

# Minería de Datos y Machine Learning

Módulo II - Aprendizaje Supervisado: Regresión

## software descripción



En los últimos años, los rápidos avances en las tecnologías de recopilación y almacenamiento de datos han dado lugar a conjuntos de datos que son "grandes" en muchos sentidos de la palabra. **La Minería de Datos** está relacionada con el análisis de estos grandes conjuntos de datos con la finalidad de proporcionar ideas, patrones, modelos descriptivos y predictivos que permitan extraer y generar conocimiento para las organizaciones. Algunas de las tareas más comunes son la clasificación, agrupamiento, descubrimiento de reglas de asociación y patrones de respuestas, detección de anomalías, etc.

Para poder realizar estas tareas, la Minería de Datos combina los aspectos teóricos fundamentales de áreas tan diversas como la estadística, las matemáticas, la ingeniería y las ciencias de la información con muchas aplicaciones prácticas y relevantes de la vida real. Dentro de esas disciplinas destaca el aprendizaje automatizado o **machine learning**, que está teniendo profundos efectos en muchas industrias diferentes, desde servicios financieros hasta ventas minoristas y publicidad; convirtiéndose rápidamente en una herramienta fundamental para tomar mejores decisiones en los negocios: decisiones basadas en datos, no en instintos o conjeturas.

Este módulo del curso tiene como objetivo presentar los fundamentos del aprendizaje estadístico, discutir estrategias para evaluar la eficiencia predictiva de los modelos supervisados y presentar los principales modelos de regresión.

## publico objetivo

Analistas de datos. Profesionales relacionados a la inteligencia de negocios, investigación de mercados e interesados en el área de Ciencia de Datos. Académicos e investigadores. . Público en general que requiera extraer conocimiento desde diferentes fuentes de información.

## logros de aprendizaje

Al finalizar este curso, el participante conocerá los fundamentos del aprendizaje supervisado y aplicará modelos de regresión para inferir y predecir. De manera específica el participante estará en capacidad de:

- Conocer las diferentes técnicas y modelos del aprendizaje estadístico y automatizado, que forman parte del proceso de aprendizaje supervisado y que son esenciales para la minería de datos.
- Aplicar computacionalmente los modelos aprendidos en diferentes campos como la industria, el comercio, la banca, los seguros, biología, etc.
- Comprender el desarrollo de los algoritmos de las principales técnicas de regresión.
- Presentar de manera efectiva los resultados obtenidos.

## contacto



+51970227123



info@perustat.com



<http://www.perustat.com>



<fb://perustat>



# contenidos

<b>Tema 1</b>	<b>Fundamentos del Aprendizaje Supervisado</b>	(8 horas)
	<ul style="list-style-type: none"><li>• Aprendizaje estadístico<ul style="list-style-type: none"><li>- Ideas generales</li><li>- Aprendizaje estadístico vs. Machine Learning.</li><li>- Aprendizaje estadístico vs. Ciencia de Datos.</li><li>- Principales problemas.</li></ul></li><li>• Predicción e Inferencia</li><li>• Modelos paramétricos y no paramétricos.</li><li>• Precisión vs. interpretación: Métodos inflexibles y métodos flexibles.</li><li>• Compromiso sesgo-varianza.</li><li>• Revisión de conceptos y notación básica<ul style="list-style-type: none"><li>- Vectores aleatorios.</li><li>- Momentos de una distribución.</li><li>- Valor esperado.</li><li>- Matriz de varianza-covarianza y correlaciones.</li><li>- Combinaciones lineales.</li><li>- Distribución normal multivariada.</li></ul></li></ul>	
<b>Tema 2</b>	<b>Regresión Lineal Múltiple</b>	(8 horas)
	<ul style="list-style-type: none"><li>• Definición y notación.</li><li>• Modelo lineal clásico.</li><li>• Estimación: Mínimos Cuadrados Ordinarios y Máxima Verosimilitud.</li><li>• Inferencia.</li><li>• Evaluación y selección de modelos.</li><li>• Evaluación de supuestos y análisis de residuos.</li><li>• Predicción.</li><li>• Predictores cualitativos.</li></ul>	
<b>Tema 3</b>	<b>Evaluación de la performance</b> predictiva	(8 horas)
	<ul style="list-style-type: none"><li>• Error de entrenamiento vs. error de prueba.</li><li>• Funciones de pérdida.</li><li>• Particionamiento del conjunto de datos.</li><li>• Técnicas de evaluación por remuestreo<ul style="list-style-type: none"><li>- LOOCV (Leave One Out cross-validation)</li><li>- KCV (k-fold cross-validation)</li><li>- Bootstrap.</li></ul></li><li>• Hiperparámetros.</li><li>• Selección de modelos.</li></ul>	
<b>Tema 4</b>	<b>Selección de Variables</b>	(4 horas)
	<ul style="list-style-type: none"><li>• Selección de Subconjuntos<ul style="list-style-type: none"><li>- Selección de los mejores subconjuntos</li><li>- Selección Stepwise.</li><li>- Elección del modelo óptimo.</li></ul></li><li>• Métodos de regularización<ul style="list-style-type: none"><li>- Regresión Ridge</li><li>- Regresión Lasso</li></ul></li></ul>	
<b>Tema 5</b>	<b>Extensiones al modelo</b> lineal	(4 horas)
	<ul style="list-style-type: none"><li>• Regresión Polinomial.</li><li>• Regresión usando Splines.</li><li>• Regresión local.</li><li>• GAMs (Modelos Aditivos Lineales)</li></ul>	

## contacto



+51970227123



info@perustat.com



<http://www.perustat.com>



fb://perustat



## metodología

La metodología del curso se basa en una combinación de clases teóricas y análisis de casos prácticos en la computadora, con la finalidad de que el participante comprenda la metodología, la motivación, los supuestos, las fortalezas y las debilidades de los métodos tratados en el curso. Cada sección del curso está motivada por un conjunto de datos en particular, de tal forma que el participante gane experiencia trabajando con una amplia variedad de fuentes de datos similares a los que usa en la realidad. Los contenidos están estructurados en 8 sesiones con un total de 32 horas académicas (24 horas cronológicas).

## evaluación y asistencia

- **Certificado de asistencia:** Para obtener este certificado debe de mantener un porcentaje mínimo de 75% de asistencia a las clases.
- **Certificado de aprobación:** El 100% de la calificación final se obtiene sobre la base de cuestionarios o listas de ejercicios al final de cada sesión. Para recibir el certificado de aprobación el participantes deben obtener al menos un 70% de los puntos posibles y contar con el porcentaje mínimo de asistencia a clases.

Serán otorgados certificados a nombre de PeruStat Analytics S.A.C. que acredita a los participantes del curso. La certificación que se otorga es excluyente.

## materiales

Material preparado por el equipo de capacitación con los contenidos del curso el cuál será entregado a los participantes en medios físicos o digitales.

## referencias

- Bishop, C. (2011). Pattern Recognition and Machine Learning. Springer.
- Clarke, B., Fokoue, E. y Zhang, H. (2009). Principles and Theory for Data Mining and Machine Learning. Springer Verlag.
- Gareth, J., Witten, D., Hastie, T. y Tibshirani, R. (2013). An Introduction to Statistical Learning: with Applications in R. Springer.
- Hastie, T., Tibshirani, R. y Friedman. J. (2016). The Elements of Statistical Learning: Data Mining, Inference, and Prediction. Springer.
- Kuhn, M. y Johnson, K. (2013). Applied Predictive Modeling. Springer Verlag.
- Lantz, B. (2015). Machine Learning with R: Expert techniques for predictive modeling to solve all your data analysis problems Packt Publishing; 2 edition .
- Larose, D.T. (2006). Data Mining Methods and Models. Wiley Interscience.
- Ledolter, J. (2013). Data Mining and Business Analytics with R. John Wiley & Sons.
- Nisbet, R., Elder IV, J. y Miner, G. (2009). Handbook of Statistical Analysis and Data Mining Applications. Academic Press.
- Putler, D. y Krider, R. E. (2012). Customer and Business Analytics: Applied Data Mining for Business Decision Making Using R. Chapman and Hall/CRC.
- Shmuelil, G., Bruce, P., Yahav, I., Patel, N. y Lichtendahl, K. (2017). Data Mining for Business Analytics: Concepts, Techniques, and Applications in R. John Wiley & Sons.
- Véliz, C. (2018). Aprendizaje Automático. Análisis para la Minería de Datos y Big Data. Departamento Académico de Ciencias. Pontificia Universidad Católica del Perú.
- Wu J. y Coggeshall, S. (2012). Foundations of Predictive Analytics. Chapman and Hall/CRC.

### contacto



+51970227123



info@perustat.com



<http://www.perustat.com>



fb://perustat

