Aplicación de Inteligencia Artificial en avalúos masivos.









XI Simposio CPCI

Cancun, Mx, 5 al 7 de setiembre de 2018

Mgter. Juan Pablo Carranza

Jefe de Modelización y Métodos Estadísticos Proyecto de Estudio Territorial Inmobiliario Gobierno de la Provincia de Córdoba, Argentina.









¿En qué consiste el aprendizaje estadístico?

Algoritmo: Población: Muestra: Conjunto de reglas matemáticas Sólo que nos ayudan a identificar las conocemos los Valores de la relaciones entre Inputs y "Inputs". variable de interés Outputs en la muestra. (Output) Características urbanas que, a priori, tienen Predicción: relación con la Sobre los Inputs conocidos en la población se variable de interés aplican las reglas definidas por el algoritmo para (Inputs) predecir el Output.



Dirección General de CATASTRO Ministerio de FINANZAS

de AS CORDOBA PROTOES

Juan Pablo Carranza

Abordaje de la estadística clásica.

Se impone una forma funcional al problema de estudio:

$$y = b_0 + b_1 X_1 + ... + b_n X_n + x$$

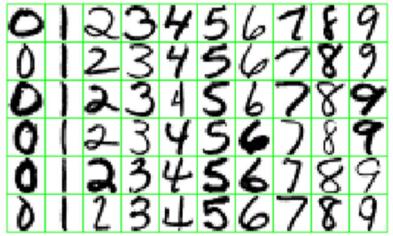
Abordaje algorítmico, aprendizaje estadístico.

Se respeta la estructura de la información, la no-linealidad, las propiedades emergentes propias de fenómenos caóticos.



La [IA] cada vez más a nuestro alrededor... algunas aplicaciones:

Ejemplo clásico! Digitalizando el mundo físico.



Tomado de Hastie, Tibshirani, Friedman: "The Elements of Statistical Learning"... Recomiendo!!!

En la actualidad? Foco en la interacción con el ser humano.



Tomado de lizuka (et. al.): "Globally and Locally Consistent Image Completion".



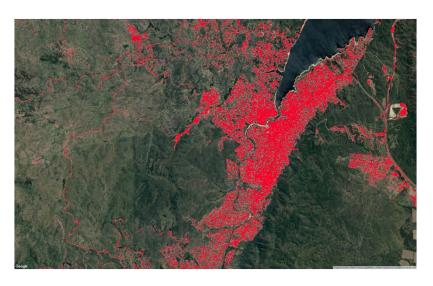
Dirección General de CATASTRO Ministerio de FINANZAS



Juan Pablo Carranza

La [IA] cada vez más a nuestro alrededor... acercándonos al estudio territorial:





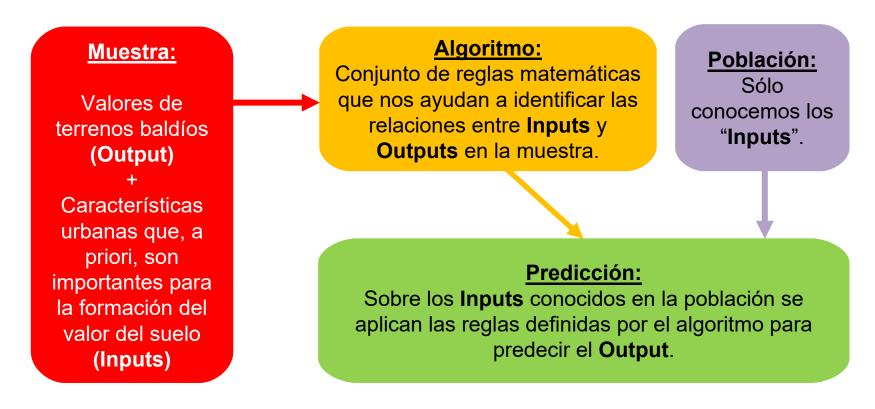
En base a muestras, el algoritmo puede clasificar las áreas construidas en una ciudad







¿Cómo se aplica al estudio del valor del suelo?





Dirección General de CATASTRO | Ministerio de FINANZAS

de 'AS CORDOBA PROTOES

Juan Pablo Carranza

¿Cuál será nuestro "output"? No todas las muestras de mercado son iguales!

Recurrimos a la **econometría espacial** para homogeneizar valores:

$$log(y) = b_0 + b_1 X + b_2 W y + b_3 W u + h$$

Diferenciando la expresión con respecto a X:

$$b_1 @ +gy/y) / gX$$



¿Cuáles serán nuestros "inputs"?

Catastrales de entorno:

densidad construida en el entorno.

disponibilidad de baldíos en el entorno.

tamaño promedio del lote en el entorno.

Dinámica inmobiliaria en el entorno.

Distancias:

al centro.

a vías principales.

a vías secundarias.

a zonas de bajo perfil inmobiliario

a zonas de alto perfil inmobiliario

al río.

a vías de FFCC.

a la ruta.

a zonas de depreciación.

etc...

Satelitales de entorno:

área construida

área no construida

dimensión fractal

índices de fragmentación









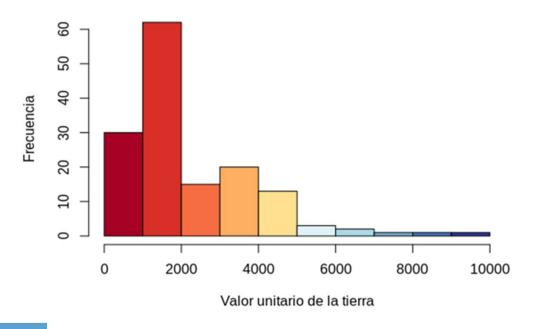
Terminando de conformar nuestra base de datos... ¿Cómo sabemos que tenemos datos de buena calidad?

Estadística clásica: Eliminar outliers.

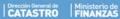
¡NO!

El problema no sigue una distribución normal e imponer esa condición limita el análisis.

Histograma del valor del suelo en San Francisco, Córdoba











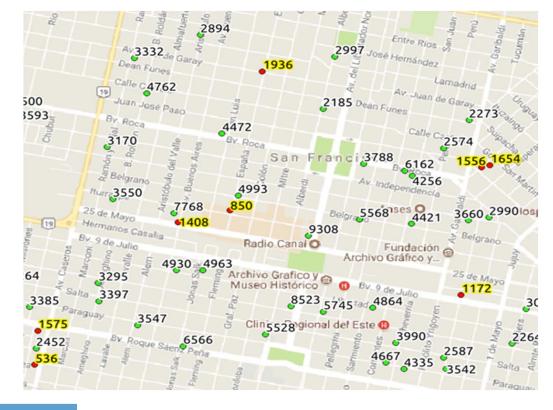
Terminando de conformar nuestra base de datos... ¿Cómo sabemos que tenemos datos de buena

calidad?

Sí eliminamos Outliers Espaciales (o inliers) mediante el **índice de Moran local**

$$I = rac{rac{N}{S_0} \sum_i \sum_j W_{ij} Z_i Z_j}{\sum_i Z_i^2}$$

Outlier espacial: dato atípico en su entorno.





Pero... ¿qué es un algoritmo? Por ejemplo: Un árbol de regresión

$$R_1(j,s) = \{X | X_j \le s\} \text{ and } R_2(j,s) = \{X | X_j > s\}$$

$$\min_{j,s} \left[\min_{c_1} \sum_{x_i \in R_1(j,s)} (y_i - c_1)^2 + \min_{c_2} \sum_{x_i \in R_2(j,s)} (y_i - c_2)^2 \right]$$

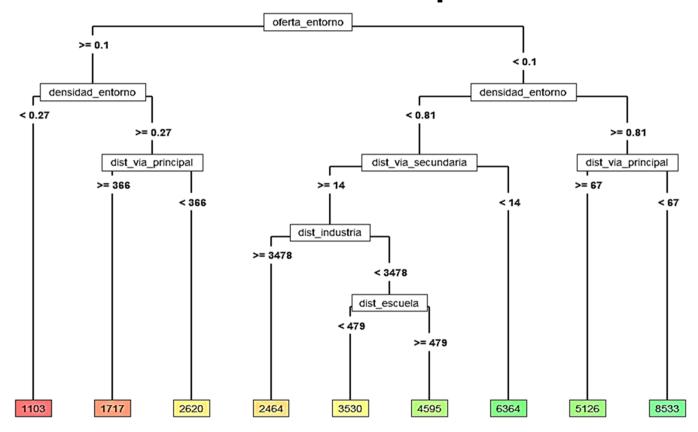
$$\hat{c}_1 = \text{ave}(y_i | x_i \in R_1(j,s)) \text{ and } \hat{c}_2 = \text{ave}(y_i | x_i \in R_2(j,s))$$

Esto es un árbol de regresión! Su misión es separar grupos de terrenos baldíos en grupos lo más homogéneos posibles y lo más heterogéneos con el resto.





Visualmente es más fácil de comprender...





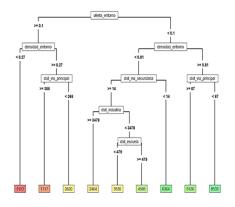
Dirección General de CATASTRO Ministerio de FINANZAS

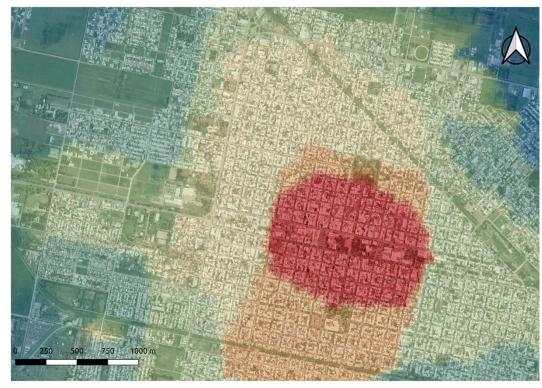


Juan Pablo Carranza

Not such an intelligent intelligence! Damn it! **Problema: Overfitting.**

La predicción del árbol de regresión devuelve unas pocas zonas con valor homogéneo







Dirección General de CATASTRO | Ministerio de FINANZAS



¿Cómo solucionamos el Overfitting del árbol de regresión? Inventamos muchos árboles para intentar simular las situaciones no observadas en la muestra.

Pero... ¿Cómo hacemos si ya hemos usado todos los datos y todos los inputs? # Usamos menos datos (con reposición) en cada árbol (bootstrap) # Usamos menos inputs en cada nodo.

¿Y cómo unificamos las predicciones de todos estos árboles? Hay muchas formas. La más sencilla: promediamos las predicciones de cada uno de los árboles para cada uno de los puntos a predecir.

Haciendo esto aplicamos una técnica llamada: RANDOM FOREST



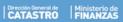
Dirección General de CATASTRO | Ministerio de FINANZAS

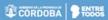


Aplicando Random **Forest** pasamos de esta situación....



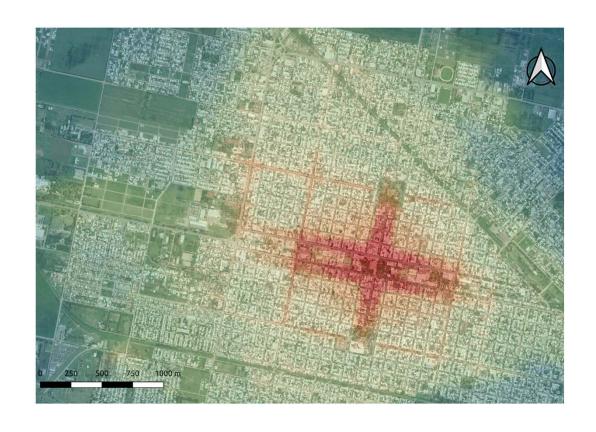






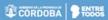
... A una estimación mucho más consistente.

Random Forest generaliza mejor.









Ventajas de Random Forest

- ☐ Reduce la varianza en la estimación (es decir, es menos propenso al overfitting).
- ☐ Logra una calidad predictiva mucho más elevada. Generaliza mejor!

Desventajas de Random Forest

☐ Ya no hay sólo un árbol que nos permita comprender esquemáticamente cómo se conforma el valor del suelo en función de los inputs utilizados.

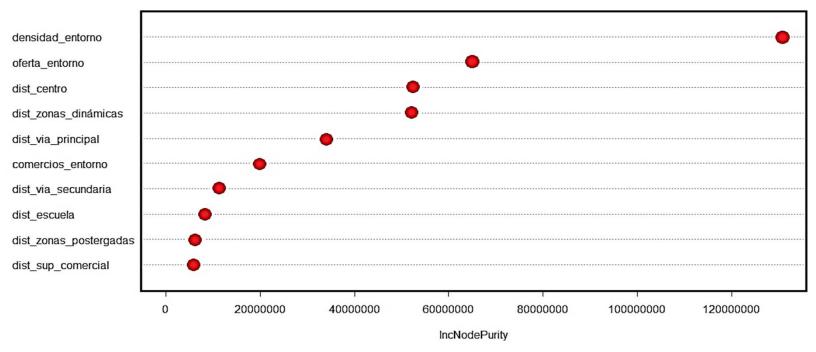
Su utilización dependerá de las características del problema de investigación.



Ministerio de FINANZAS CORDOBA STODOS

Sin embargo, sí podemos conocer con mucha precisión cuál es la importancia relativa de cada input en la predicción:

Importancia relativa de variables









¿Cómo medimos el error de predicción de nuestras estimaciones?

Validación Cruzada	Grupos									
Acción:	1	2	3	4	5	6	7	8	9	10
Sale el grupo 1	1	2	3	4	5	6	7	8	9	10
	1	Se es	tima e	el mod	lelo c	on los	dato	s de l	os gr	upos
Se predice el gru	po 1	y se n	nide e	l erro	de p	redico	ión <	Į		
Sale el grupo 2	1	2	3	4	5	6	7	8	9	10
Sale el grupo 3	1	2	3	4	5	6	7	8	9	10
Sale el grupo 4	1	2	3	4	5	6	7	8	9	10
Sale el grupo 5	1	2	3	4	5	6	7	8	9	10
Sale el grupo 6	1	2	3	4	5	6	7	8	9	10
Sale el grupo 7	1	2	3	4	5	6	7	8	9	10
Sale el grupo 8	1	2	3	4	5	6	7	8	9	10
Sale el grupo 9	1	2	3	4	5	6	7	8	9	10
		2	3	4	5	6	7	8	9	10

Se calculan diferentes medidas de error...

☐ Error relativo promedio en valor absoluto:

$$ERP = \frac{\sum_{i=1}^{n} \left(\frac{|\widehat{y_i} - y_i|}{y_i} \right)}{n}$$

☐ Coeficiente de variación (IAAO):

$$CV = \frac{\sum_{i=1}^{n} \left(\frac{|\widehat{y_i} - y_i|}{y_i} - \frac{\sum_{i=1}^{n} \left(\frac{|\widehat{y_i} - y_i|}{y_i} \right)}{n} \right)}{\sum_{i=1}^{n} \left(\frac{|\widehat{y_i} - y_i|}{y_i} \right)}$$



Evaluación comparativa de diferentes modelos aplicados (caso Ciudad de San Francisco):





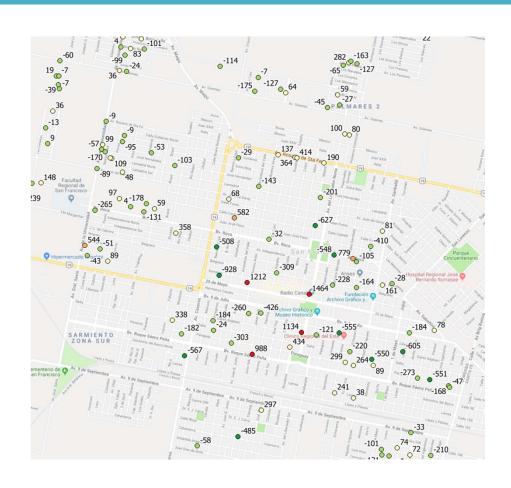




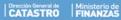
Incorporación de la Auto-Correlación espacial:

La estimación realizada mediante Random Forest tiene un error intrínseco. Para reducir su impacto, se realiza un Kriging de los errores y esta nueva estimación se suma a la estimación original.

Estimación final = Estimación Random Forest + Kriging Ordinario del error.





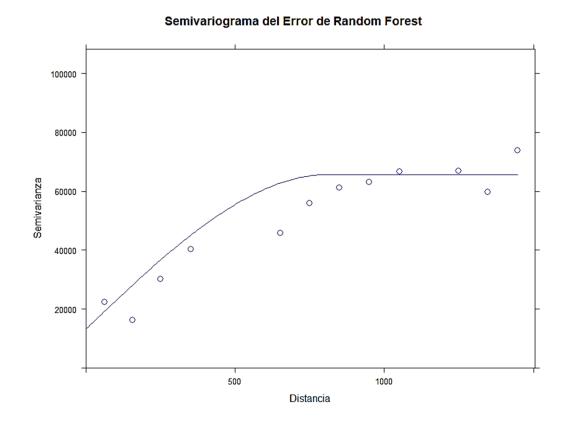




Interpolación del ajuste del error:

A los fines de ajustar el error de la predicción inicial se utiliza la dependencia espacial de las observaciones, capturadas mediante el semivariograma.

Ley de Tobler: Todas las cosas están relacionadas entre sí, pero las cosas más próximas en el espacio tienen una relación mayor que las distantes.





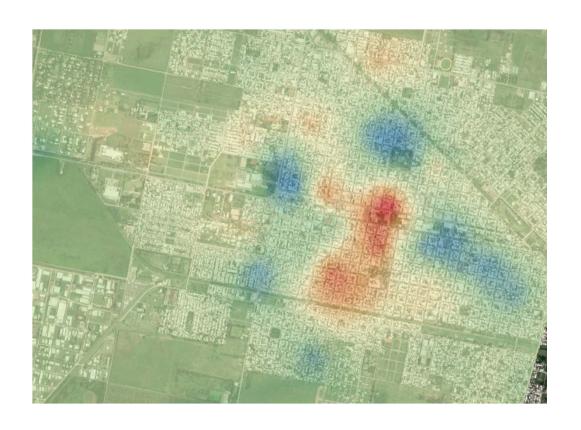




Interpolación del ajuste del residuo:

En rojo: Ajusta la predicción hacia arriba.

En azul: Ajusta la predicción hacia abajo.









Estimación final Random Forest + **Kriging ordinario:**

A la predicción original se le suma el Kriging del Error.











Evaluación comparativa de diferentes modelos aplicados (caso Ciudad de San Francisco):





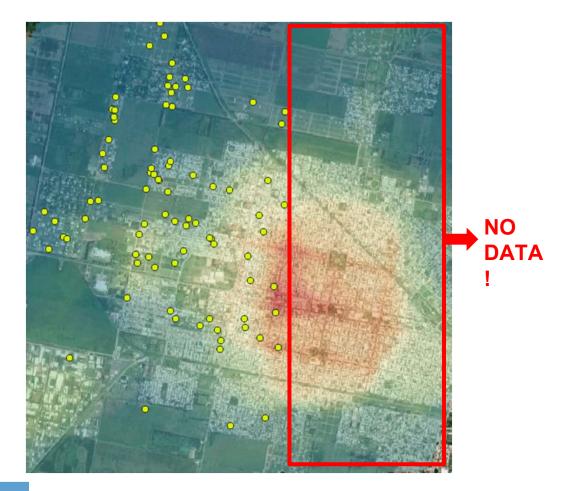


Potencia del aprendizaje estadístico:

Capacidad de generalización!

(menos costoso y más preciso)

Ejemplo: Estimación realizada con datos sólo en la mitad oeste de la ciudad.





Dirección General de CATASTRO Ministerio de FINANZAS

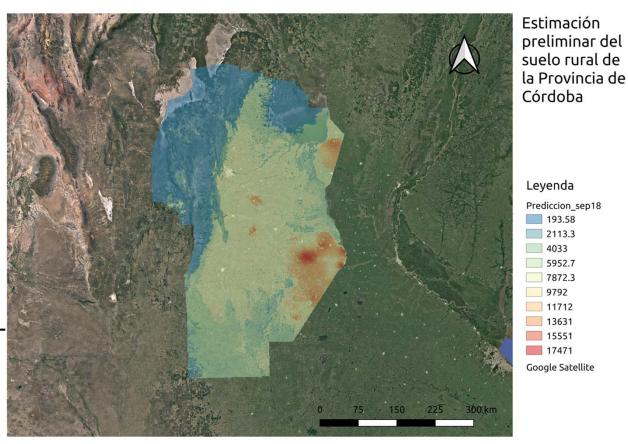


Avances en la estimación del suelo rural:

Métodos más adecuados:

Supported Vector Machine (error 18%)

Random Forest (error 20%-31%)











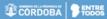
Variables utilizadas en la estimación....

TIPO	VARIABLE	TIPO	VARIABLE			
Suelo	Capacidad de Uso de la Tierra		Distancia a asentamientos urbanos			
	Indice de Productividad		Distancia a red vial pavimentada			
	Coberturas/Usos de la tierra		Distancia a red de Energía Eléctrica			
	Pendiente		Distancia a red de Gas Natural			
	Altura		Distancia a centros de acopio			
	Suceptibilidad a inundación y/o anegamiento	Infraestructura	Distancia a Balanzas Públicas			
	Suceptibilidad a erosión eólica		Distancia a puerto San Lorenzo			
	Deficiencia de Húmedad		Acceso a riego			
	NDVI		Acceso a riego complentario			
Hídrología	Distancia a cursos de agua		Distancia a obras hídricas			
	Disponibilidad de agua subterránea		Producción Tambera			
	Profundidad del nivel freático					
Climáticos	Precipitación media anual	Estructura productiva	Producción Ganadera			
	Régimen de Temperaturas		Actividad turística			
	Radiación Solar		Explotación minera			
	Bioclimaticas	Económicas	Rendimientos zonales por localidad			
	Vulnerabilidad de sequía	Economicas	Arrendamientos zonales por localidad			



CATASTR

Ministerio de FINANZAS



Juan Pablo Carranza





