

MASTER		Master en Data Science para Finanzas	
ASIGNATURA		<i>Data Science para la Gestión de Información No Estructurada</i>	
Nº de ECTS		2	
Nº de horas docentes		Parte II: 15h (2 ECTS, 10 sesiones)	
Nº de horas actividades académicas dirigidas			
Profesor responsable de la asignatura		Francisco José Izquierdo Catalán	
Cuatrimestre		2º Cuatrimestre	

### 1.- DESCRIPCIÓN GENERAL DE LA ASIGNATURA Y OBJETIVOS DE DOCENCIA:

El objetivo de esta asignatura es introducir a los alumnos en las necesidades que aparecen en el entorno financiero con la gestión de información en lenguaje natural, así como la introducción a las diferentes herramientas y técnicas de minería de textos. Se cubrirán de forma teórica y práctica aspectos desde la captura y procesamiento de diferentes fuentes, la clasificación de la información, extracción de entidades y finalmente la visualización y explotación de la misma. Para cada uno de los diferentes puntos se verá la forma de abordarlos mediante programación. Así mismo se realizará una instalación de una herramienta de Minería de Textos de Código Abierto.

Durante las horas docentes se cubrirán los siguientes puntos:

- Introducción a la Gestión de Información No Estructurada y Minería de Textos
- Pre-procesamiento de la información y extracción de entidades
- Modelos de almacenamiento y representación de textos
- Clasificación
- Clusterización y modelado de tópicos
- Análisis semántico. Análisis de sentimientos
- UIMA y herramientas de Minería de Textos de Código Abierto.
- Chatbots y asistentes virtuales
- Deep Learning en NLP

### 2.- FORMA DE EVALUACIÓN PREVISTA:

La evaluación de la asignatura se compondrá en un 40% de la entrega y evaluación de las actividades académicas dirigidas, y en un 60% del resultado del examen final de la asignatura. Para aprobar la asignatura **será imprescindible obtener al menos un 5 en el examen final**; las actividades académicas dirigidas no serán reevaluables. La asistencia a clase es obligatoria, admitiéndose hasta un 20% de ausencias sin justificación; será criterio del profesor admitir o no la justificación; una asistencia menor del 80% supondrá la pérdida del derecho a examen en convocatoria ordinaria. Durante el curso se propondrán 4 actividades académicas dirigidas. Estas actividades se propondrán al final de la clase semanal, y se deberán entregar al final de la segunda semana posterior, es decir si la actividad se propone un miércoles, el alumno contará con los 4 días de esa semana y dos semanas adicionales. **El plazo de entrega no es prorrogable.**

Actividades académicas dirigidas	40%
Prueba objetiva final	60%

PROGRAMA DETALLADO		
Nº de sesión	Detalle del contenido docente: temas, casos prácticos, actividades académicas dirigidas que se verán en dicha sesión,...	Lecturas recomendadas o referencias bibliográficas relativas a los conceptos-temas desarrollados en la sesión
1	<p>Sesión 1: Introducción a la Minería de Textos</p> <ul style="list-style-type: none"> <li>• Historia</li> <li>• Objetivos</li> <li>• Campos de aplicación de la Minería de Textos</li> <li>• Aplicación de Minería de Textos en Economía y Finanzas</li> </ul>	<ol style="list-style-type: none"> <li>1. Ronen Feldman y James Sanger (2007). <i>The Text Mining Handbook, Advanced Approaches in Analyzing Unstructured Data</i>. Cambridge University Press</li> <li>2. Presentación y notas de clase</li> </ol>
2	<p>Sesión 2: Pre-procesamiento y extracción de entidades</p> <ul style="list-style-type: none"> <li>• Pre-procesamiento y normalización</li> <li>• Detección de idioma</li> <li>• Tokenización</li> <li>• Lematización</li> <li>• Análisis gramatical</li> <li>• Extracción de Entidades</li> </ul>	<ol style="list-style-type: none"> <li>1. Sarkar D. (2019). <i>Text Analytics with Python</i>. Apress, Berkeley, CA</li> <li>2. Presentación y notas de clase</li> </ol>
3	<p>Sesión 3: Modelos de almacenamiento y representación de textos</p> <ul style="list-style-type: none"> <li>• Bag of Words</li> <li>• Bag of N-grams</li> <li>• El modelo TF-IDF</li> <li>• Word2Vec</li> <li>• GloVe</li> <li>• FastText</li> </ul>	<ol style="list-style-type: none"> <li>1. Sarkar D. (2019). <i>Text Analytics with Python</i>. Apress, Berkeley, CA</li> <li>2. Presentación y notas de clase</li> </ol>
4	<p>Sesión 4: Clasificación</p> <ul style="list-style-type: none"> <li>• Medidas de frecuencias</li> <li>• Ley de Zipf</li> <li>• Concepto de Clasificación</li> <li>• Clasificación Bayesiana</li> <li>• Regresión Logística</li> <li>• Support Vector Machine</li> <li>• Random Forest</li> </ul>	<ol style="list-style-type: none"> <li>1. Gareth James, Daniela Witten, Trevor Hastie y Robert Tibshirani (2013). <i>An Introduction to Statistical Learning with applications in R</i>. Springer</li> <li>2. Ted Kwartler (2017). <i>Text Mining in Practice with R</i>. Wiley</li> <li>3. Ashish Kumar y Avinash Paul (2016). <i>Mastering Text Mining with R</i>. Packt Publishing Ltd</li> </ol>

		<ol style="list-style-type: none"> <li>Sarkar D. (2019). Text Analytics with Python. Apress, Berkeley, CA</li> <li>Presentación y notas de clase</li> </ol>
5	<p>Sesión 5: Clustering y modelado de tópicos</p> <ul style="list-style-type: none"> <li>Latent Dirichlet Allocation</li> <li>Latent Semantic Analysis</li> <li>Medidas de Similitud</li> <li>K-means</li> </ul>	<ol style="list-style-type: none"> <li>Ronen Feldman y James Sanger (2007). <i>The Text Mining Handbook, Advanced Approaches in Analyzing Unstructured Data</i>. Cambridge University Press</li> <li>Ashish Kumar y Avinash Paul (2016). <i>Mastering Text Mining with R</i>. Packt Publishing Ltd</li> <li>Sarkar D. (2019). Text Analytics with Python. Apress, Berkeley, CA</li> <li>Presentación y notas de clase</li> </ol>
6	<p>Sesión 6: Análisis semántico y análisis de Sentimientos</p> <ul style="list-style-type: none"> <li>Conceptos</li> <li>Relaciones semánticas</li> <li>Desambiguación</li> <li>Representación del análisis semántico</li> <li>Modelos de análisis de sentimiento no supervisados</li> <li>Modelos supervisados</li> </ul>	<ol style="list-style-type: none"> <li>Ashish Kumar y Avinash Paul (2016). <i>Mastering Text Mining with R</i>. Packt Publishing Ltd. 2016</li> <li>Sarkar D. (2019). Text Analytics with Python. Apress, Berkeley, CA</li> <li>Presentación y notas de clase</li> </ol>
7	<p>Sesión 7: UIMA y herramientas de código abierto</p> <ul style="list-style-type: none"> <li>Arquitectura UIMA</li> <li>Soluciones de código abierto</li> <li>Instalación</li> <li>Carga de documentación</li> <li>Búsqueda</li> <li>Otras funcionalidades</li> </ul>	<ol style="list-style-type: none"> <li>Página de documentación del Proyecto UIMA: <a href="https://uima.apache.org/documentation.html">https://uima.apache.org/documentation.html</a></li> <li>Apache SolR Tutorial: <a href="http://lucene.apache.org/solr/guide/7_4/solr-tutorial.html">http://lucene.apache.org/solr/guide/7_4/solr-tutorial.html</a></li> <li>Rafał Kuć (2013). Apache SolR 4 Cookbook. Packt Publishing Ltd</li> <li>Presentación y notas de clase</li> </ol>
8	<p>Sesión 8: Chatbots y asistentes virtuales</p> <ul style="list-style-type: none"> <li>Concepto de asistente virtual</li> <li>Intenciones, entidades y diálogos</li> <li>Diálogos complejos</li> </ul>	<ol style="list-style-type: none"> <li>Chiara Martino (2019). Conversation Design Workflow: How to design your chatbot in 10(basic) steps: <a href="https://medium.com/@chiara.martino/conversation-design-workflow-how-to-design-your-chatbot-in-10-basic-steps-721652b056d">https://medium.com/@chiara.martino/conversation-design-workflow-how-to-design-your-chatbot-in-10-basic-steps-721652b056d</a></li> </ol>

	<ul style="list-style-type: none"> <li>• Integraciones</li> </ul>	2. Presentación y notas de clase
9	Sesión 9: Deep learning en NLP <ul style="list-style-type: none"> <li>• Word-embedding</li> <li>• Text Classification en DL</li> </ul>	1. Sarkar D. (2019). <i>Text Analytics with Python</i> . Apress, Berkeley, CA 2. Presentación y notas de clase
10	Sesión 10: Revisión general anterior al examen. Resolución de dudas	

INFORMACION ADICIONAL	
Bibliografía básica	<ol style="list-style-type: none"> <li>1. Bholat et al (2015). <i>Text Mining for Central Banks</i>. CCBS, Bank of England</li> <li>2. Sarkar D. (2019). <i>Text Analytics with Python</i>. Apress, Berkeley, CA</li> <li>3. Manning, Raghavan, and Schütze (2009). <i>An Introduction to Information Retrieval</i>. Cambridge University Press</li> <li>4. Murphy (2012). <i>Machine Learning: a Probabilistic Perspective</i>. MIT Press</li> <li>5. Hércules Antonio do Prado y Edilson Ferneda (2008). <i>Emerging Technologies of Text Mining: Techniques and Applications</i>. Information Science Reference</li> <li>6. Héctor Cuesta (2013). <i>Practical Data Analysis</i>. Packt Publishing Ltd</li> <li>7. Michael W. Berry y Malu Castellanos (2008). <i>Survey of Text Mining: Clustering, Classification and Retrieval</i>. Springer</li> </ol>
Bibliografía Complementaria	<ol style="list-style-type: none"> <li>1. Scott R. Baker, Nicholas Bloom y Steven J. Davis. <i>Measuring Economic Policy Uncertainty</i> (Enero de 2013). Chicago Booth Research Paper No. 13-02. Disponible en SSRN: <a href="http://ssrn.com/abstract=2198490">http://ssrn.com/abstract=2198490</a> or <a href="http://dx.doi.org/10.2139/ssrn.2198490">http://dx.doi.org/10.2139/ssrn.2198490</a></li> <li>2. Paul C. Tetlock. <i>Giving Content to Investor Sentiment: The Role of Media in the Stock Market</i>. <i>Journal of Finance</i>, Forthcoming. Disponible en SSRN: <a href="http://ssrn.com/abstract=685145">http://ssrn.com/abstract=685145</a> o <a href="http://dx.doi.org/10.2139/ssrn.685145">http://dx.doi.org/10.2139/ssrn.685145</a></li> </ol>

	<p>3. Tim Loughran y Bill McDonald. <i>When is a Liability not a Liability? Textual Analysis, Dictionaries, and 10-Ks</i> (Marzo de 2010). <i>Journal of Finance</i>, Forthcoming. Disponible en SSRN: <a href="http://ssrn.com/abstract=1331573">http://ssrn.com/abstract=1331573</a></p> <p>4. Stephen Hansen,, Michael McMahon y Andrea Prat. (2014) <i>Transparency and deliberation within the FOMC: a computational linguistics approach</i>. CFM discussion paper series, CFM-DP2014-11. Centre For Macroeconomics, London, UK.</p> <p>5. G. Hoberg, y G. M. Phillips. (), '<i>Product Market Synergies and Competition in Mergers and Acquisitions: A Text-Based Analysis</i>', 2010. <i>The Review of Financial Studies</i>, Vol. 23, No. 10, pages 3773-3811.</p> <p>6. Scott Hendry y Alison Madeley. <i>Text Mining and the Information Content of Bank of Canada Communications</i> (Noviembre, 2010). Disponible en SSRN: <a href="http://ssrn.com/abstract=1722829">http://ssrn.com/abstract=1722829</a> o <a href="http://dx.doi.org/10.2139/ssrn.1722829">http://dx.doi.org/10.2139/ssrn.1722829</a></p>
<b>Actividades Complementarias</b>	
<b>Localización del profesor</b>	francisco.izquierdo@cunef.edu