

Selección de candidatos informativos en datos astronómicos

Javiera Astudillo

Profesores: Karim Pichara (DCC, PUC); Pavlos Protopapas (IACS, Harvard)

1. Aprendizaje de Máquinas
2. Algoritmos y técnicas (algunos)
3. Trabajo de Investigación

Aprendizaje de Máquinas

¿Qué es?

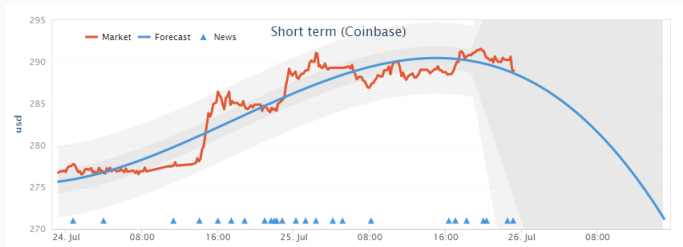
Aprender de los datos.

- Haciendo uso de herramientas de probabilidades y estadística.
- Sin indicarle qué aprender (descubrimientos de patrones no vistos).

Ejemplo: clasificación



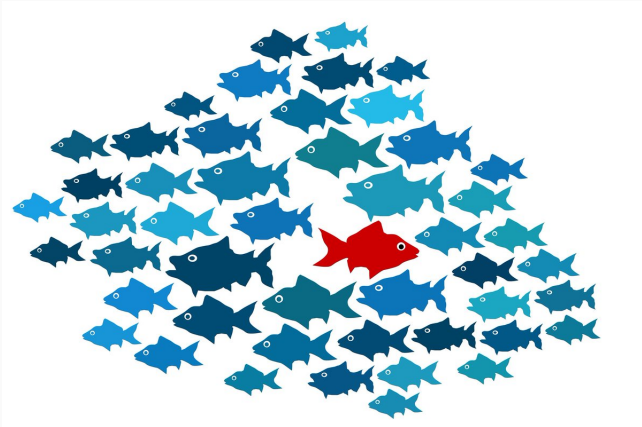
Ejemplo: regresión



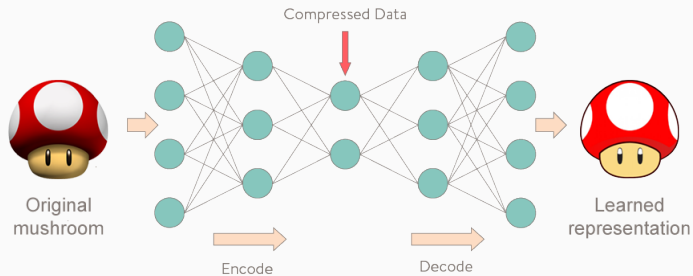
Ejemplo: clusterización



Ejemplo: detección de outliers

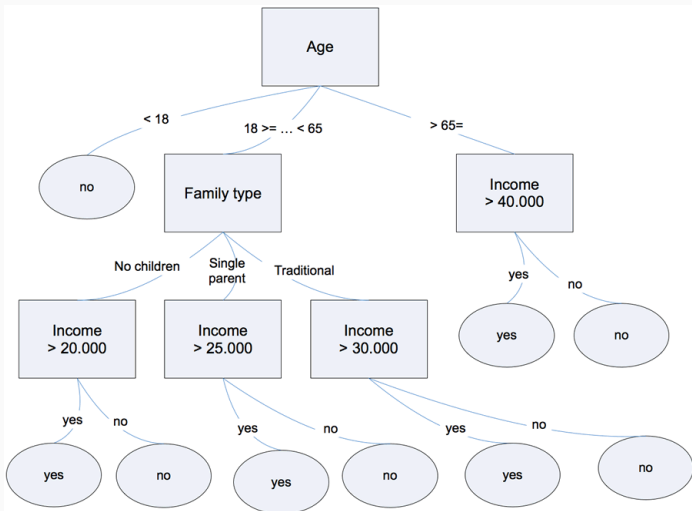


Ejemplo: reducción de dimensionalidad

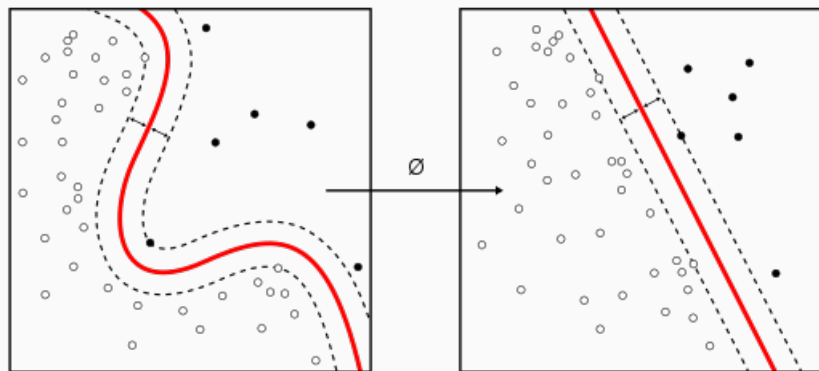


Algoritmos y técnicas (algunos)

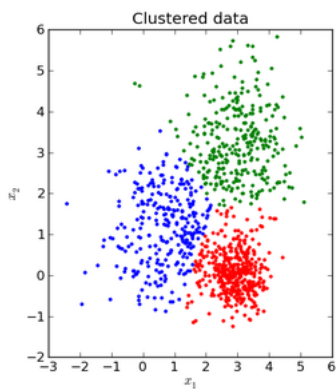
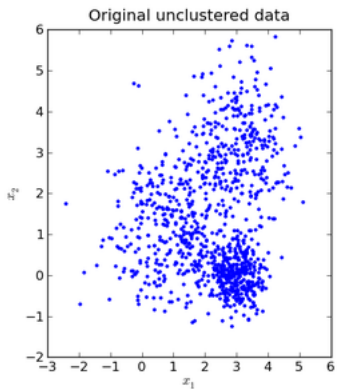
Árboles de decisión: predicción, regresión



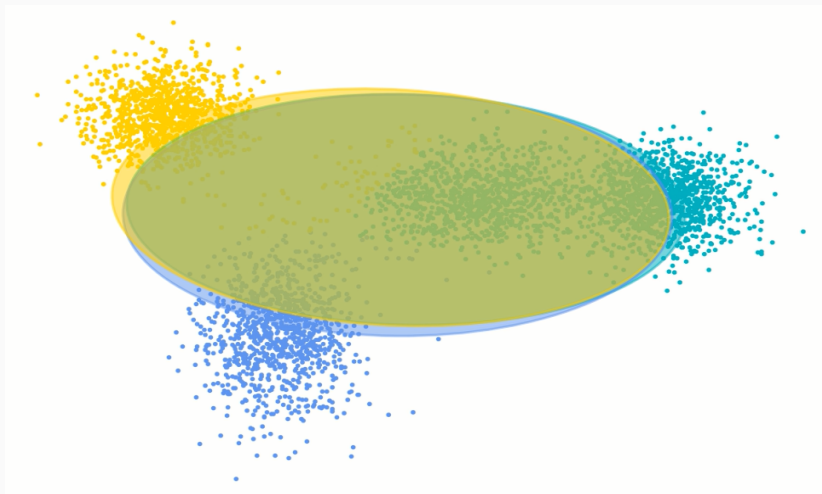
Support Vector Machine: clasificación



Mezcla de Gaussianas (GMM)

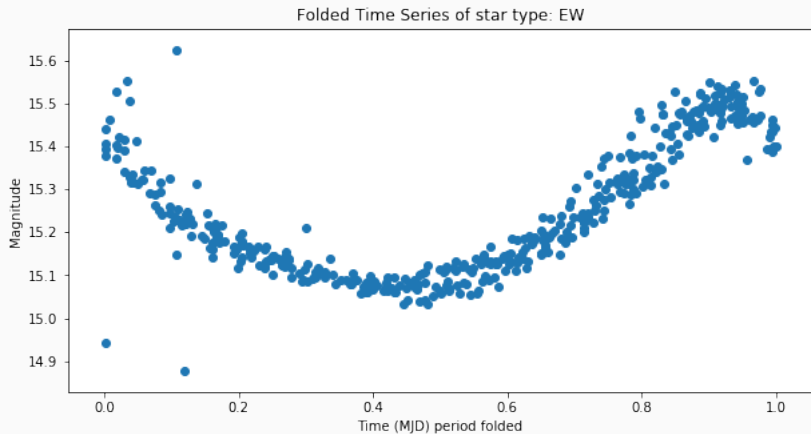


GMM (Gaussian Mixture Model)

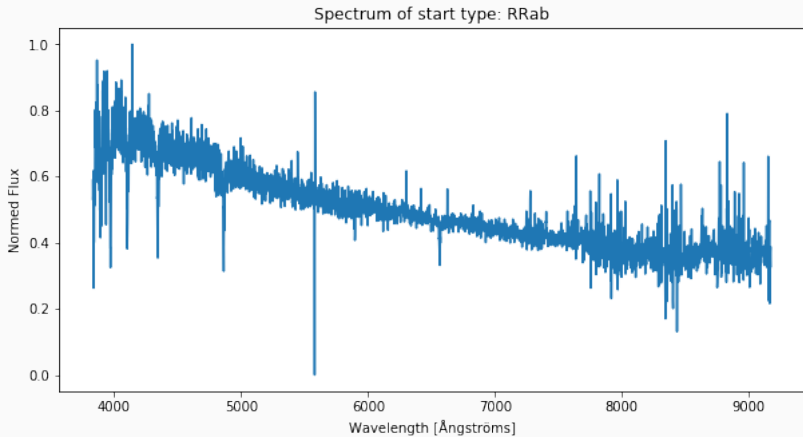


Trabajo de Investigación

Datos: Series de Tiempo



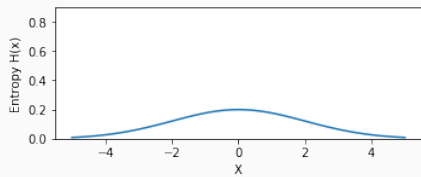
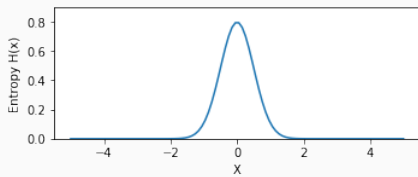
Datos: Espectros



- Recursos limitados
- Seleccionar objetos para *follow-up*
- ¿Criterio?

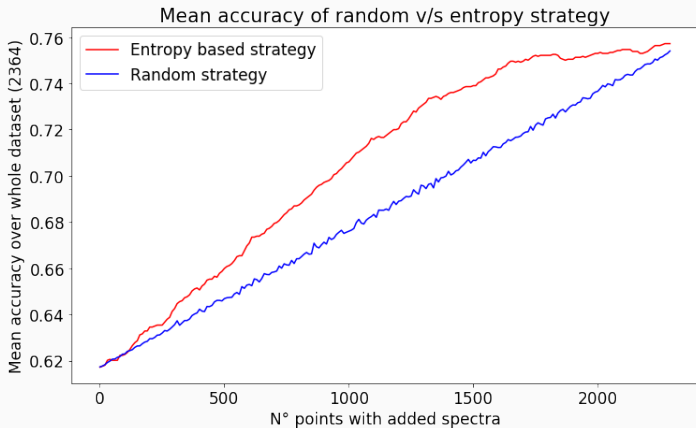
- ¿Qué es información?
 - En este trabajo: mejor clasificación.
 - Estimación del desempeño de la clasificación: Entropía

Entropía

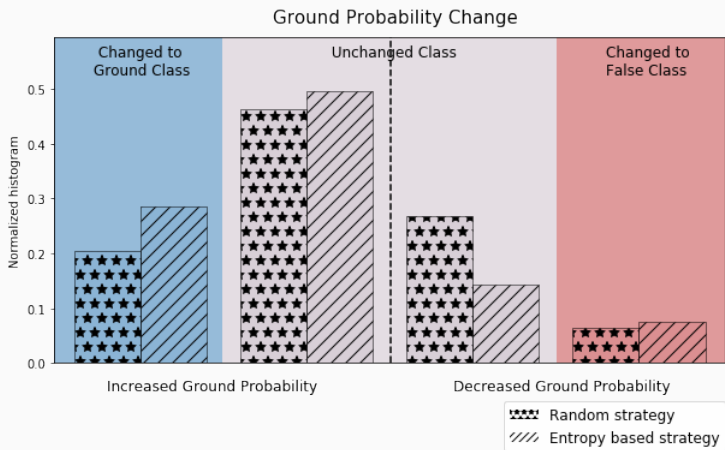


Metodología (idea)

- Extraemos características de series de tiempo y espectros a través de Auto Encoders variacionales (VAE).
- Entrenamos clasificadores solo con series de tiempo.
- Entrenamos clasificadores con series de tiempo y espectros.
- Estimamos el espectro de un objeto, dada su serie de tiempo.
- Estimamos el cambio de entropía en la clasificación con los modelos anteriores.
- Seleccionamos candidatos que presenten un gran cambio en la entropía.



Resultados

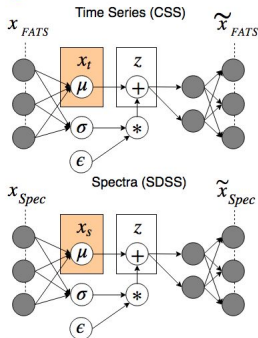


- Primer approach a abordar este problema.
- Criterio de selección mejor que aleatorio.
- Hay una correlación entre entropía y desempeño de clasificación.
- Aún quedan varias mejoras como trabajo futuro, tales como el criterio de mejora y los diversos modelos usados.

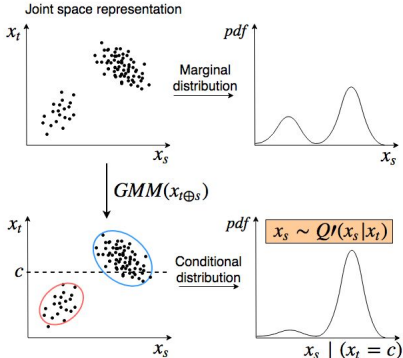
Preguntas?

Metodología

1 Feature Extraction



2 Conditional Distribution



3 Train Classifier

$$x_t \xrightarrow{\text{Random Forest}} P(y | x_t)$$

$$x_{t \oplus s} \xrightarrow{\text{Random Forest}} Q(y | x_{t \oplus s}) \rightarrow Q(y | x_t) = \int Q(y | x_{t \oplus s}) * Q'(x_s | x_t) dx_s \stackrel{\text{MCMC}}{\approx} \sum Q(y | x_{t \oplus s}) * Q'(x_s | x_t)$$

4 Candidate Selection

$$\Delta H(y|x_t) = H_Q(y|x_t) - H_P(y|x_t) \geq 0 \rightarrow \text{Add } x_s \rightarrow y = \operatorname{argmax}_y P(y|x_t)$$

$$< 0 \rightarrow \text{Add } x_s \rightarrow y = \operatorname{argmax}_y Q(y|x_{t \oplus s})$$

$$x_t \oplus x_s = x_{t \oplus s}$$