





Diplomado Introducción Analítica a la Ciencia de Datos

Coordinador: M. en C. Jaime Vázquez Alamilla.

Cuerpo docente

- 1. Actuario Eduardo Selim Martínez Mayorga. Profesor de Asignatura, Facultad de Ciencias, UNAM.
- 2. Maestro José Salvador Zamora Muñoz. Profesor de Asignatura, Facultad de Ciencias, UNAM.
- 3. Doctor Carlos Erwin Rodríguez Hernández-Vela. Investigador de Carrera, Instituto de Investigaciones en Matemáticas Aplicadas y en Sistemas, UNAM.
- 4. Doctor Jorge Ignacio González Cázares. Investigador de Carrera, Instituto de Investigaciones en Matemáticas Aplicadas y en Sistemas, UNAM.
- 5. Doctor Alan Riva Palacio Cohen. Investigador de Carrera, Instituto de Investigaciones en Matemáticas Aplicadas y en Sistemas, UNAM.
- 6.Doctora María Fernanda Gil Leyva Villa. Investigadora de Carrera, Instituto de Investigaciones en Matemáticas Aplicadas y en Sistemas, UNAM.
- 7.Doctora Claudia Ivonne Juárez Gallegos. Profesora de Asignatura, Facultad de Ciencias, UNAM.
- 8.Doctora Sandra Palau Calderón. Investigadora de Carrera, Instituto de Investigaciones en Matemáticas Aplicadas y en Sistemas, UNAM.

TEMARIO

MÓDULO 1. Programación para Ciencia de Datos (24 horas)

TEMA	Subtemas	No. de horas	Ponente(s)
I. Fundamentos de R	 R como calculadora Objetos básicos de R Sentencias de control if y else switch ifelse Ciclos Funciones adhoc 	3	Eduardo Selim Martínez Mayorga
II. Iteraciones en R	 Ciclos for Ciclos while Iteración la familia de funciones apply Iteración con librería {purrr} Funciones de mapeo Comparación con la familia apply Manejo de listas 	6	Eduardo Selim Martínez Mayorga
III. Tidyverse y flujo de ciencia de datos	 Importación Limpieza Transformación Modelación Visualización Comunicación 	9	Eduardo Selim Martínez Mayorga
IV. Sistemas de manejo de versiones	 Git GitHub GitLab 	3	Eduardo Selim Martínez Mayorga
V. Algunas integraciones populares	 Integración con python vía {reticulate} Nociones de programación en paralelo en R 	3	Eduardo Selim Martínez Mayorga

Módulo 2. Lectura y limpieza de datos (20 horas)

TEMA	Subtemas	No. de horas	Ponente(s)
I. Lectura de diversas fuentes de datos en R	 SQL y MariaDB noSQL: Mongo, Casandra, STATA, SAS, SPSS PDF, EXCEL XML JSON 	6	Eduardo Selim Martínez Mayorga
II. Datos web y web scraping	 Un centavo de HTML Librería {rvest} Selector Gadget de Google Chrome Introducción a expresiones regulares 	6	Eduardo Selim Martínez Mayorga
II. Principios de datos limpios	 tidyr Transformación y limpieza con dplyr Transformación y limpieza con data.table Transformación y limpieza con janitor 	8	Eduardo Selim Martínez Mayorga

Módulo 3. Análisis exploratorio de datos (26 horas)

TEMA		Subtemas	No. de horas	Ponente(s)
I. Medidas de tendencia central		 Media Media truncada Mediana Cuantiles 	3	Eduardo Selim Martínez Mayorga
II. Medidas de variación	1. 2.	Rango Rango intercuartil	3	Eduardo Selim Martínez Mayorga

	3.	Varianza y desviación estándar		
III. Medidas de asociación		 Covarianza y correlación Correlación de Spearman rankorder Relative risk y odds ratio V de Cramer Faceting 	6	Eduardo Selim Martínez Mayorga
IV. Datos categóricos	1. 2. 3.	Conteos Frecuencias relativas Tabla de contingencia	4	Eduardo Selim Martínez Mayorga
V. Análisis exploratorio de datos automatizado	1. 2. 3. 4.	dataMaid DataExplorer SmartEDA skimr	4	Eduardo Selim Martínez Mayorga
VI. Graficación y visualización de datos	1. 2. 3. 4.	ggplot2 Plotly Leaflet (mapas) Grafos	6	Eduardo Selim Martínez Mayorga

Módulo 4: Elementos de probabilidad y estadística para ciencia de datos (18 horas)

TEMA	Subtemas	No. de horas	Ponente(s)
1.Fundamentos de Probabilidad	 Fenómenos aleatorios, incertidumbre y variabilidad Fundamentos de probabilidad Variables aleatorias y sus distribuciones Límite de promedios de variables aleatorias 	6	Carlos Erwin Rodríguez Hernández-Vela
2. Fundamentos de Inferencia estadística	 De probabilidad a estadística Inferencia no paramétrica Inferencia clásica Inferencia bayesiana Ejemplos con R 	6	Carlos Erwin Rodríguez Hernández-Vela

3. Introducción a Modelos Lineales	 Modelo de regresión lineal simple Modelo de regresión lineal múltiple Ejemplos con R 	6	Carlos Erwin Rodríguez Hernández-Vela
---------------------------------------	--	---	--

Módulo 5: Aprendizaje supervisado (44 horas)

TEMA	Subtemas	No. de horas	Ponente(s)
1.Introducción al aprendizaje supervisado	 Datos en entrenamiento y prueba Validación cruzada y estimación del error de predicción Selección de modelos 	10	Jorge Ignacio González Cázares
2. Modelos lineales	 Regresión Categorical inputs Regresión regularizada (ridge, LASSO, etc.) Regresión no lineal Árboles de regresión Modelos lineales generalizados Clasificación Regresión logística Análisis de discriminante lineal Bayes ingenuo Árboles de clasificación 	12	Alan Riva Palacio Cohen
3. Árboles de decisión y bosques aleatorios	 Bagging Ensemble methods Boosting: adaboost & xgboost 	12	Jorge Ignacio González Cázares

4. Redes neuronales	 Redes superficiales y profundas Redes convolucionales y recurrentes Autoencoders y redes generativas Modelos de difusión Support vector machines 	10	Alan Riva Palacio Cohen
------------------------	--	----	-------------------------

Módulo 6. Aprendizaje no-supervisado (30 horas)

TEMA	Subtemas	No. de horas	Ponente(s)
I. Clustering	Determinación del número de clusters Vecinos más cercanos Clustering jerárquico	10	José Salvador Zamora Muñoz
II. Reducción de dimensionalidad	Descomposición SVD PCA (Componentes principales) t-SNE	10	José Salvador Zamora Muñoz
III. Análisis de factores	Análisis de factores exploratorio Análisis de factores confirmatorio	10	José Salvador Zamora Muñoz

Módulo 7: Optimización en aprendizaje de máquina (40 horas)

TEMA	Subtemas	No. de horas	Ponente(s)
1. Introducción a la optimización	 Modelos y funciones de pérdida Panorama de optimización Búsqueda en rejilla (grid search) Métodos de primer orden 	10	Sandra Palau Calderón

	5. Métodos de segundo orden		
2. Descenso del gradiente y sus variaciones	 Descenso del gradiente Descenso del gradiente estocástico Variaciones adaptativas Variaciones aleatorizadas 	10	María Fernanda Gil Leyva Villa
3. Optimización convexa y no- convexa	 Optimización de funciones convexas Optimización de funciones noconvexas Efecto de la inicialización en la optimización Efecto de la tasa de aprendizaje en la optimización Efecto de otros hiperparámetros en la optimización 	10	Sandra Palau Calderón
4. Otros paradigmas de optimización.	 Aprendizaje en línea Aprendizaje por refuerzo 	10	María Fernanda Gil Leyva Villa

Módulo 8: Comunicación de resultados (40 horas)

TEMA	Subtemas	No. de horas	Ponente(s)
1. Generación de documentos con RMarkdown	 Sintaxis básica de lenguaje de marcado ligero Generación de reportes HTML/PDF/WORD Generación de presentaciones Generación de notebooks de R 	6	Claudia Ivonne Juárez Gallegos
2. Generación de dashboard simples (flexdashboards)	 Introducción a flexdashboard Estructura general de un dashboard Personalización y diseño de dashboards 	8	Claudia Ivonne Juárez Gallegos

3. Documentos de larga extensión con bookdown	 Introducción a bookdown Estructura general de de un documento con bookdown Referencias y citas 	6	Claudia Ivonne Juárez Gallegos
4. Introducción a blogdown	 Introducción a blogdown Creación y personalización de sitio web Publicación del sitio web 	8	Claudia Ivonne Juárez Gallegos
5. Aplicaciones interactivas y dashboards interactivos con Shiny	 Introducción a Shiny Elementos para la visualización dinámica Publicación de la aplicación 	12	Claudia Ivonne Juárez Gallegos