

### PROGRAMA DE CURSO

Código	Nombre			
IN 5528	INTRODUCCIÓN A BIG DATA			
Nombre en Inglés				
Introduction to Big Data				
SCT	Unidades Docentes	Horas de Cátedra	Horas Docencia Auxiliar	Horas de Trabajo Personal
6	10	3,0	1,5	5,5
Requisitos			Carácter del Curso	
IN3501 Tecnologías de Información y Comunicaciones para la Gestión o CC3201 Bases de Datos y IN4901 Práctica Profesional I o AUTOR			Electivo de la carrera Ingeniería Civil Industrial	
Competencias a la que tributa el curso				
<p>Competencias de Egreso</p> <ul style="list-style-type: none"> <li>• Concebir soluciones a los problemas que surgen en las organizaciones, utilizando los conocimientos provenientes de las tecnologías de información y comunicaciones.</li> <li>• Empezar e innovar en el desarrollo de soluciones a problemas de ingeniería, demostrando iniciativa y capacidad de toma de decisión.</li> <li>• Comunicar ideas y resultados de trabajos profesionales o de investigación, en forma escrita y oral.</li> <li>• Gestionar su auto-aprendizaje en el desarrollo del conocimiento de su profesión, adaptándose a los cambios del entorno.</li> </ul>				
Propósito del Curso				
<p>Hoy en día, se están obteniendo grandes cantidades de datos tanto estructurados como no estructurados, ante esto, nace la necesidad de entregar a los alumnos de ingeniería industrial conocimientos relacionados con Big Data para que puedan concebir nuevas soluciones a problemas que puedan surgir en todo tipo de organizaciones. Para llevar a cabo lo anterior, el cuerpo docente impartirá tanto clases teóricas como prácticas, además de diversas actividades que realizarán los estudiantes durante el curso haciéndolos agentes activos durante la realización del curso.</p>				
Resultados de Aprendizaje				
<p>Al finalizar el curso, el estudiante es capaz de utilizar la arquitectura Lambda para poder construir una solución a un problema de Big Data. Los conocimientos adquiridos para este fin son los siguientes:</p> <ul style="list-style-type: none"> <li>• Principios de los sistemas de datos y dar una visión general de la Arquitectura Lambda.</li> <li>• <i>Batch Layer</i> de la arquitectura Lambda. Se aprenderá sobre el modelamiento de un <i>master dataset</i>, utilizando procesamiento por <i>batch</i> para crear vistas arbitrarias de sus datos, así como las ventajas y desventajas entre el procesamiento gradual y por lotes.</li> <li>• <i>Serving Layer</i>, la que proporciona baja latencia de acceso a las vistas producidas por la <i>Batch Layer</i>. Se aprenderá acerca de las bases de datos especializadas que sólo se escriben de forma masiva y que éstas son dramáticamente más simples que las tradicionales,</li> </ul>				

dándoles un excelente rendimiento, funcionamiento, y las propiedades de robustez.

- *Speed Layer*, la que compensa la alta latencia de la *Batch Layer* para poder entregar resultados actualizados para las consultas. Se aprenderá acerca de las bases de datos NoSQL y procesamiento en stream.
- Procesamiento en *batch* incrementales, las variantes de la arquitectura básica Lambda, y cómo obtener el máximo provecho de sus recursos.

Metodología Docente	Evaluación General
<p>Este curso tiene una connotación teórico-práctica. Está compuesto por cátedras y auxiliares, las cuales condujeran a la construcción de una solución de <i>Big Data</i>. Se aplicará la teoría a la resolución de casos simples, que permitirán al alumno prepararse para desarrollar con éxito las distintas actividades del curso. Esta metodología tiene como objetivo primario entregar la práctica necesaria para llevar con éxito el diseño e implementación de las tareas a realizar durante el curso. En este sentido, se realizarán controles para evaluar el aprendizaje logrado en el curso, además de otras actividades complementarias como lecturas o presentaciones sobre temáticas relevantes para el curso.</p>	<p>El curso consiste de 2 notas, tareas (NT) y controles (NC). El cálculo de esas notas se efectúa de la siguiente forma:</p> <ul style="list-style-type: none"> <li>• <math>NC = \text{Promedio de controles } (\sum Ci)/n</math>, donde <math>Ci</math> son las notas de los controles.</li> <li>• <math>NT = \text{Promedio de las entregas parciales } (\sum wi*Pi)/n</math>, donde <math>Pi</math> son las notas de las tareas y <math>wi</math> la ponderación que tiene cada una de ellas.</li> <li>• El alumno puede eximirse de dar el examen si el promedio actual de controles es mayor o igual a 5.5 y la nota de tareas (NT) es mayor o igual que 5.5. En este caso, la nota final corresponde al promedio simple entre NC y NT</li> <li>• En caso de que el alumno rinda el examen, la nota final se calcula de la siguiente forma. <math>(0,6*NC+0,4*EX)*0,5+NP*0,5</math></li> <li>• La condición para aprobar el curso es: <math>NP \geq 4.0</math> y <math>NC \geq 4.0</math></li> </ul>

### UNIDADES TEMÁTICAS

Número	Nombre de la Unidad	Duración en Semanas
1	Introducción a Big Data	2
Contenidos	Indicador de Logro	Referencias a la Bibliografía
1. Un nuevo paradigma para Big Data <ol style="list-style-type: none"> <li>Escalando con bases de datos tradicionales.</li> <li>Bases de datos NoSQL.</li> <li>Primeros principios.</li> <li>¿Qué es lo que deseamos en un sistema Big Data?</li> <li>Problemas con arquitecturas completamente incrementales.</li> <li>Arquitectura Lambda.</li> <li>Tendencias en tecnología.</li> </ol>	Aprendizaje de los principios de los sistemas de datos y dar una visión general de la Arquitectura Lambda.	1,2,3

Número	Nombre de la Unidad	Duración en Semanas
2	<i>Batch Layer</i>	6
Contenidos	Indicador de Logro	Referencias a la Bibliografía
1. Modelo de datos para Big Data <ol style="list-style-type: none"> <li>Las propiedades de los datos.</li> <li>Modelo basado en hechos para representar la data.</li> <li>Esquemas gráficos.</li> </ol> 2. Almacenamiento de datos en la <i>Batch Layer</i> <ol style="list-style-type: none"> <li>Requisitos de almacenamientos para el <i>dataset</i> maestro.</li> <li>Escogiendo una solución de almacenamiento.</li> <li>Como funciona un sistema de archivos distribuidos.</li> <li>Almacenando un <i>dataset</i> maestro en un sistema de</li> </ol>	Aprendizaje sobre el modelamiento de un <i>master dataset</i> , utilizando procesamiento por <i>batch</i> para crear vistas arbitrarias de sus datos, así como las ventajas y desventajas entre el procesamiento gradual y por lotes.	1,2,4

<p>archivos distribuidos.</p> <p>3. <i>Batch Layer</i></p> <ol style="list-style-type: none"> <li>Computación en la <i>Batch Layer</i></li> <li>Algoritmos recomputados vs algoritmos incrementales.</li> <li>Escalabilidad en la <i>Batch Layer</i>.</li> <li>MapReduce: un paradigma para Big Data.</li> </ol>		
--	--	--

Número	Nombre de la Unidad	Duración en Semanas
3	<i>Serving Layer</i>	2
Contenidos	Indicador de Logro	Referencias a la Bibliografía
<p>1. <i>Serving Layer</i></p> <ol style="list-style-type: none"> <li>Métricas de rendimiento para la <i>Serving Layer</i>.</li> <li>Solución para la el problema de normalización/desnormalización.</li> <li>Requisitos para una base de datos <i>Serving Layer</i>.</li> <li>Contrastando con una solución completamente incremental.</li> </ol>	<p>Aprendizaje acerca de las bases de datos especializadas que sólo se escriben de forma masiva y que éstas son dramáticamente más simples que las tradicionales, dándoles un excelente rendimiento, funcionamiento, y las propiedades de robustez.</p>	1,2,5

Número	Nombre de la Unidad	Duración en Semanas
4	<i>Speed Layer</i>	5
Contenidos	Indicador de Logro	Referencias a la Bibliografía
<p>1. Vistas en tiempo real</p> <ol style="list-style-type: none"> <li>Computando en vistas en tiempo real.</li> <li>Almacenando vistas en tiempo real.</li> </ol>	<p>Aprendizaje de las bases de datos NoSQL y procesamiento en stream. Además de procesamiento en <i>batch</i> incrementales, las variantes</p>	1,2,5

<ul style="list-style-type: none"> <li>c. Desafíos de la computación incremental.</li> <li>d. Actualizaciones asíncronas vs síncronas.</li> </ul> <p>2. Procesamiento en <i>stream</i> y <i>queuing</i>.</p> <ul style="list-style-type: none"> <li>a. <i>Queuing</i>.</li> <li>b. Procesamiento en <i>stream</i>.</li> </ul> <p>3. Procesamiento en <i>stream</i> con <i>micro-batch</i></p> <ul style="list-style-type: none"> <li>a. Logrando la semántica exactamente de una vez.</li> <li>b. Conceptos fundamentales.</li> <li>c. Extendiendo diagrama para procesamiento en <i>micro-batch</i>.</li> <li>d. Terminando la <i>Speed Layer</i>.</li> </ul> <p>4. Arquitectura Lambda en profundidad.</p> <ul style="list-style-type: none"> <li>a. Definición de un sistema de datos.</li> <li>b. <i>Batch</i> y <i>Serving Layers</i>.</li> <li>c. <i>Speed Layer</i>.</li> <li>d. <i>Query Layer</i>.</li> </ul>	<p>de la arquitectura básica Lambda, y cómo obtener el máximo provecho de sus recursos.</p>	
--	---	--

Bibliografía General	
<ol style="list-style-type: none"> <li>1. Nathan Marz y James Warren. "Big Data: Principles and best practices of scalable realtime data systems". Manning 2015.</li> <li>2. Foster Provost y Tom Fawcett. "Data Science for Business: What you need to know about data mining and data-analytic thinking". 2013</li> <li>3. Bernard Marr. "Big Data: Using SMART Big Data, Analytics and Metrics To Make Better Decisions and Improve Performance". 2015</li> <li>4. Tom White. "Hadoop: The Definitive Guide". O'Reilly 2015</li> <li>5. Mat Brown. "Learning Apache Cassandra". 2015</li> </ol>	

Vigencia desde:	Otoño 2016
Elaborado por:	Juan Velásquez y Felipe Vildoso
Validado por:	
Revisado por:	Comisión de Docencia DII