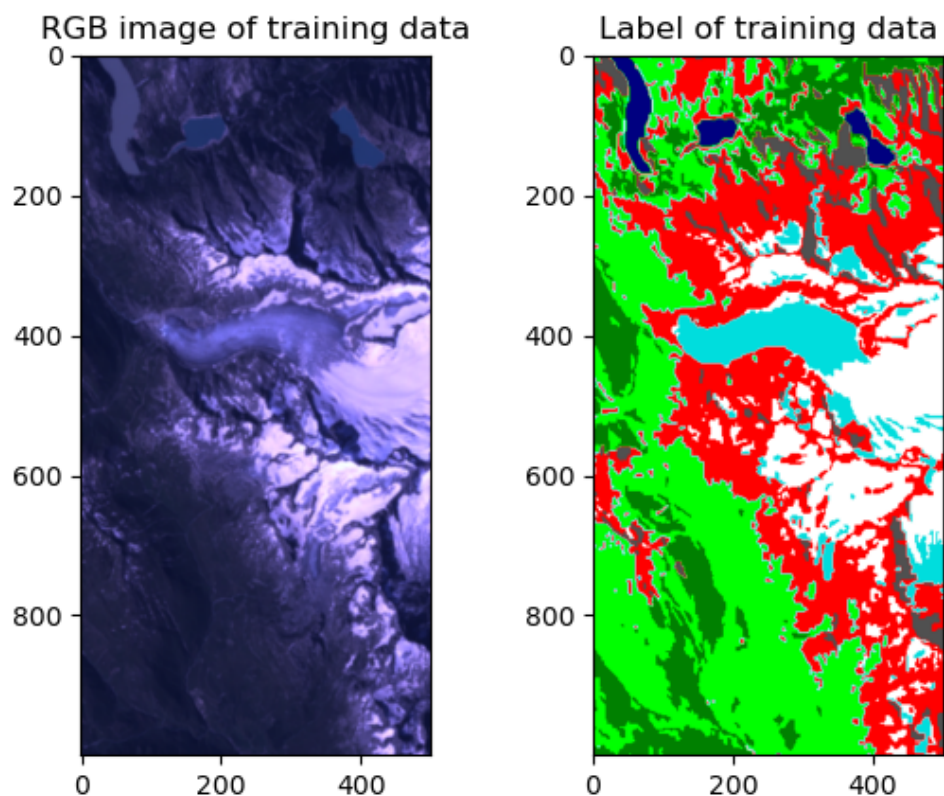


## Introduction

In this exercise, we will use the maximum likelihood method to perform a Gaussian Naive Bayes classification. Two clips of Landsat satellite images with 7 different spectral channels are given. The first clip is used as training data and the second clip is used as test data. The goal is to classify the pixels of the test data into one of the 7 classes. The classes are: water, forest, vegetation, ice, snow, rock, shadow.

## Training data and test data

Firstly we want to visualize the training data and the ground truth label of the training data.



We can see that is hard to distinguish the classes by only looking at the spectral channels. However, we can see that the water and the shadow classes are quite different from the other classes. The water class has a very low reflectance in all spectral channels. The shadow class has a very low reflectance in the visible spectral channels. The other classes have similar reflectance in the visible spectral channels. However, they have different reflectance in the infrared spectral channels. Therefore, we can expect that the infrared spectral channels are more important for the classification.

## Maximum likelihood method

The maximum likelihood method is a method to estimate the parameters of a probability distribution. The probability distribution is assumed to be a Gaussian distribution. The probability density function of a Gaussian distribution is

$$p(\mathbf{x}|\omega_i) = \frac{1}{(2\pi)^{\frac{N}{2}} \sqrt{|\mathbf{K}_{xxi}|}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \mathbf{m}_i)^T \mathbf{K}_{xxi}^{-1}(\mathbf{x} - \mathbf{m}_i) \right\}$$

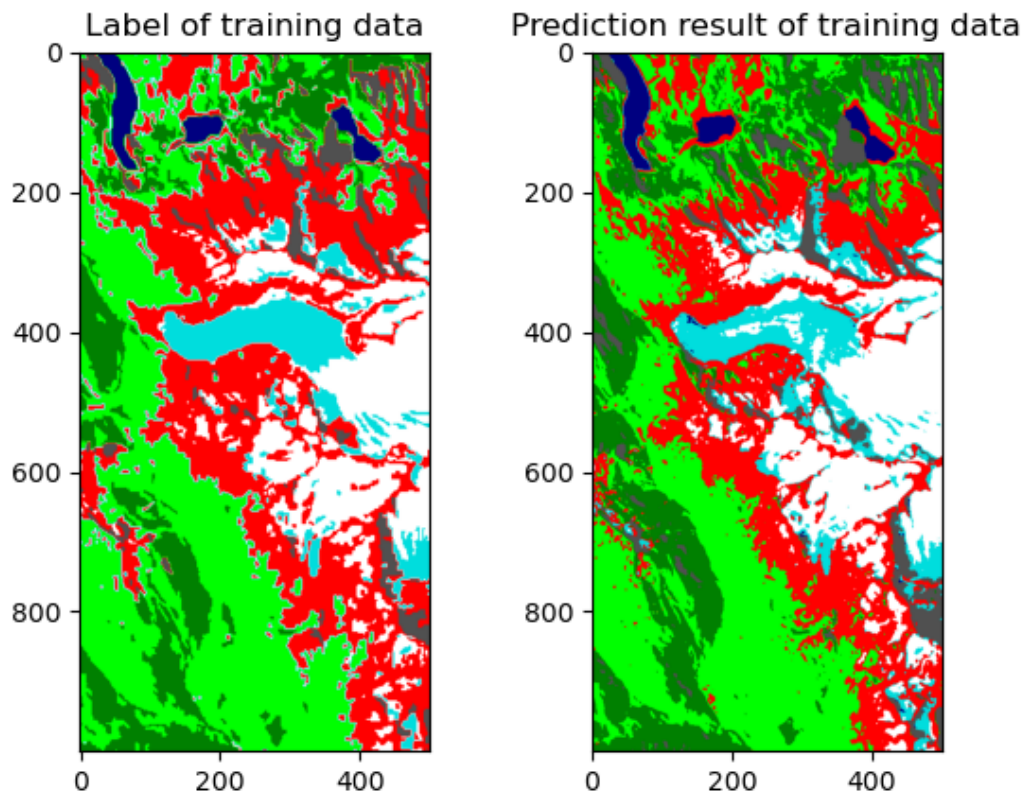
where  $\mathbf{x}$  is a vector of the spectral channels,  $\omega_i$  is the class,  $N$  is the number of spectral channels,  $\mathbf{m}_i$  is the mean vector of the class  $\omega_i$  and  $\mathbf{K}_{xxi}$  is the covariance matrix of the class  $\omega_i$ .

## Training the classifier

Now we have to train our classifier with the training data. We use the maximum likelihood method to estimate the parameters of the Gaussian Naive Bayes classifier.

```
1 clf=GaussianNB()  
2 clf.fit(X_train,y_train)  
3 y_pred=clf.predict(X_train)
```

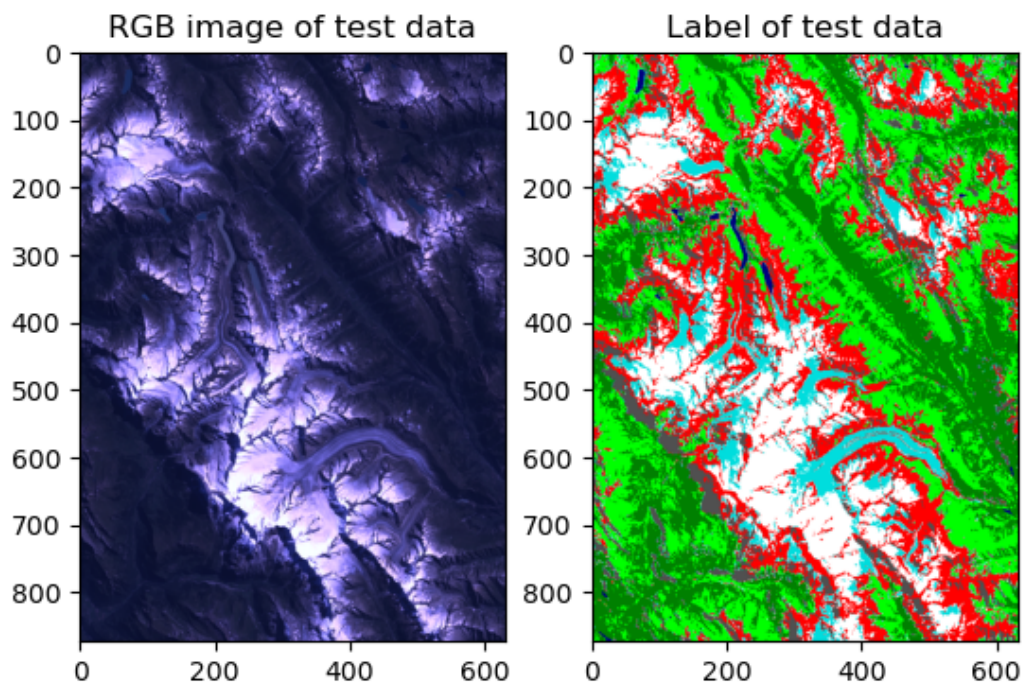
Now we want to visualize the classification result of the training data.



We can see that the water and rock classes are classified very well. However, the other classes for example vegetation and forest, ice and snow are not classified very well. The reason is that the other classes have similar reflectance in the spectral channels. Therefore, it is hard to distinguish them.

## Prediction on the test data

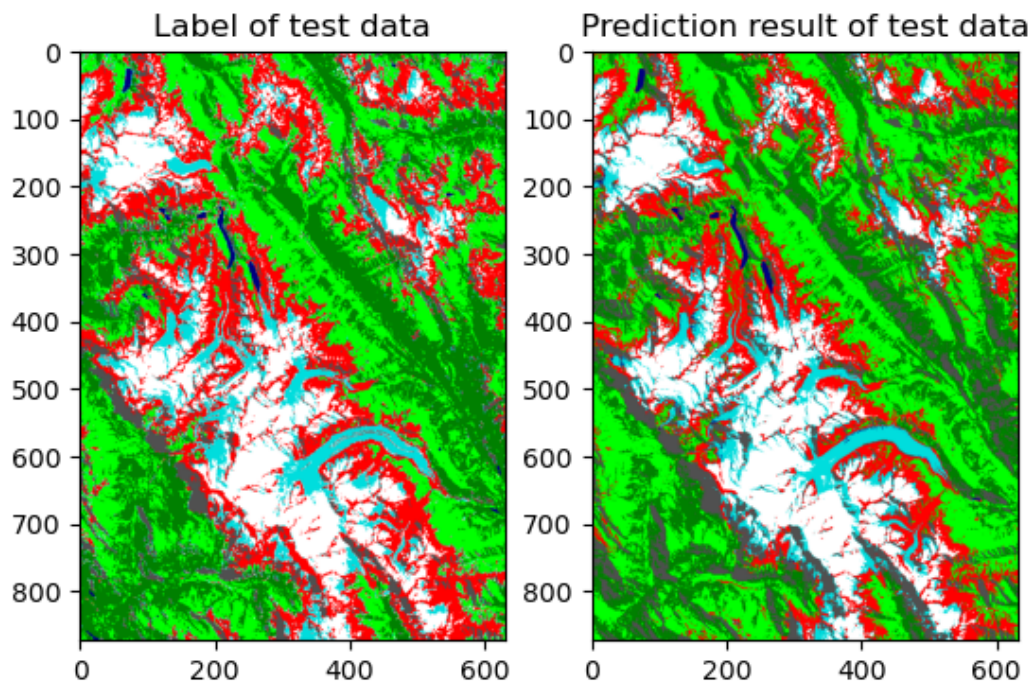
Now we want to visualize the test data and the ground truth label of the test data.



We can predict the classes of the test data with our trained classifier.

```
1 y_pred=clf.predict(X_test)
```

Now we want to visualize the classification result of the test data.



The same as the training data, the classification of the water and rocks are better than other classes. In the test data, the snow and ice are slightly better classified than in the train image. However, the vegetation and forest are still not classified very well.

## Confusion matrix

We can also observe the confusion matrix of the classification result. The confusion matrix is a matrix which shows the number of samples of the true class and the predicted class. The diagonal elements of the confusion matrix are the number of samples which are correctly classified. The off-diagonal elements of the confusion matrix are the number of samples which are incorrectly classified.

There are several accuracies which can be calculated from the confusion matrix. The user's accuracy is the number of correctly classified samples divided by the number of samples which are predicted to be in the class. The producer's accuracy is the number of correctly classified samples divided by the number of samples which are actually in the class. The overall accuracy is the number of correctly classified samples divided by the total number of samples.



Class		Reference data								Accuracy
		Water	Forest	Vegetation	Ice	Snow	Rock	Shadow	Total	User's
Predict	Water	<b>1481</b>	39	12	619	189	94	91	2525	0.586
	Forest	34	<b>81218</b>	11094	2083	1074	7305	1832	104640	0.776
	Vegetation	0	17381	<b>112113</b>	5362	2123	13203	288	150470	0.745
	Ice	3	20	23	<b>29524</b>	7611	6353	369	43903	0.672
	Snow	0	0	0	8515	<b>70736</b>	2507	0	81758	0.865
	Rock	6	94	5717	4447	5252	<b>86618</b>	12710	102495	0.845
	Shadow	401	14564	457	2253	1575	12710	<b>33353</b>	65313	0.511
	Total	1925	113316	129416	52803	88560	128790	36294	<b>551104</b>	
Accuracy	Producer's	0.769	0.717	0.866	0.559	0.799	0.672	0.919		<b>0.769</b>

From the user's accuracy, we can see that the snow and rock classes are classified very well. The forest, vegetation, ice classes are classified quite well. The accuracies of water and shadow are just over 50%. We achieve an overall accuracy of 76.9%.

## Kappa coefficient

The kappa coefficient is a measure of the agreement between the predicted class and the true class. The kappa coefficient is defined as

$$\kappa = \frac{p_o - p_e}{1 - p_e} = \frac{n \sum_{i=1}^m n_{ii} - \sum_{i=1}^m n_{i+} n_{+i}}{n^2 - \sum_{i=1}^m n_{i+} n_{+i}}$$

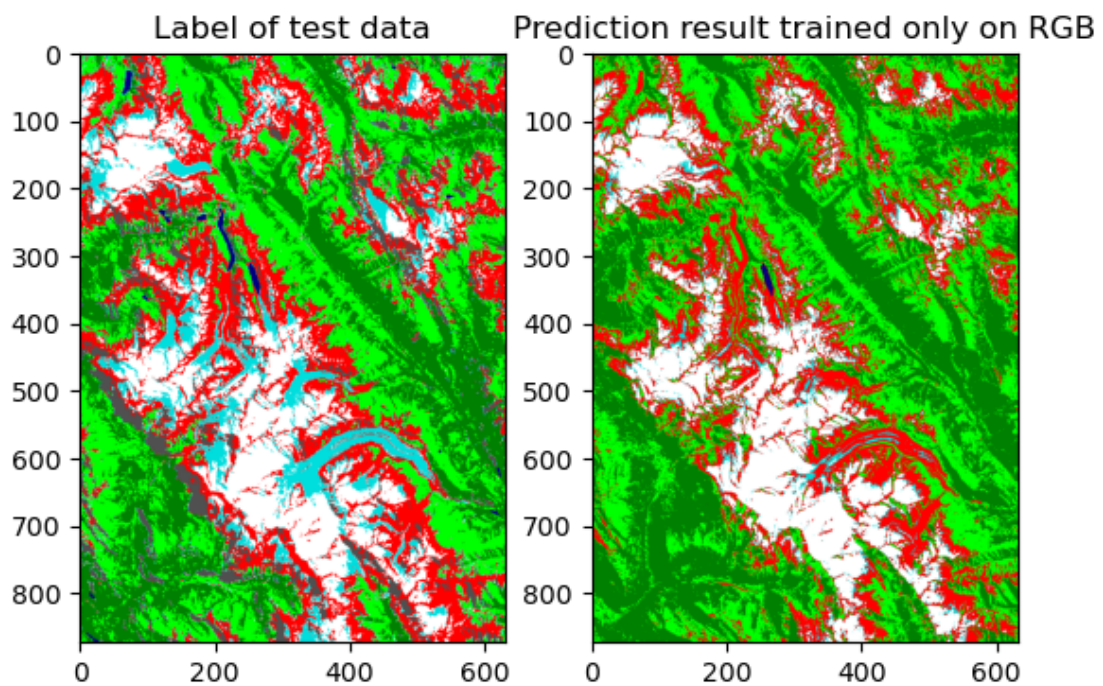
where  $p_o$  is the observed agreement,  $p_e$  is the expected agreement,  $n$  is the total number of samples,  $n_{ii}$  is the number of samples which are correctly classified,  $n_{i+}$  is the number of samples which are predicted to be in the class  $i$ ,  $n_{+i}$  is the number of samples which are actually in the class  $i$  and  $m$  is the number of classes.

Kappa coefficient: 0.6967297666414108

The result of kappa coefficient is 0.6967. The kappa coefficient is between 0 and 1. The kappa coefficient is 1 if the predicted class is the same as the true class. The kappa coefficient is 0 if the predicted class is the same as the true class by chance. The kappa coefficient is negative if the predicted class is worse than the true class. Therefore, we can say that our classifier is much better than random guessing.

## Train only on RGB channels

Now we want to train our classifier only on the RGB channels. We want to see if the infrared channels are important for the classification.



Class		Reference data								Accuracy
		Water	Forest	Vegetation	Ice	Snow	Rock	Shadow	Total	User's
Predict	Water	<b>349</b>	345	112	0	0	1119	0	1925	0.181
	Forest	6	<b>100562</b>	12576	0	0	172	0	113316	0.887
	Vegetation	9	29693	<b>92876</b>	0	0	6838	0	129416	0.718
	Ice	5	2771	7041	<b>9826</b>	9467	23693	0	52803	0.186
	Snow	2	1338	3815	2868	<b>70922</b>	9615	0	88560	0.801
	Rock	1	6854	33914	5716	4546	<b>77753</b>	6	128790	0.604
	Shadow	0	30079	5670	0	0	524	<b>21</b>	36294	$5.786 \cdot 10^{-4}$
	Total	372	171642	156004	18410	84935	119714	27	<b>551104</b>	
Accuracy	Producer's	0.938	0.586	0.595	0.534	0.835	0.650	0.778		<b>0.639</b>