



**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación



**UNIVERSITAT POLITÈCNICA
DE CATALUNYA
BARCELONATECH**



AI and Predictive Analytics in Data-Center Environments

Distributed Computing using Spark

SparkML (Hands On)

Josep Ll. Berral @BSC

Intel Academic Education Mindshare Initiative for AI

Hands-On: SparkML

- SparkML
 - Training models
 - Evaluate models
 - Use models for inference

Last remarks for SparkML

- Transforming “tabular” DataFrames to “libsvm” format
 - We use a “Vector Assembler”
 - `from pyspark.ml.feature import VectorAssembler`
 - `from pyspark.ml.linalg import Vectors`
 - `df = spark.read.csv("/home/vagrant/hs/ss13husa.csv", header = True, mode="DROPMALFORMED", inferSchema = True)`
 - `slice1 = df.select("SERIALNO", "PUMA", "DIVISION").limit(10)`
 - `assembler = VectorAssembler(inputCols = ["SERIALNO", "PUMA", "DIVISION"], outputCol = "features")`
 - `output = assembler.transform(slice1)`
 - `output.select("features").show()`

Summary

- Basic examples of SparkML
 - Train, evaluate and use machine learning models