



Facultad de  
Ciencias  
UNAM



## 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	<ol style="list-style-type: none"><li>1. R como calculadora</li><li>2. Objetos básicos de R</li><li>3. Sentencias de control</li><li>4. if y else</li><li>5. switch</li><li>6. ifelse</li><li>7. Ciclos</li><li>8. Funciones adhoc</li></ol>	3	Eduardo Selim Martínez Mayorga
II. Iteraciones en R	<ol style="list-style-type: none"><li>1. Ciclos for</li><li>2. Ciclos while</li><li>3. Iteración la familia de funciones apply</li><li>4. Iteración con librería {purrr}</li><li>5. Funciones de mapeo</li><li>6. Comparación con la familia apply</li><li>7. Manejo de listas</li></ol>	6	Eduardo Selim Martínez Mayorga
III. Tidyverse y flujo de ciencia de datos	<ol style="list-style-type: none"><li>1. Importación</li><li>2. Limpieza</li><li>3. Transformación</li><li>4. Modelación</li><li>5. Visualización</li><li>6. Comunicación</li></ol>	9	Eduardo Selim Martínez Mayorga
IV. Sistemas de manejo de versiones	<ol style="list-style-type: none"><li>1. Git</li><li>2. GitHub</li><li>3. GitLab</li></ol>	3	Eduardo Selim Martínez Mayorga
V. Algunas integraciones populares	<ol style="list-style-type: none"><li>1. Integración con python vía {reticulate}</li><li>2. Nociones de programación en paralelo en R</li></ol>	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	<ol style="list-style-type: none"><li>1. SQL y MariaDB</li><li>2. noSQL: Mongo, Casandra,</li><li>3. STATA, SAS, SPSS</li><li>4. PDF, EXCEL</li><li>5. XML</li><li>6. JSON</li></ol>	6	Eduardo Selim Martínez Mayorga
II. Datos web y web scraping	<ol style="list-style-type: none"><li>1. Un centavo de HTML</li><li>2. Librería {rvest}</li><li>3. Selector Gadget de Google Chrome</li><li>4. Introducción a expresiones regulares</li></ol>	6	Eduardo Selim Martínez Mayorga
II. Principios de datos limpios	<ol style="list-style-type: none"><li>1. tidy</li><li>2. Transformación y limpieza con dplyr</li><li>3. Transformación y limpieza con data.table</li><li>4. Transformación y limpieza con janitor</li></ol>	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	<ol style="list-style-type: none"><li>1. Media</li><li>2. Media truncada</li><li>3. Mediana</li><li>4. Cuantiles</li></ol>	3	Eduardo Selim Martínez Mayorga
II. Medidas de variación	<ol style="list-style-type: none"><li>1. Rango</li><li>2. Rango intercuartil</li></ol>	3	Eduardo Selim Martínez Mayorga

	3.	Varianza y desviación estándar		
III. Medidas de asociación	1. 2. 3. 4. 5.	Covarianza y correlación Correlación de Spearman rank-order 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	1. Fenómenos aleatorios, incertidumbre y variabilidad 2. Fundamentos de probabilidad 3. Variables aleatorias y sus distribuciones 4. Límite de promedios de variables aleatorias	6	Carlos Erwin Rodríguez Hernández-Vela
2. Fundamentos de Inferencia estadística	1. De probabilidad a estadística 2. Inferencia no paramétrica 3. Inferencia clásica 4. Inferencia bayesiana 5. Ejemplos con R	6	Carlos Erwin Rodríguez Hernández-Vela

3. Introducción a Modelos Lineales	<ol style="list-style-type: none"> <li>1. Modelo de regresión lineal simple</li> <li>2. Modelo de regresión lineal múltiple</li> <li>3. Ejemplos con R</li> </ol>	6	Carlos Erwin Rodríguez Hernández-Vela
------------------------------------	---	---	---------------------------------------

### Módulo 5: Aprendizaje supervisado (44 horas)

TEMA	Subtemas	No. de horas	Ponente(s)
1. Introducción al aprendizaje supervisado	<ol style="list-style-type: none"> <li>1. Datos en entrenamiento y prueba</li> <li>2. Validación cruzada y estimación del error de predicción</li> <li>3. Selección de modelos</li> </ol>	10	Jorge Ignacio González Cázares
2. Modelos lineales	<ol style="list-style-type: none"> <li>1. Regresión <ul style="list-style-type: none"> <li>• Categorical inputs</li> <li>• Regresión regularizada (ridge, LASSO, etc.)</li> <li>• Regresión no lineal</li> <li>• Árboles de regresión</li> <li>• Modelos lineales generalizados</li> </ul> </li> <li>2. Clasificación <ul style="list-style-type: none"> <li>• Regresión logística</li> <li>• Análisis de discriminante lineal</li> <li>• Bayes ingenuo</li> <li>• Árboles de clasificación</li> </ul> </li> </ol>	12	Alan Riva Palacio Cohen
3. Árboles de decisión y bosques aleatorios	<ol style="list-style-type: none"> <li>1. Bagging</li> <li>2. Ensemble methods</li> <li>3. Boosting: adaboost &amp; xgboost</li> </ol>	12	Jorge Ignacio González Cázares

4. Redes neuronales	<ol style="list-style-type: none"> <li>1. Redes superficiales y profundas</li> <li>2. Redes convolucionales y recurrentes</li> <li>3. Autoencoders y redes generativas</li> <li>4. Modelos de difusión</li> <li>5. Support vector machines</li> </ol>	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	<ol style="list-style-type: none"> <li>1. Modelos y funciones de pérdida</li> <li>2. Panorama de optimización</li> <li>3. Búsqueda en rejilla (grid search)</li> <li>4. Métodos de primer orden</li> </ol>	10	Sandra Palau Calderón

	5. Métodos de segundo orden		
2. Descenso del gradiente y sus variaciones	<ol style="list-style-type: none"> <li>1. Descenso del gradiente</li> <li>2. Descenso del gradiente estocástico</li> <li>3. Variaciones adaptativas</li> <li>4. Variaciones aleatorizadas</li> </ol>	10	María Fernanda Gil Leyva Villa
3. Optimización convexa y no-convexa	<ol style="list-style-type: none"> <li>1. Optimización de funciones convexas</li> <li>2. Optimización de funciones no-convexas</li> <li>3. Efecto de la inicialización en la optimización</li> <li>4. Efecto de la tasa de aprendizaje en la optimización</li> <li>5. Efecto de otros hiperparámetros en la optimización</li> </ol>	10	Sandra Palau Calderón
4. Otros paradigmas de optimización.	<ol style="list-style-type: none"> <li>1. Aprendizaje en línea</li> <li>2. Aprendizaje por refuerzo</li> </ol>	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	<ol style="list-style-type: none"> <li>1. Sintaxis básica de lenguaje de marcado ligero</li> <li>2. Generación de reportes HTML/PDF/WORD</li> <li>3. Generación de presentaciones</li> <li>4. Generación de notebooks de R</li> </ol>	6	Claudia Ivonne Juárez Gallegos
2. Generación de dashboard simples (flexdashboards)	<ol style="list-style-type: none"> <li>1. Introducción a flexdashboard</li> <li>2. Estructura general de un dashboard</li> <li>3. Personalización y diseño de dashboards</li> </ol>	8	Claudia Ivonne Juárez Gallegos

3. Documentos de larga extensión con bookdown	<ol style="list-style-type: none"> <li>1. Introducción a bookdown</li> <li>2. Estructura general de de un documento con bookdown</li> <li>3. Referencias y citas</li> </ol>	6	Claudia Ivonne Juárez Gallegos
4. Introducción a blogdown	<ol style="list-style-type: none"> <li>1. Introducción a blogdown</li> <li>2. Creación y personalización de sitio web</li> <li>3. Publicación del sitio web</li> </ol>	8	Claudia Ivonne Juárez Gallegos
5. Aplicaciones interactivas y dashboards interactivos con Shiny	<ol style="list-style-type: none"> <li>1. Introducción a Shiny</li> <li>2. Elementos para la visualización dinámica</li> <li>3. Publicación de la aplicación</li> </ol>	12	Claudia Ivonne Juárez Gallegos