

Diplomado Data Science - Machine Learning e Inteligencia Artificial

Temario

Hamdi Raissi, PhD Universidad de Lille, Francia, Profesor
Auxiliar PUCV

Patricio Videla, Profesor de planta PUCV, Jefe de docencia del Instituto de
Estadística

Mario Guzmán, Data
Scientist

Software: R, Python, Spark, SQL. Se incluye también 50USD de uso de la plataforma cloud de Amazon Web Services (AWS). No se necesita conocimientos previos de los softwares dado que una introducción será hecha para cada software ocupado. Los códigos listos para el uso y comentados en la clase.

Fechas: Presentación del diplomado 3 de diciembre a las 19hrs

Clases: 7, 10, 14, 17, 21, 28 de diciembre, 4, 7, 11, 14, 18, 21 de enero, 22 y 25 de febrero, 1, 4, 8, 11, 15, 18, 22, 25, 29 de marzo, 1, 5 y 7 de abril. Todas las clases son de 3 horas y empiezan a las 19hrs en modalidad “online”.

Charlas profesionales:

- A.** 20 de enero: Modelos VAR Aplicados en Economía (3 horas), Renata Abbott, Banco Central de Chile.
- B.** 31 de marzo: Aplicación de los modelos predictivos churn a recursos humanos (3 horas), Marleen van Kalmthout, Evalueserve.

Los conceptos presentados en clase serán cada vez ilustrados con datos reales o simulados.

Temas Básicos

1. Estadística descriptiva y introducción a R (2 horas)
 - a. Como utilizar R, funciones básicas, estrategias para elegir los paquetes R.
 - b. Estadísticas descriptivas y su visualización.
 - c. Tipos de variables en los datos.
2. Toma de decisión en un entorno aleatorio (2 horas)

- a. Test estadístico.
 - b. Intervalos de confianza para pronósticos.
3. Análisis de asociación de variables (5 horas)
- a. Estrategias para medir la correlación entre variables: Pearson, Spearman o Kendall?
 - b. Modelos lineales simples: Estimación MCO, Diagnóstico de bondad. Test de normalidad.
 - c. One way ANOVA y two way ANOVA, razón de correlación.
4. Reducción de la dimensión: Análisis por Componentes Principales (ACP) (3 horas)

Temas Avanzados

1. Modelos lineales múltiples (12 horas)
- a. Estimación MCO, diagnóstico de bondad (t-test, test de Fisher) y tipos de predicción (individual y del fenómeno estudiado).
 - b. Test de homogeneidad poblacional de Chow
 - c. Identificación de las variables pertinentes (C_p de Mallows, Criterios de información, algoritmos de selección forward, stepwise y backward). Como introducir las variables categóricas en un modelo lineal.
 - d. Problema de colinealidad y soluciones (regresión PCR, regresión Ridge)
 - e. Datos outliers (atípicos): detección y diagnóstico (leverages, residuos studentizados, distancia de Cook, DFBETAS). Solución con la estimación robusta de Theil-Sen y Siegel.
 - f. Heteroscedasticidad y autocorrelación: diagnóstico (test de Durbin Watson, tests de Breusch-Pagan) y estimación MCG.
2. Métodos numéricos de alto nivel computacional (6 horas)
- a. Introducción a EC2 de AWS.
 - b. Métodos bootstrap.
 - c. Experimentos de Monte Carlo.
3. Modelos para datos temporales (6 horas)
- a. Modelamiento univariado de datos temporales a través de series de tiempo AR, MA y ARMA.
 - b. Una caja de herramientas para el modelamiento de datos temporales:
 - i. Identificación: Autocorrelaciones (ACF), Autocorrelaciones parciales (PACF), Criterios de información
 - ii. Estimación: Menos Cuadrados Ordinarios (MCO), Máximo de verosimilitud

- iii. Diagnóstico
- c. Modelos SARIMA

4. Modelización de rendimientos financieros (6 horas)

- i. Hechos estilizados de las series de tiempo
 - a. Reagrupación de los valores extremos (Volatility clustering)
 - b. Leptocurticidad
 - c. Asimetría
- ii. Modelos GARCH y extensiones
- iii. Detección de la naturaleza financiera de datos dependientes
- iv. Medir los riesgos en finanza:
 - a. Valor en Riesgo (Value-at-Risk, VaR), VaR condicional
 - b. Técnicas bootstrap y Monte Carlo para mediciones de riesgos a horizonte más grande que uno
 - c. Backtesting de las medidas de riesgo
- v. Big data aplicada a la finanza: Uso de Elastic Cloud Computing (EC2) de AWS Amazon.

5. Introducción a SQL (3 horas)

- a. Modelos relacionales.
- b. Transformación de la información.
- c. Conexión con diferentes bases de datos.
- d. Depuración.
- e. Estudio de caso.

6. Introducción a Spark (3 horas)

- a. Tratamiento de data frame.
- b. Análisis descriptivo.
- c. Categorización de bases.
- d. Rutinas de Spark.

7. Algoritmo de K-medias (3 horas)

- a. Medidas de similaridades.
- b. Identificación del número de conglomerados.
- c. Métricas de validación.

8. Árboles de decisión (3 horas)

- a. Clasificación del árbol.
- b. Requisitos y supuestos de los datos.
- c. Interpretación de los resultados.

- d. Predicción y Evaluación.
 - e. Aplicación de un caso real en R.
9. Random Forest (3 horas)
- a. Introducción al Random Forest.
 - b. Entrenamiento de un modelo Random Forest.
 - c. Evaluación de out-of-bag error.
 - d. Evaluación del rendimiento del modelo Random Forest.
 - e. Estudio de caso en R.
10. Modelo de Regresión Logística (3 horas)
- a. Presentación del modelo e interpretación.
 - b. Validación de supuestos.
 - c. Ajuste del Modelo e interpretación de resultados.
 - d. Estudio de caso aplicado en R: Evaluación y Construcción.
11. Máquinas de vectores de soporte (3 horas)
- a. Definición de hiperplano de separación.
 - b. Clasificador de margen máximo.
 - c. SVM para clasificador linealmente separable.
 - d. SVM para clasificador linealmente no separable.
 - e. Extensión de las máquinas de vectores de soporte.
 - f. Métricas de validación.
12. Redes neuronales (3 horas)
- a. Arquitectura de una red.
 - b. Perceptrón.
 - c. Función de activación.
 - d. Back-propagation.
 - e. Métricas de validación.
13. Text mining (3 horas)
- a. Homologación de textos en base a cercanía de textos.
 - b. Arquitectura del web scraping.
 - c. Aplicaciones de web scraping y cercanía de textos.
12. Manejo de herramientas de AWS (6 horas)
- a. Introducción a S3.
 - b. Gestión de permisos con IAM.
 - c. Redes virtuales en la nube VPC.
 - d. Introducción a SageMaker.
 - e. Rutinas de modelos de ML en SageMaker.

13. Sistemas de recomendación (3 horas)
 - a. Filtros colaborativos
 - b. Sistema basado en usuarios e items.
 - c. Aplicaciones de sistemas de recomendación.