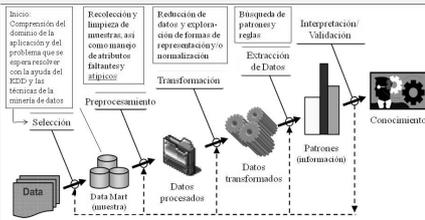


### 1.1 ANÁLISIS PREDICTIVO



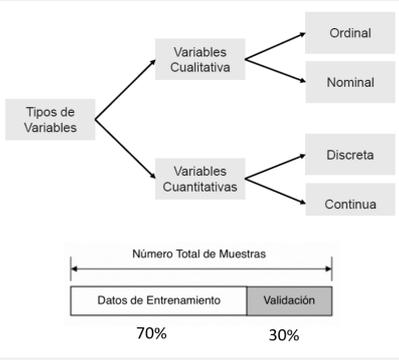
El análisis predictivo consiste en la tecnología que aprende de la experiencia para predecir el futuro comportamiento de individuos para tomar mejores decisiones. INFORMS Institute for Operations Research and the Management Sciences

### 1.4 PROCESO KDD



Proceso de Extracción del Conocimiento conocido como Knowledge Discovery in Databases KDD

### 1.5 KDD - Selección de los Datos

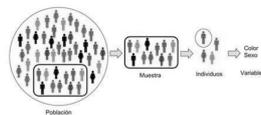


Exploración de los datos : Base Original -> Tablón de Datos (Minable)

Muestra de entrenamiento : Un subconjunto para entrenar un Modelo

Muestra de validación : Un subconjunto para probar el modelo entrenado

### 2.2 CONCEPTOS ESTADÍSTICA DESCRIPTIVA



La población, es el conjunto total de objetos o personas de interés en un estudio. *una característica relevante es todos sus elementos deben cumplir con un conjunto predefinido de características.*

La Muestra es el subconjunto de la población, la cual se utiliza para estudiar las características de la población en general. Estas deben ser: Aleatorias y representativas.

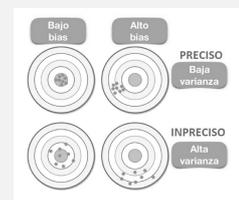
Variable Aleatoria cualquier característica que tome dos o más valores en una población.

### 3 PREDICIÓN LINEAL

### MCO Minimo cuadrado Ordinario

### 5.2 EVALUACION DE MODELOS BINARIOS

### 1.2 BIAS O SESGO Y VARIANZA



#### ¿Qué es el Bias?

El bias o sesgo puede ser pensado como un modelo que no ha tenido en cuenta toda la información disponible en el set de datos, lo que dificulta predicciones precisas.

#### ¿Qué es la varianza?

La Varianza es una medida de dispersión que se utiliza para representar la variabilidad de un conjunto de datos respecto de la media aritmética de los mismo.

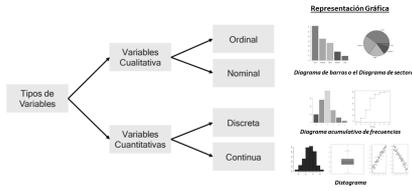
### 1.6 KDD -Preprocesamiento

- Análisis Descriptivo Univariado** Evaluar mediana de tendencia central ( Mediana, Moda), de dispersión. Los datos null se replazan por 0
- Análisis Descriptivo Multivariado** Análisis de Correlaciones, Gráficos de Dispersión, Etc.
- Análisis Descriptivo Temporal (de Proporciones)** Debe contar con estabilidad temporal de proporciones.
- Análisis Descriptivo Temporal (de Predicción)** Busca predecir el comportamiento de una variable en particular.

### MODELOS KDD - SEMMA - CRISP

Modelo	Descripción	Características	Aplicaciones
SEMMA	Selección, Preparación, Construcción, Evaluación	Procesamiento de datos	Modelado predictivo
CRISP	Selección de objetivos, Análisis de datos, Preparación de datos, Construcción de modelos, Evaluación de modelos, Implementación de modelos	Modelado predictivo	Modelado predictivo

### 2.3 TIPOS DE VARIABLES ALEATORIAS



#### Variable Cualitativa

**Ordinal** presenta modalidades no numéricas, en las que existe un orden por Ej: notas en un examen.

**Nominal** presenta modalidades no numéricas que no admiten un criterio de orden. Ej: El estado civil, con las siguientes modalidades: sol, cas, sep, divor y viudo.

#### Variable Cuantitativa

**Discreta** puede asumir un número contable de valores: Ej: # hijos en una familia.

**Continua** puede asumir un número incontable de valores. Ej: peso de una persona.

### 4 REGRESIÓN LOGÍSTICA

**Regresión Logística**

La regresión logística es una herramienta estadística y de minería de datos que predice la probabilidad de un resultado que sólo puede tener dos valores: los dicotómicos: Sí/No, Verdadero/Falso, etc. La predicción se basa en el uso de una función predictora lineal y un modelo de clasificación de máxima verosimilitud. El análisis de regresión logística se encuentra en el conjunto de Modelos Lineales Generalizados (GLM) por sus raíces en inglés que usa como función de enlace la función logit, la cual se aplica a continuación:

Por lo tanto, se debe representar los siguientes conceptos:

- La función logit que es un derivado de la función sigmoide.
- Las variables independientes pueden ser cuantitativas o cualitativas, en cualquier caso éstas se deben transformar a variables dummy.
- La ecuación modelo es  $Y = 1 / (1 + e^{-X\beta})$ .

**Propiedades de la Función:**

Se observa que la regresión logística produce una curva sigmoide, que se limita a valores entre 0 y 1. La regresión logística se utiliza para predecir la probabilidad de que un suceso ocurra en función de un conjunto de variables. Por lo tanto, los resultados se expresan como una probabilidad. El valor de la probabilidad se puede interpretar como la probabilidad de que ocurra un suceso.

### 1.3 SOBREAJUSTE E INFRAJUSTE

	Underfitting	Just right	Overfitting
<b>Symptoms</b>	<ul style="list-style-type: none"> <li>High training error</li> <li>Training error close to test error</li> <li>High bias</li> </ul>	<ul style="list-style-type: none"> <li>Training error slightly lower than test error</li> </ul>	<ul style="list-style-type: none"> <li>Very low training error</li> <li>Training error much lower than test error</li> <li>High variance</li> </ul>
<b>Regression illustration</b>			
<b>Classification illustration</b>			
<b>Deep learning illustration</b>			
<b>Possible remedies</b>	<ul style="list-style-type: none"> <li>Complexify model</li> <li>Add more features</li> <li>Train longer</li> </ul>		<ul style="list-style-type: none"> <li>Perform regularization</li> <li>Get more data</li> </ul>

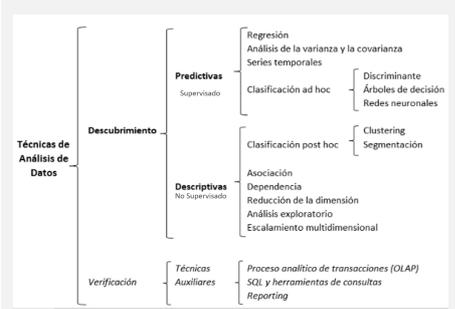
#### Overfitting

Se produce cuando un modelo modela demasiado bien los datos de entrenamiento., por lo que no es capaz de generalizar, y cuando le lleguen nuevos datos obtendrá pésimos resultados.

#### Underfitting

Se produce cuando nuestro modelo no es capaz de identificar patrones. Por lo que tendrá siempre pésimos resultados.

### 1.7 KDD - Minería de Datos



### Tipos de aprendizaje de datos

