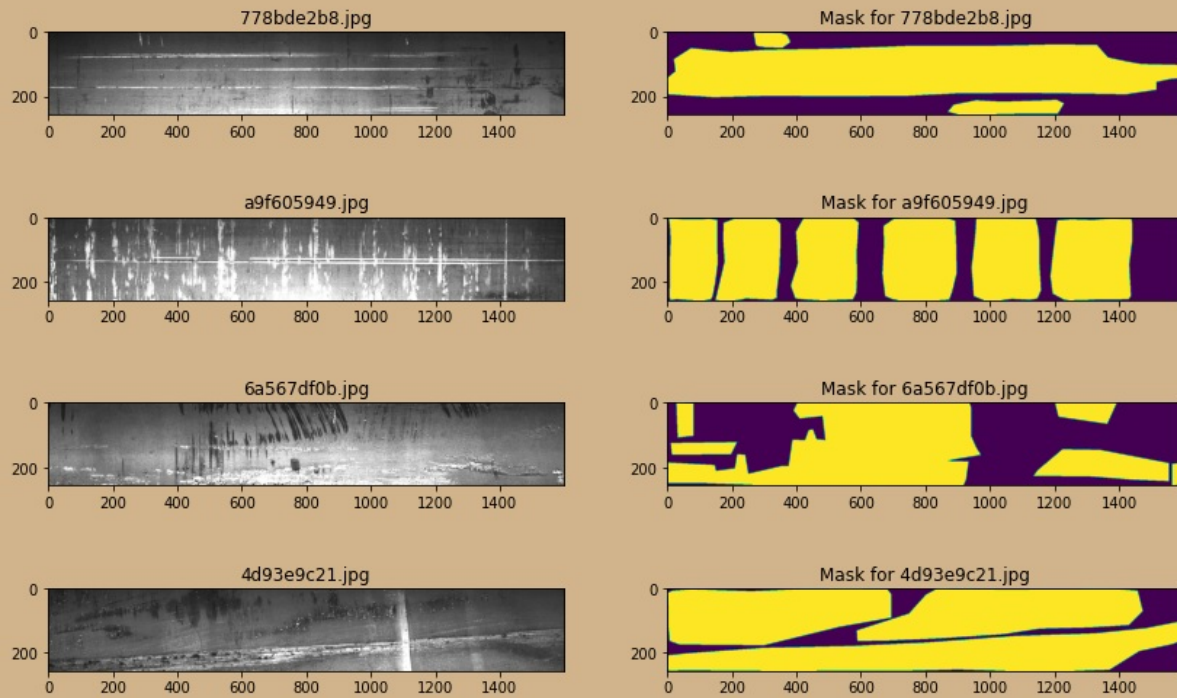
sns.kdeplot(rle_1_area,label='Defect_1')
sns.kdeplot(rle_2_area,label='Defect_2')
sns.kdeplot(rle_3_area,label='Defect_3')
sns.kdeplot(rle_4_area,label='Defect_4')
plt.legend()
ax.set_facecolor("tan")
fig.set_facecolor("tan")
plt.ylabel('Density',fontsize=13)
plt.xlabel('Mask Area',fontsize=13)
plt.show()
```



- Masks with large areas seem very suspicious to me, so I will try to plot few images with large mask areas picked by random index

```
In [ ]:
rle_defect=train[train['defect_3']==1]
rle=rle_defect['rle_3']
rle_3_area=mask_areas(rle)
rle_defect['rle_3_area']=rle_3_area
rle_defect=rle_defect[rle_defect['rle_3_area']>200000]
rle_defect=rle_defect[['image_id','rle_3']]
rle_defect=rle_defect.sample(n=4)
rle_defect=rle_defect.reset_index()
plot_mask(rle_defect,3)
```

Defect_3_Images



- Large masks seem to be okay except for the fact that these masks seem to contain a lot of empty space without any defects

```
In [ ]:
```

```
In [ ]:
```