**Machine Learning for Soil and Crop Management**
**Professor. Somsubhra Chakraborty**
**Agriculture and Food Engineering Department**
**Indian Institute of Technology, Kharagpur**
**Lecture 28**
**Use of ML for Portable Proximal Soil and Crop Sensors (Contd.)**
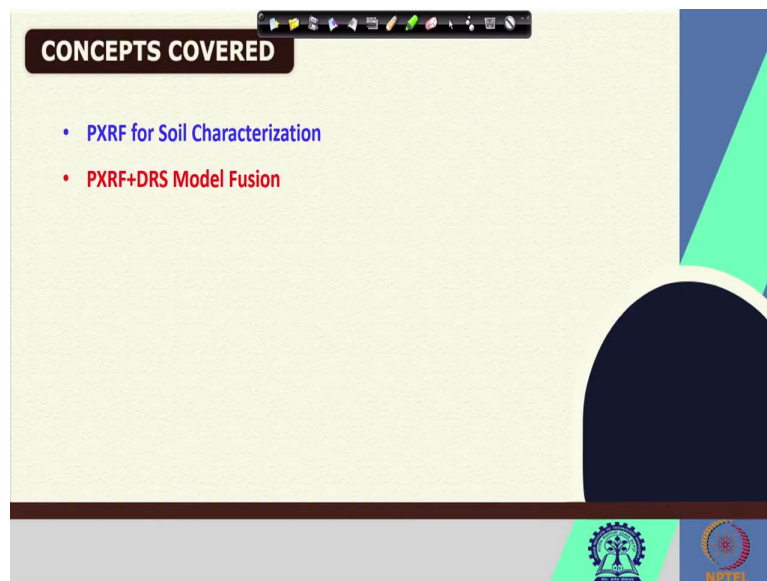
(Refer Slide Time: 00:21)



Welcome friends to this third lecture of week 6 of this NPTEL online certification course of Machine Learning for Soil and Crop Management. And in this week, we are talking about the application of machine learning for portable proximal soil and crop sensors. In our first two lectures, we have discussed about the framework of smart soil sensing and how it is related with site specific nutrient management.

And also, we have seen the broad classification and examples of proximal soil sensors and the properties they can measure. We have also seen the principle of portable XRF, working principle a portable XRF. And in the previous lecture specifically we focused on evolution of PXRF for characterizing different soil properties and how it was used. And also, how it moved from the application of simple statistical linear regression relationship to the machine learning based application for characterizing different soil properties.
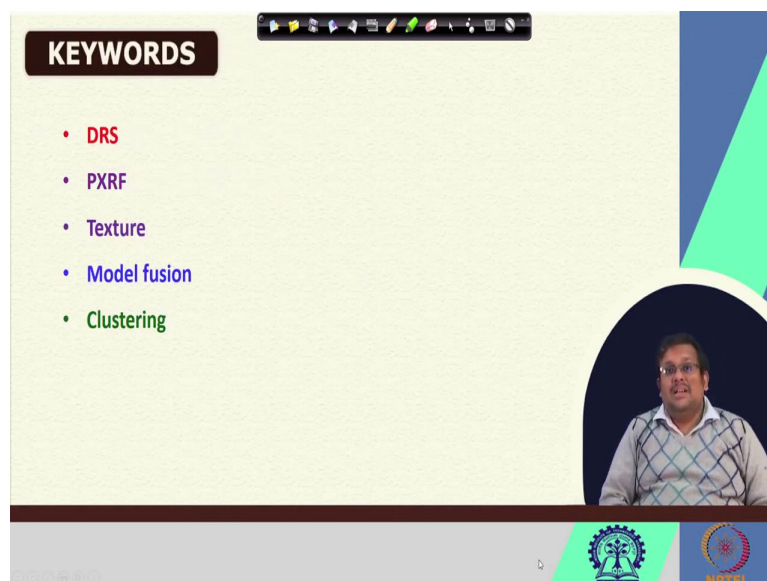
We have seen its application for soil pH, EC, CEC, compost pH, compost CEC measurement, and also, we have seen their application for measurement of heavy metals in the industrial sites. So, today, we will start from there, and then we will see some other applications. And then, we will also discuss the sensor fusion aspect, where we can combine multiple sensors to augment or to get the synergistic prediction of soil properties.

(Refer Slide Time: 02:37)



So, these are the major two topics for today's lecture or this lecture number 28. First of all, we are going to discuss the PXRF for soil characterization, and then we will see the PXRF plus DRS model fusion.

(Refer Slide Time: 02:55)



Now, these are the some of the keywords which we are going to discuss today in this lecture, DRS, PXRF, Texture, texture by texture, we mean soil texture, and also model fusion, and also clustering.
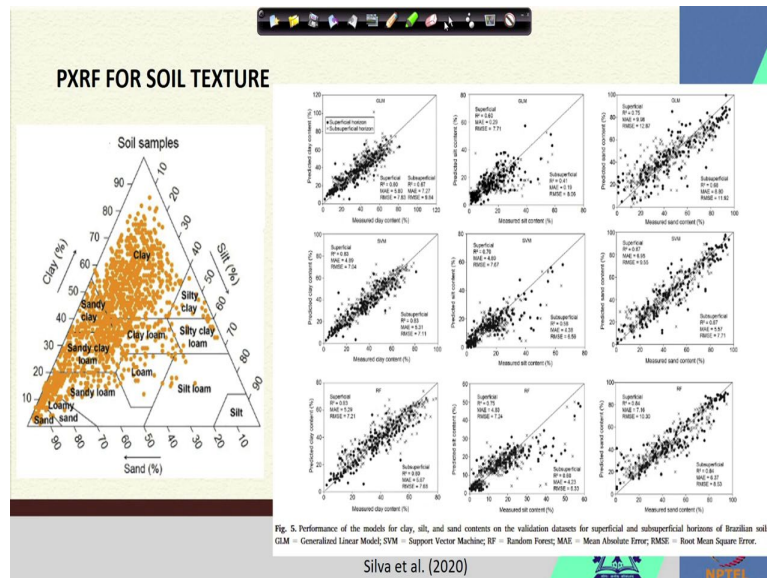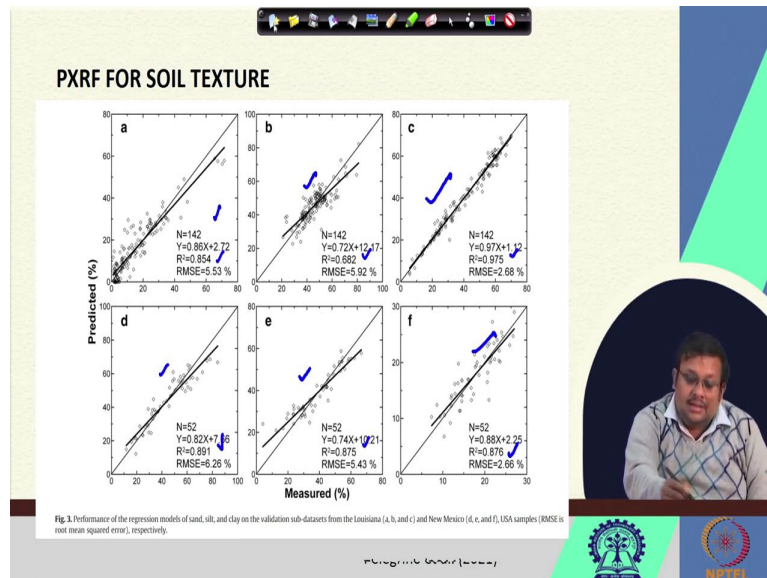
So, another important application for PXRF is soil texture. And soil texture it is a very important parameter which controls different soil physical-chemical properties. So, in this study in 2020, the Brazilian scientist have proved that PXRF elemental content can be also used for measurement of textual parameters like sand, silt and clay using different models, machine learning models.

So, they have used the textural PXRF elemental content for predicting the sand, silt and clay. As you can see here in this plot, this is clay content by GLM or generalized linear modeling, then generalized linear model here this is a support vector machine model and this represent random forest model. Similarly, this is the GLM model for silt content, this is the SVM model for silt content and this is the random forest model for the silt content.

Similarly, this is the GLM model for sand content, this is support vector machine for sand content and this is the random forest model for sand content. In all the cases you can see the model are, they have modelled both the superficial horizons, as well as the sub superficial horizons.

And you have seen that more or less this type PXRF can be utilized for prediction, rapid prediction of soil texture, instead of relying on time consuming particle size analysis, we can use the PXRF elemental content as a proxy for measurement of clay, silt and sand, because the presence of sand, silt and clay is generally they indicate the presence of certain elements. And the coexistence of these elements with sand, silt and clay particles shows their applicability for predicting the soil particle size analysis. So, this is one important application.

(Refer Slide Time: 06:12)



Fig. 3. Performance of the regression models of sand, silt, and clay on the validation sub-datasets from the Louisiana (a, b, and c) and New Mexico (d, e, and f), USA samples (RMSE is root mean squared error), respectively.

Similarly, Zhu et al in 2011 also proved the utility of PXRF for rapid measurement of soil texture. So, here you can see that these are sand, silt and clay content, these three plots represent sand, silt and clay content for one location and this is for another location. And they have utilized the stepwise linear regression model for producing this equation and you can see for sand content they have got 0.85 to 0.89 R square, for silt content 0.68 to 0.87 and for clay content they got very high that is 0.97 to 0.87 R squared values, so that this is the first application of PXRF for prediction of soil texture.

(Refer Slide Time: 07:31)



PXRF FOR SOIL NUTRIENTS

Table 3 Equations and coefficient of determination ($R^2$) of models generated by portable X-ray fluorescence (pXRF) spectrometry data to predict soil nutrient content in Brazil

| Soil nutrient | Method | Model | $R^2$ |
|---|---|---|---|
| $Ca^{2+}$ | LR | $Ca^{2+} = 1.5586 + 0.0006Ca$ | 0.84 |
| | PR | $Ca^{2+} = 1.87 + 0.0005Ca + 0.000000007Ca^2$ | 0.84 |
| | PwR | $Ca^{2+} = 2.2184e^{0.0001Ca}$ | 0.70 |
| | SMLR | $Ca^{2+} = 4.48 + 0.249Al + 2.32Ca + 0.60Ni + 0.31Zr$ | 0.78 |
| $K^+$ | LR | $K^+ = 157.03 + 0.0003 K$ | 0 |
| | PR | $K^+ = 155.53 + 0.001 K - 0.00000002 K^2$ | 0 |
| | PwR | $K^+ = 138.89e^{0.00005K}$ | 0 |
| | SMLR | $K^+ = 158.16 - 18.99Al + 44.77Si + 19.44P - 19.83Mn + 73.96Fe - 16.50Ni$ | 0.17 |
| P | LR | $P = 10.047 + 0.0306P$ | 0.39 |
| | PR | $P = 8.6829 - 0.0453P + 0.00005P^2$ | 0.87 |
| | PwR | $P = 1.9622e^{0.0015P}$ | 0.56 |
| | SMLR | $P = 10.88 - 10.36Si + 13.19P - 11.18Ti + 15.19 V - 11.82Cr$ | 0.37 |

LR linear regression, PR 2nd degree polynomial regression, Pw power regression, SMLR stepwise multiple linear regression
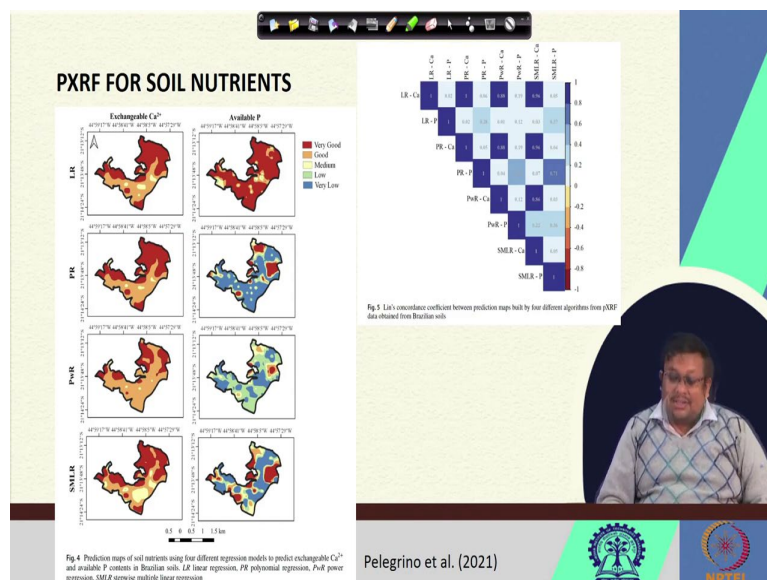
Pelegrino et al. (2021)

And another very recent application which was published last year in precision agriculture journal and there we have proved that px serif can be utilized for prediction of certain soil

available nutrients. Here, you can see the exchangeable calcium and also available potassium and available phosphorus were predicted using different types of models like linear regression, then power regression.

So, this is linear regression, then second degree polynomial regression, then power regression model, then stepwise multiple linear regression model. So, these models were tried in for three different parameters and it was seen that for calcium we are getting very high R square values and highest R square you can get for LR and second-degree polynomial.

Whereas, in case of phosphorus also, we are getting high R square values using the PR regression. So, that shows that utilizing the PXRF elemental content, it is possible to predict the nutrient content available nutrient content not only you can predict the nutrient content, but also you can produce the special variability map of nutrient.
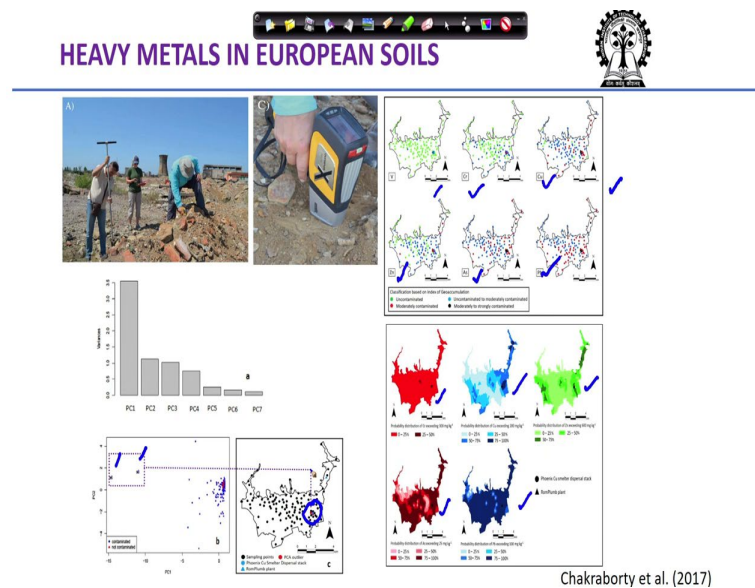
(Refer Slide Time: 09:15)



So, similarly, they have done the special variability map of the calcium, exchangeable calcium and available phosphorus using these four different models. As you can see here, linear regression model in the first panel followed by the second-degree polynomial regression, then power regression and in the last there will be stepwise multiple linear regression.

Here also, we can see that Lin's concordance coefficient between prediction maps, map built by four different algorithms from PXRF data obtained from the Brazilian soil. And we can see how these values are coming related with each other. So, that shows an important

application because whatever we do using PXRF from the soil science point of view, its utility can be well justified, if we can show its application for nutrient measurement.

So, if we can use this instrument as an alternative to the traditional weight chemistry-based measurement of available nutrient, that will certainly make a paradigm shift and that will change and that will change the dependency of the farming community towards the traditional soil testing processes.

(Refer Slide Time: 10:49)



Chakraborty et al. (2017)

PXRF has also been used for heavy metal contamination in polluted areas. So, in this research, in 2017, we have utilized this PXRF instrument as you can see in this picture to on site scan the soil samples in an industrial area, abandoned industrial area of eastern Europe in Romania.
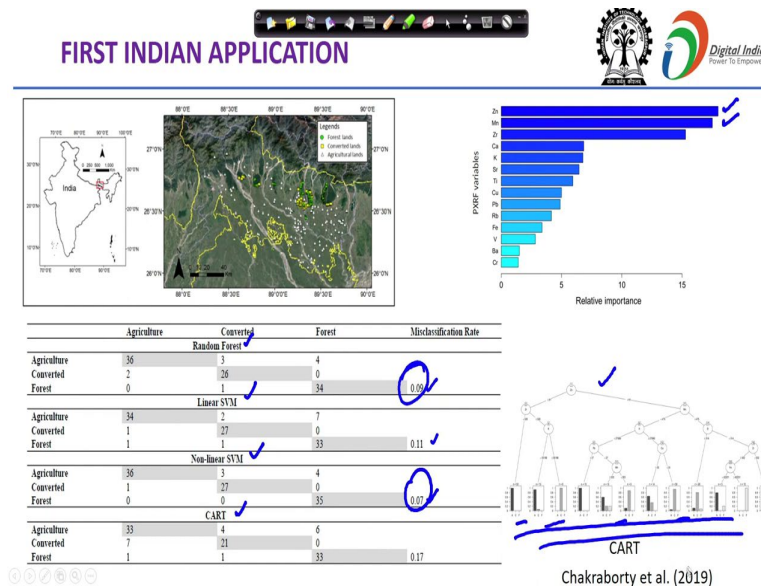
And after collecting the scan results, we did the principal component analysis and using the principal component analysis, we try to explore the relationship between the day, between the elements heavy metals to cluster them, and also to identify their source. We are also calculated several pollution index. And for example, index of geo-accumulation as you can see in this map.

So, these index of geo-accumulation were created for Vanadium, Chromium, Copper, Lead, Arsenic and Zinc. And we can classify the samples based on their index of geo-accumulation and after classifying we have identified the reason behind this classification categories. So, we seen that the presence of the sampling points nearby the industrial sites, the major reason behind their classification as highly contaminated sites.

So, not only that, but also, we use the Indicative Kriging approach for producing the special variability map of different heavy metals like here you can see Chromium, then Copper, then Zinc, then Arsenic and Lead. So, these special variability maps are also been created for identification of the pollution hotspots.

So, also the PCA score, from the PCA score plot, we try to identify the outlets and we try to locate those outlets and we have seen that those outlets are very close to the smelter dispersal stack. So, that represents that this is, that shows that the nearby presence of this sampling size to the original contamination source is responsible for their abnormally high concentration. So, that shows the utility of these PXRF for heavy metal contamination mapping.

(Refer Slide Time: 13:57)



And first Indian application was seen in the year 2019 where we have used around 500 soil samples collected from the northern part of eastern India, I mean in the state of West Bengal. The northern part is hilly region dominated by forests, so we collected 500 soil samples from three different land use type. So, these land use where forest land use, also agriculture and converted. Converted means, these lands where actually forest land 20 years back, but then they slowly started converting into agricultural fields.

So, we collected around 500 samples from all these three land use types. And then we use the different classification methods using the random forest, their linear support vector machine, nonlinear support vector machine and classification and regression tree. And we compared their misclassification rate. You can see the misclassification rate of the random forest and also the nonlinear SVM are very close. And that shows the misclassification rate is quite small, that shows in other words, the higher classification accuracy.

So, that means that PXRF can be used in future for assigning the proper class of land use to any unknown samples coming from those three land use types. So, not only we have compared their misclassification rate, but also the random forest relative importance of PXRF variables were identified, whereas Zinc and Manganese were the two most influential parameters.

Not only that, we also utilize this classification regression tree to form the rules to segregate the classes based on the PXRF elements. So, this is the classification and regression tree. These forms the rules of segregating the samples or assigning the samples into one of these terminal nodes. I would request you to please go through this literature and see this application. So, to get more comprehensive information about this type of classification. And this was the first Indian application of PXRF in soil.

(Refer Slide Time: 17:12)



Also, PXRF elemental content was used to differentiate the rural and urban lakes in the United States. So, the lakes or playas they call it, they have different elemental content and their variation in the rural areas as well as urban areas due to variation in their land use. And as a result of that, we are interested to see whether PXRF could be able to separate the urban samples from the rural samples or not. So, principal component analysis was executed.
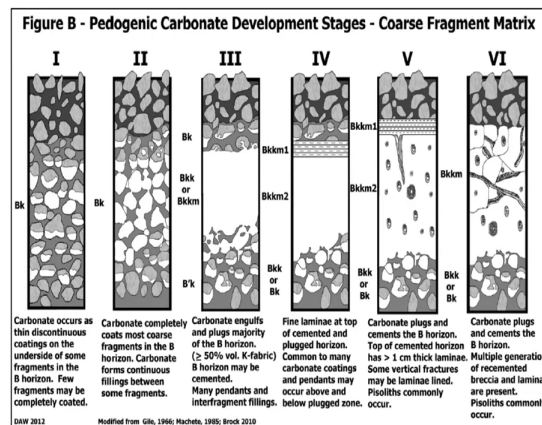
And then classification was done using the principal component 1 and principal component 2. Principal component 1 and 2 cumulatively produce around 80 percent of the variability more than 80 percent of the variability. And we have seen that classification using the principal component based on the principal component score where more or less uniformly or efficiently separated the samples coming from the rural playas as well as the urban lakes.

This is the confusion matrix showing the random forest-based classification of the playas or lakes of a county in Texas and these cells these. So, here is the rural versus rural there are 90 plus correctly classified samples, then urban versus urban, they are 96, there are some hybrid playas also, which 15 correctly classified.

So, they are a total 219 samples and from there we can see that majority of the, that is 91.78 percent of the samples were correctly classified. So, that also shows that PXRF is able to segregate the sample based on their geochemistry. So, that also made an important paradigm shift in the soil characterization, specially when we consider the geochemical parameters.

(Refer Slide Time: 19:42)





Another very important application for PXRF was to segregate the carbonate stage development. Now, there are different carbonate stage development like here you can see 1 to

6 there are paedogenic carbonate development stages for both fine and coarse fragment metrics. So, there are 6 different stages. So, here you can see stage 1 and stage 5.

(Refer Slide Time: 20:13)



Carbonate Stage Development

- Is there a correlation between developmental stage and $CaCO_3$ content?
- If so, can PXRF be useful in helping to determine developmental stage?
- 75 samples collected across four states
  - Texas, New Mexico, Colorado, Kansas
  - Represented all six developmental stages
  - Samples scanned as intact aggregates and ground powders (<2mm)
  - Avg. developmental stage determined independently by a panel of five experienced pedologists (Soil Survey Staff)

So, the question was, is there any correlation between the development stage and the calcium carbonate content? And if so, can PXRF be useful in helping to determine the development stage. So, 75 samples were collected across four states like Texas, New Mexico, Colorado and Kansas and which represented all six development stages and sample scanned as intact aggregates and ground powders.

So, we first ground the, we first scan the soil samples as intact aggregates as well as we subsequently ground the soil samples into powders and sieve them with 2 millimetre sieve and then we can we again scan them using PXRF.

(Refer Slide Time: 21:11)



So, average development stage are determined independently by a panel of 5 experienced radiologist, they also did, but we have seen that panelist only unanimously agreed to only 22 percent of the samples, evaluated ex-situ.
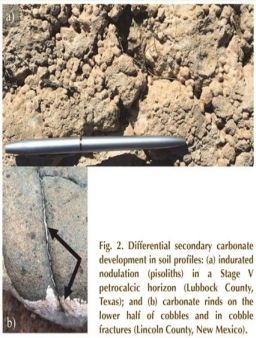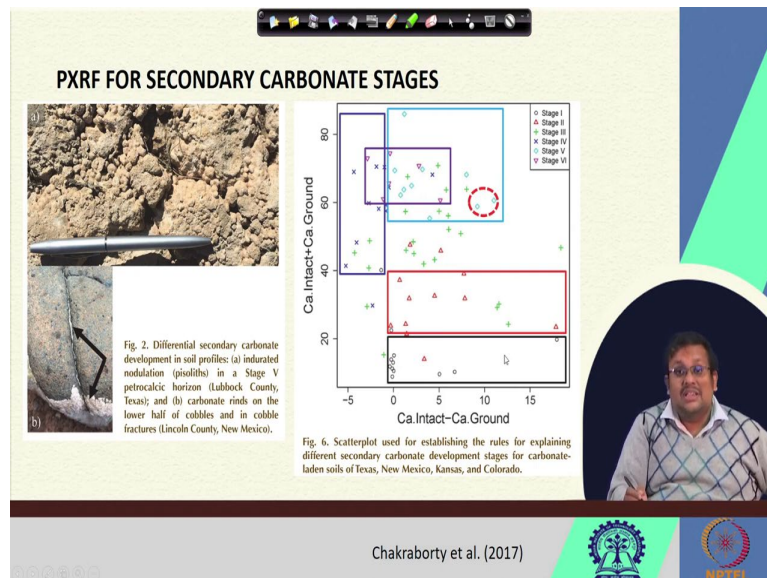
(Refer Slide Time: 21:20)



So, that means that in the most of the cases there is a disagreement between the Phonologists while describing these 6 development stages of calcium carbonate. So, we thought that let us use the portable XRF for identifying these development stages of the carbonate or developments stages of the secondary carbonate. So, there are 6 development secondary stages of the carbonate and you can see here we have utilized mathematical formulas to

identify or to establish certain rules for identification of the 6 secondary carbonate level stage, development stage first, second, third, fourth, fifth and sixth.

So, the rules were developed based on our judgment from one of our figures, which is created because Calcium Intact and Calcium Ground were highly correlated. And some stages where we have seen that these Calcium Intact soil and Calcium in ground soil were highly correlated, and some stages were above the diagonal line.
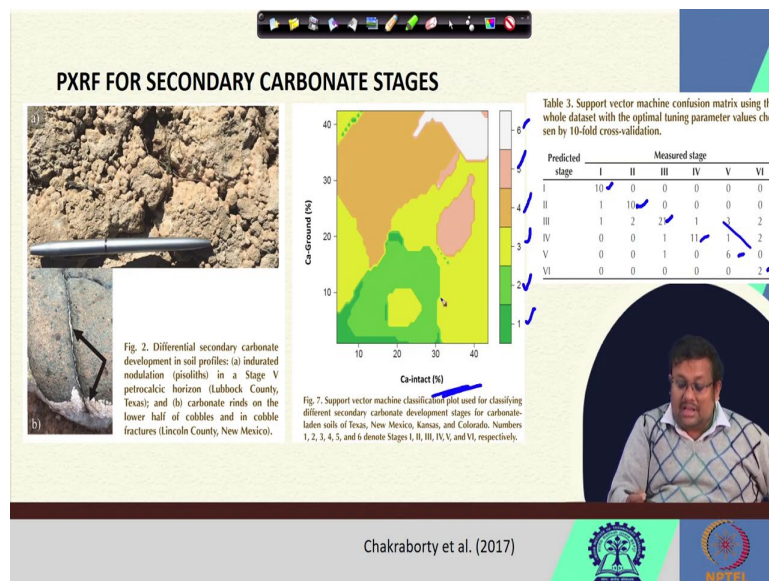
So, we utilized for this rule formation, you can see one thing guys, for rule formation, we have utilized this Calcium Intact and Calcium Ground only and then we formed these rules. So, based on the calcium content only using these rules, we can separate these stages. So, that shows another good application of PXRF for identification of the development stage.

(Refer Slide Time: 23:03)



And this is the scatter plot used for establishing the rules for explaining different secondary carbonates development stages based on Calcium Carbonate, Calcium Intact minus Calcium Carbonate Calcium Ground in the x-axis and Calcium Intact plus Calcium Ground in the y-axis. And you can see that utilizing these two axes, you can separate out or you can produce the boundary for the samples belongs to each of these 6 stages.
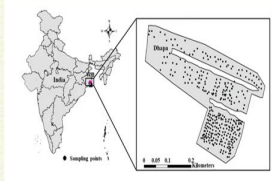
Fig. 2. Differential secondary carbonate development in soil profiles: (a) indurated nodulation (pisoliths) in a Stage V petrocalcic horizon (Lubbock County, Texas); and (b) carbonate rinds on the lower half of cobbles and in cobble fractures (Lincoln County, New Mexico).

Fig. 7. Support vector machine classification plot used for classifying different secondary carbonate development stages for carbonate-laden soils of Texas, New Mexico, Kansas, and Colorado. Numbers 1, 2, 3, 4, 5, and 6 denote Stages I, II, III, IV, V, and VI, respectively.

Table 3. Support vector machine confusion matrix using the whole dataset with the optimal tuning parameter values chosen by 10-fold cross-validation.

Chakraborty et al. (2017)

And based on this, we have utilized a support vector machine confusion matrix and we have seen the support, based on the support vector machine calculation, we can see the most of the stages were correctly classified using the PXRF based elements. So, here you can see these are the measured stages and predicted stages and this diagonal line is showing the correctly classified sample. So, you can see these correctly classified samples ultimately shows that the PXRF is very much helpful for quantitatively identifying these stages.

And these SVM classification plot we have already seen it previously in one of our lecture. So, this support vector machine classification plot used was used for classifying the different secondary carbonate development states. So, you can see here which are colour coded 1, 2, 3, 4, 5, 6; 6 stages 6 colour and then they are colour coded to explain the stages based on the Calcium Intact and Calcium Ground. Calcium Intact sample as well as Calcium in ground sample.

One of my PhD student has recently utilize this PXRF also in India for landfill soil characterization. So, we collected around 335 samples from a landfill soil adjacent agricultural fields and we characterized the soil samples using PXRF and we focused on 7 heavy metals like Chromium, Manganese, Copper, Zinc, Arsenic, Lead and Mercury. So, these 7 heavy metals were considered and we have compared their abundance with the local background value as well as the threshold levels, geochemical threshold levels.

Now, this is the landfill site and these are the agricultural, surrounding agricultural fields from which we have collected the samples. And the major goal, the major reason for undertaking this research was to establish whether PXRF can be used to predict the pollution
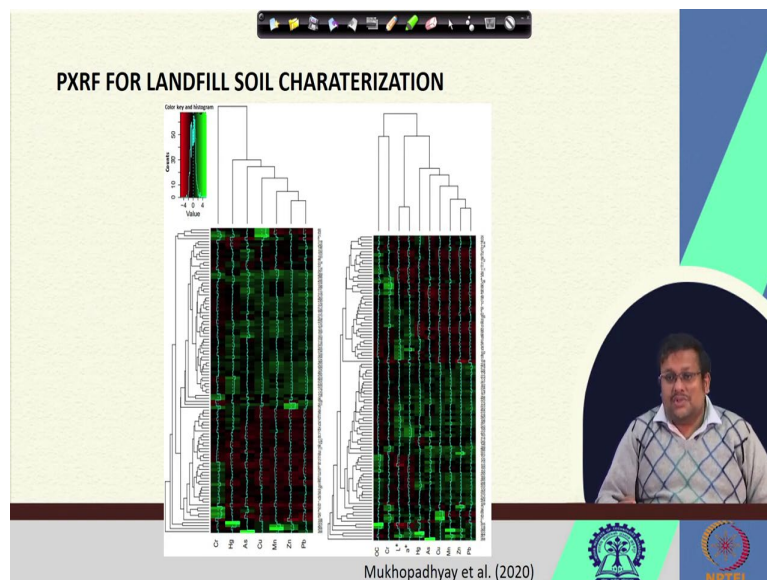
hotspots in the surrounding agricultural zones, so that we can judiciously plan our agricultural practice or we can avoid those pollution hotspots for subsequent growing of the crops.

So, you can see this is the soil type and these are composed of several debris, mainly they are composted products. So, after collecting the soil samples, we scanned them and then using the scan results, we utilize several clustering patterns, we have explored several clustering patterns. So, here you can see, we did the k medoids clustering and based on the solute width, in the solute plot, we have identified two clusters.

So, here it is cluster 1 and this is cluster 2. So, actually this was the area. So, this area was, this is the upper patch of the area, this is the lower patch of the area, and we have identified two clusters based on k medoids clustering. And identifying we have located that the samples from this cluster one belongs to this upper region and sample for cluster two belongs to this lower region.

And the reason behind this clustering is the relative closeness of this lower patch to the main active dump site. So, of course, the concentration was relatively higher, concentration of heavy metals were relatively higher in the lower patch than that of upper patch.

(Refer Slide Time: 28:01)



So, we have also utilized these the heatmap, these heat maps to identify the correlation between the different heavy metals and their relative clustering, which will be utilized, which can be very much informative for identifying their source. We call it bi-clustering heatmaps. So, using this bi-clustering heatmaps, we are able to cluster certain elements, and then which were actually helpful for identifying their source.

(Refer Slide Time: 28:43)



Mukhopadhyay et al. (2020)

Not only we have cluster, but we have also utilizing these PXRF elemental content, we produce the indicator kriging maps, where before all these seven elements where we can see the probability of exceeding its respective threshold values 0 to 25 percent probability, 25 to 50 percent probability. And in this way, we have identified this 75 to 100 percent probability, which were actually the lower patch as I have already told you.

So, that shows that our PXRF method was actually able to identify this pollution hotspot zones mainly in the lower patch for these heavy metals, so that we can using this indicator kriging. So, in this way it will be helpful for future planning of agricultural practices in this zone.

So, we should avoid growing the leafy vegetables in these zones, which are known as hyper accumulator of the heavy metals from the soil. So, this is how PXRF is helping the farmers to make, a policymaker to make the informed decision quickly and rapidly without wasting their time for traditional soil analysis.
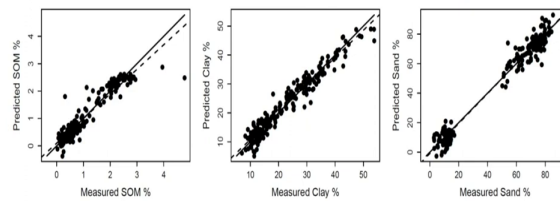
(Refer Slide Time: 30:16)



And also, a couple of years back we have seen that these PXRF and DRS when we combine their data together, they produce the synergistic results. So, these synergistic results is we call it sensor fusion or model fusion. And we got the 2 plus 2 equal to 5 synergistic effect because there are some kind of complementarity we have found in this two data set from PXRF and also DRS. And the application was useful for prediction of soil organic matter, clay content and sand content.

(Refer Slide Time: 31:01)



And the fusion of this PXRF and DRS was performed in different fashion. And we have secured a United State patent for these using the sensor fusion technology. We will talk about this sensor fusion more in our next lecture, and we will describe this thing in details. So,

guys, I hope that you have gained some good information in this lecture. There are certain other aspects of this sensor fusion which we are going to talk about in our next lecture.

(Refer Slide Time: 31:43)

## REFERENCES

- Aldabaa, A.A.A., D.C. Weindorf, S. Chakraborty, A. Sharma, and B. Li. 2015. Combination of proximal and remote sensing methods for rapid soil salinity quantification. Geoderma 239-240:34-46
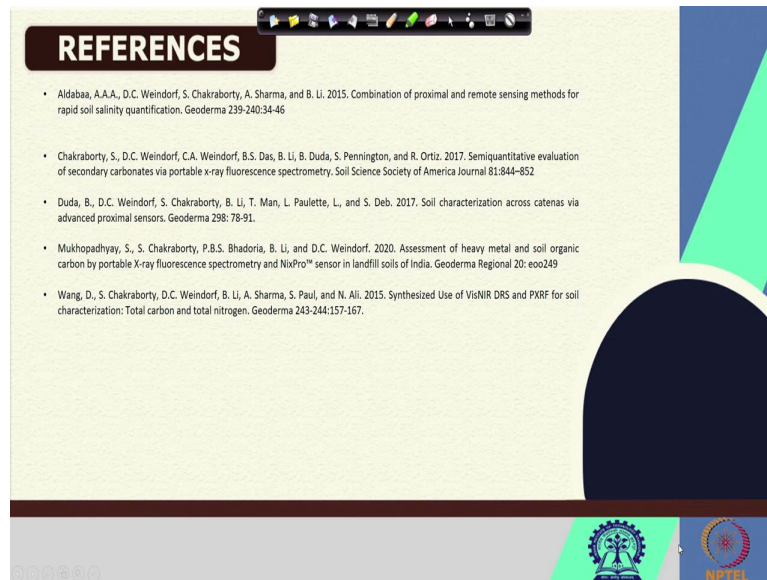
- Chakraborty, S., D.C. Weindorf, C.A. Weindorf, B.S. Das, B. Li, B. Duda, S. Pennington, and R. Ortiz. 2017. Semiquantitative evaluation of secondary carbonates via portable x-ray fluorescence spectrometry. Soil Science Society of America Journal 81:844–852

- Duda, B., D.C. Weindorf, S. Chakraborty, B. Li, T. Man, L. Paulette, L., and S. Deb. 2017. Soil characterization across catenas via advanced proximal sensors. Geoderma 298: 78-91.

- Mukhopadhyay, S., S. Chakraborty, P.B.S. Bhadoria, B. Li, and D.C. Weindorf. 2020. Assessment of heavy metal and soil organic carbon by portable X-ray fluorescence spectrometry and NixPro™ sensor in landfill soils of India. Geoderma Regional 20: eoo249

- Wang, D., S. Chakraborty, D.C. Weindorf, B. Li, A. Sharma, S. Paul, and N. Ali. 2015. Synthesized Use of VisNIR DRS and PXRF for soil characterization: Total carbon and total nitrogen. Geoderma 243-244:157-167.

And these are some of the references which I have used. And so, in a nutshell, we have seen how PXRF is useful for prediction of several soil properties and how the evolution of the PXRF that have happened from the simple statistical methods to complex statistical methods and how the expansion of PXRF application happened in different domains of the soil, starting from the routine soil parameters to land classification to pollution identification.

And also, we have seen the synthesis of the PXRF data with the DRS data and how they produce better results than using the individual sensor alone. We will start from here we will talk more about the sensor fusion. And then, I will show you another important portable proximal sensor called Nix sensor, which is also making a paradigm shift in the proximal soil analysis.

And I will also show you how we are using different machine learning approaches using the Nix sensor based results for better prediction of different soil properties rapidly and cost effectively. So, thank you guys. Let us meet in our next class. And we will start from here. We will talk more about this sensor fusion, and I will show you some application and then we will go to other sensors like Nix and so on so forth. Thank you and let us meet in our next lecture.