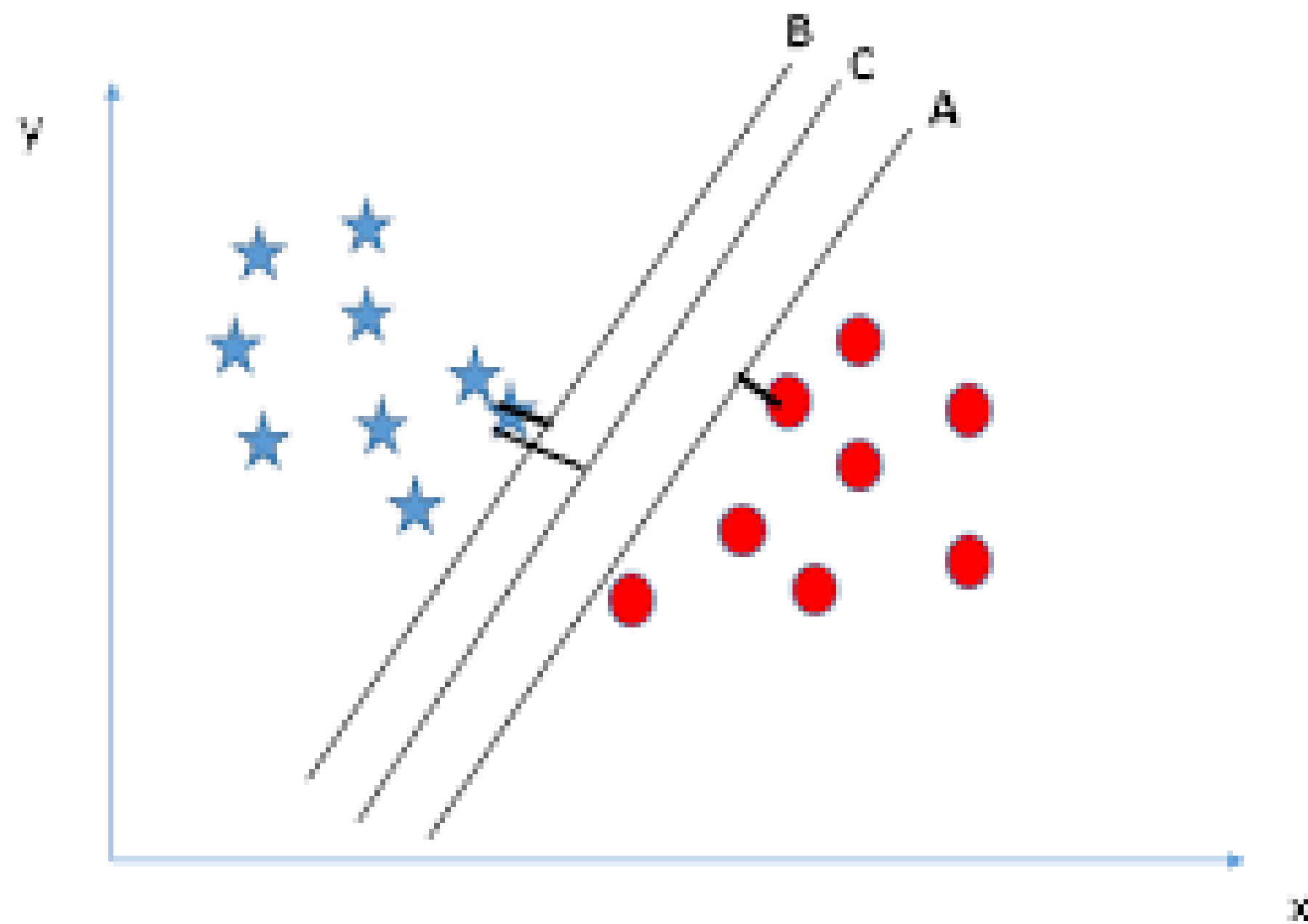


Practice 4

Support Vector Machine

Problem

- Predict whether income exceeds \$50K/yr
- Use linear SVM in mllib



- Use predefined function in `pyspark.mllib.classification`

Dataset

➤ Dataset description

- We encoded the features of data points like following

➤ 14 Statistic Features

1. >50K or <=50K
2. Age
3. workclass
4. Weight
5. Education
6. Education-num
...



➤ Preprocessed 123 features

1. >50K or <=50K
2. Age1
3. Age2
4. Age3
5. workclass1
...
124. native-country

❖ **The first column of the data matrix indicates the class labels.**

* UCI Machine Learning Repository :

<https://archive.ics.uci.edu/ml/datasets/adult>

➤ You can download the pre-processed train and test dataset on i-campus

Dataset

➤ We encoded original feature(Continuous, Categorical, etc) to have True/False value.

- 1. **Continuous feature:** Whether this value is in specific range

For example, if the value of age variable is bigger than 30 and smaller than 40, then the value of new range feature “30~40” becomes true.

- 2. **Categorical feature:** We made new features (for example, **A**, **B** and **C**) which are same with original category, and if one data had **A** category, then this data has **true** value in **A** feature but **false** value in **B** and **C** feature (OneHotEncoding).

	AGE	Categorical		10~20	20~30	30~40	A	B	C
Data 1	35	A	➡	False	False	True	True	False	False

Practice 4

1. Use predefined classes in *pyspark.mllib.classification* : *SVMwithSGD()*

Parameters for the method (default)

- *iterations = 100, step = 1.0, regParam = 0.01, regType = "l2"*

2. After training the models, calculate the F1 score, precision, recall for each label and accuracy

using test data points.

3. Due date: **May 7th 23:59**

<https://spark.apache.org/docs/latest/api/python/pyspark.mllib.html#pyspark.mllib.classification.SVMWithSGD>

<https://spark.apache.org/docs/latest/mllib-linear-methods.html>

Submission

➤ You need to submit *result.txt* file

- ✓ Write *f1 score, precision, recall* value of SVM result for *label 0*, **NOT** using predefined function but using *filter()* function
- ✓ Write *f1 score, precision, recall* value of SVM result for *label 1*, **NOT** using predefined function but using *filter()* function
- ✓ Write *accuracy* for *all labels*, using *TP(true positive), TN(true negative), FN(false negative), and FP(false positive)* values.

Label 0	Label 1
F1 Score : 0.8958	F1 Score : 0.5976
Precision : 0.8528	Precision : 0.7419
Recall : 0.9433	Recall : 0.5003
Accuracy : 0.8345	

Windows

```
Label 0
F1 Score : 0.8958
Precision : 0.8528
Recall : 0.9433
```

Linux

```
Label 1
F1 Score : 0.5976
Precision : 0.7419
Recall : 0.5003

Accuracy : 0.8345
```