

Predicción de movimientos de acciones a partir de tweets y precios históricos

Enrique Cortés, Orlando Uc

Diciembre de 2020

- Introducción.
- Formulación del problema.
- Recolección de datos.
- Descripción general del modelo.
- Componentes del modelo.
 - Market Information Encoder.
 - Variational Movement Decoder.
 - Attentive Temporal Auxiliary
- Experimentos.
- Comentarios sobre la implementación con datos de México.
- Conclusiones
- Referencias.

La predicción del movimiento de acciones ha atraído durante mucho tiempo tanto a inversores como investigadores.

Se presenta un modelo para predecir el movimiento del precio de las acciones a partir de tweets y precios históricos.

Si una compañía tiene un gran escándalo en el día de trading d_1 , generalmente, los precios de sus acciones tenderán a disminuir hasta un día d_2 .

Si un modelo es capaz de identificar la baja en los precios en el intervalo $[d_1, d_2]$, entonces se pueden desarrollar estrategias de inversión que permita aprovechar esa situación.

Lo anterior viene del hecho de que la información pública, por ejemplo, un escándalo, requiere tiempo para ser absorbida en los movimientos de los precios (Luss y d'Aspremont, 2015).

Se propone el modelo StockNet, un modelo generativo profundo para la predicción del movimiento de las acciones, que toma en consideración la estocasticidad del mercado, la información caótica del mercado (tweets) y la predicción temporalmente dependiente (precios históricos).

En comparación con otros baselines, los experimentos realizados muestran que StockNet logra un mejor desempeño.

Formulación del problema

El objetivo es predecir el movimiento del precio de una acción s en una colección de acciones preseleccionada S en un día de trading objetivo d .

Formalmente, usamos la información del mercado M compuesta por tweets, y precios históricos, en la ventana de observación $[d - \Delta d, d - 1]$ donde Δd es un tamaño de lag fijo.

Estimamos el movimiento binario donde 1 denota alza y 0 denota baja,

$$y = \begin{cases} 1 & p_d^c > p_{d-1}^c \\ 0 & p_d^c \leq p_{d-1}^c \end{cases}$$

donde p_d^c denota el precio de cierre ajustado ajustado por los movimientos corporativos que afecten el precio de las acciones, es decir, reparto de dividendos, etc.

En finanzas, las acciones se clasifican en 9 industrias: Materiales Básicos, Bienes de Consumo, Sanidad, Servicios, Servicios Públicos, Conglomerados, Financiero, Bienes Industriales y Tecnología.

Se seleccionan 88 acciones diferentes, en el periodo entre 01/01/2014 y 01/01/2016, siendo las 8 acciones de los Conglomerados y las 10 principales acciones (en tamaño de capital) de cada una de las otras 8 industrias.

Dos componentes del conjunto de datos:

- 1 Datos de Twitter,
- 2 Datos de precios históricos.

Los precios históricos de las 88 acciones fueron tomados de Yahoo Finance.

Los tweets fueron tomados desde la API oficial de twitter y fueron preprocesados a través del módulo NLTK de Python.

Descripción general del modelo

Dado que existen, por ejemplo, fines de semana y días inhábiles, se consideran como input los días de trading en lugar de los días calendarios.

Se propone una alineación del día de trading: se reordenan los corpus de tweets y los precios históricos, alineándolos con estos días de negociación T .

Específicamente, para el t -ésimo día de trading, se reconoce la información del mercado M_t en $[d_{t-1}, dt)$ y los precios históricos p_t en d_{t-1} , para predecir el movimiento y_t en d_t .

Descripción general del modelo

En la Figura 1 se muestra el proceso generativo del modelo. Se realiza el encoding de la información del mercado como una variable aleatoria $X = [x_1, \dots, x_T]$, desde la que se genera un factor latente $Z = [z_1, \dots, z_T]$ para la predicción. Se busca modelar la distribución de probabilidad condicional $p_\theta(y|X) = \int_Z p_\theta(y, Z|X)$ en lugar de $p_\theta(y_T|X)$.

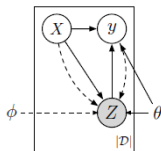


Figura: Gráfico del proceso generativo de la información observada del mercado del movimiento de las acciones.

Descripción general del modelo

StockNet, de manera general, tiene como componentes:

- 1 Market Information Encoder (MIE) que realiza el encoding de los precios históricos y los tweets a X .
- 2 Variational Movement Decoder (VMD) que infiere a Z con X , y realiza el decoding de los movimientos de las acciones y a partir de X, Z .
- 3 Attentive Temporal Auxiliary (ATA) que integra la función de costo temporal a través de un mecanismo de atención para el entrenamiento del modelo.

Arquitectura del modelo StockNet

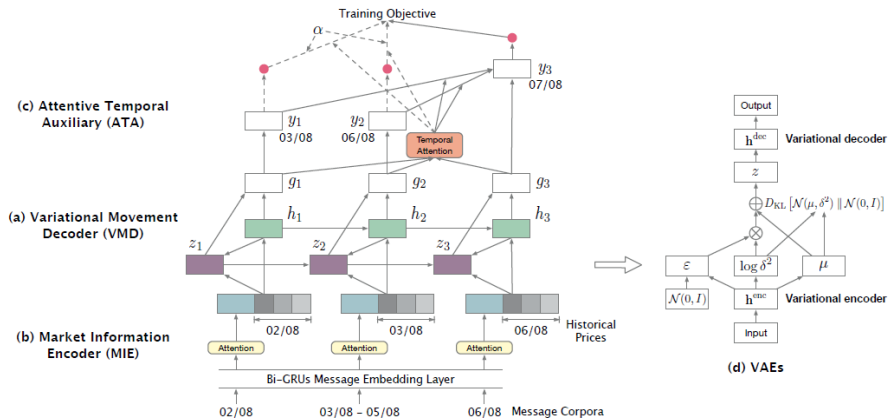


Figura: Arquitetura de la red StockNet

Componentes del modelo

Market Information Encoder

Cada input temporal se define como:

$$x_t = [c_t, p_t]$$

donde c_t y p_t son el embedding del corpus y el vector de precios históricos, respectivamente.

La capa consiste de una forward GRU y una backward GRU.

Componentes del modelo

Market Information Encoder

Las GRUs son:

$$\vec{h}_f = \overrightarrow{\text{GRU}}(e_f, \vec{h}_{f-1})$$

$$\overleftarrow{h}_b = \overleftarrow{\text{GRU}}(e_b, \overleftarrow{h}_{b+1})$$

$$m = \frac{\vec{h}_{I^*} + \overleftarrow{h}_{I^*}}{2}$$

y donde los valores ocultos \vec{h}_{I^*} y \overleftarrow{h}_{I^*} son promediados para obtener el tweet embedding m , que se guardan en una matriz de tweets embeddings M_t .

La calidad de los tweets varía drásticamente, por lo que se realiza una ponderación a través de medidas de inteligencia colectiva, de la forma:

$$u_t = \zeta(w_u^T \tanh(W_{m,u} M_t))$$

donde ζ es la función sigmoide y $W_{m,u}$, w_u son parámetros del modelo.

Componentes del mercado

Market Information Encoder

Entonces, el embedding del corpus es:

$$c_t = M_t u_t^T$$

Con respecto a los precios históricos, en lugar de ingresar el vector de precios $\tilde{p}_t = [\tilde{p}_t^c, \tilde{p}_t^h, \tilde{p}_t^l]$, se toma el precio ajustado y se normaliza con respecto a la observación anterior, de ta forma que:

$$p_t = \frac{\tilde{p}_t^c}{\tilde{p}_{t-1}^c} - 1$$

Componentes del modelo

Variational Movement Decoder

Se utilizan redes neuronales profundas para ajustar las distribuciones latentes, es decir, la distribución a priori $p_0(z_t|z_{<t}, x_{<t})$ y la posterior $p_\theta(z_t|z_{<t}, x_{<t}, y_t)$ y esquivar la complejidad de la estimación a partir de aproximación neuronal y reparametrización.

A partir de una minimización de la divergencia de Kullback-Leibler se maximiza el siguiente límite inferior recurrente variacional:

$$\begin{aligned}\mathcal{L}(\theta, \phi; X, y) &= \sum_{t=1}^T \mathbb{E}_{q_\phi(z_t|z_{<t}, x_{<t}, y_t)} \{ \log p_\theta(y_t|z_{<t}, x_{<t}) \\ &\quad - D_{KL}[q_\phi(z_t|z_{<t}, x_{<t}, y_t) || (y_t|z_{<t}, x_{<t})] \} \\ &\leq \log p_\theta(y|X)\end{aligned}$$

Componentes del modelo

Variational Movement Decoder

Como en el caso de las series de tiempo, VMD utiliza una RNN con una GRU para extraer las características y realizar el decoding de manera recurrente:

$$h_t^s = \text{GRU}(x_t, h_{t-1}^s)$$

y con la representación oculta:

$$h_t^z = \tanh(W_z^\phi[z_{t-1}, x_t, h_t^s, y_t] + b_z^\phi)$$

donde $W_{z,\mu}^\phi$, $W_{z,\mu}^\delta$, W_z^ϕ son matrices peso y b_μ^ϕ , b_δ^ϕ y b_z^ϕ son sesgos. Finalmente, se integran características determinísticas y la predicción final está dada por:

$$g_t = \tanh(W_g[x_t, h_t^s, z_t] + b_g)$$

$$\tilde{y}_t = \zeta(W_y g_t + b_y), t < T$$

donde W_g , W_y son matrices peso y b_g y b_y son sesgos. La función softmax ζ determina si la salida es un alza o una baja.

Componentes del modelo

Attentive Temporal Auxiliary

Se introduce un mecanismo de atención temporal, a través de un *information score* y un *dependency score*:

$$v'_i = w_i^T \tanh(W_{g,i} G^*)$$

$$v'_d = g_t^T \tanh(W_{g,d} G^*)$$

$$v^* = \zeta(v'_i \odot v'_d)$$

Los information score v'_i evalúan la calidad de la información de los días de trading, mientras que los dependency score v'_d capturan sus dependencias con respecto a la respuesta. Integrando los dos se obtiene el peso de atención normalizado final v^* .

Componentes del modelo

Attentive Temporal Auxiliary

A partir de v^* se construye el vector de peso temporal final $v \in \mathbb{R}^{1 \times T}$

$$v = [\alpha v^*, 1]$$

donde 1 es para la predicción principal y se adopta un peso auxiliar $\alpha \in [0, 1]$ para controlar los efectos auxiliares generales en el entrenamiento del modelo. Finalmente, se escribe la función objetivo \mathcal{F} por recomposición:

$$\mathcal{F}(\theta, \phi; X, y) = \frac{1}{N} \sum_n^N v^{(n)} f^{(n)}$$

Y se toma la derivada de \mathcal{F} con respecto a todos los parámetros del modelo $\{\theta, \phi\}$ a través de backpropagation para la actualización.

Componentes del modelo

Attentive Temporal Auxiliary

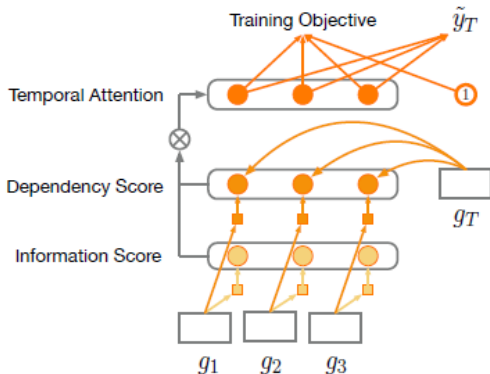


Figura: Atención temporal del modelo StockNet.

Se utilizó una ventana de observaciones de 5 días y 32 muestras aleatorias por batch. El número máximo de tokens por tweet y el máximo número de tweets en un día de trading son 30 y 40 respectivamente, tomados de manera empírica.

El tamaño de embedding de palabras es de 50 para limitar el costo computacional. Se utiliza el optimizador Adam con un learning rate inicial de 0.001. Se utiliza una tasa de dropout de 0.3 para regularización. El ajuste del modelo se implementó en Tensorflow.

Se utilizan la precisión y el Coeficiente de correlación de Matthews (MCC) que evita el sesgo de medición debido al sesgo de los datos. Dada la matriz de confusión:

$$\begin{pmatrix} tp & fn \\ fp & tn \end{pmatrix}$$

la métrica MCC se define cómo:

$$MCC = \frac{tp \times tn - fp \times fn}{\sqrt{(tp + fp)(tp + fn)(tn + fp)(tn + fn)}}$$

Para hacer un análisis detallado de todos los componentes de StockNet se realizaron cinco variaciones:

- `HedgeFundAnalyst`: La versión completa de StockNet.
- `TechnicalAnalyst`: La versión de StockNet generativa utilizando solamente los precios históricos.
- `FundamentalAnalyst`: La versión de StockNet generativa utilizando solamente la información de los tweets.
- `IndependentAnalyst`: La versión de StockNet generativa sin utilizar predicciones auxiliares.
- `DiscriminativeAnalyst`: La versión de StockNet discriminativa que optimiza directamente la verosimilitud objetivo.

Baseline models	Acc.	MCC	StockNet variations	Acc.	MCC
RAND	50.89	-0.002266	TECHNICALANALYST	54.96	0.016456
ARIMA (Brown, 2004)	51.39	-0.020588	FUNDAMENTALANALYST	58.23	0.071704
RANDFOREST (Pagolu et al., 2016)	53.08	0.012929	INDEPENDENTANALYST	57.54	0.036610
TSLDA (Nguyen and Shirai, 2015)	54.07	0.065382	DISCRIMINATIVEANALYST	56.15	0.056493
HAN (Hu et al., 2018)	57.64	0.051800	HEDGEFUNDANALYST	58.23	0.080796

Figura: Medidas de precisión y de MCC de los baselines y de las variaciones StockNet.

Experimentación

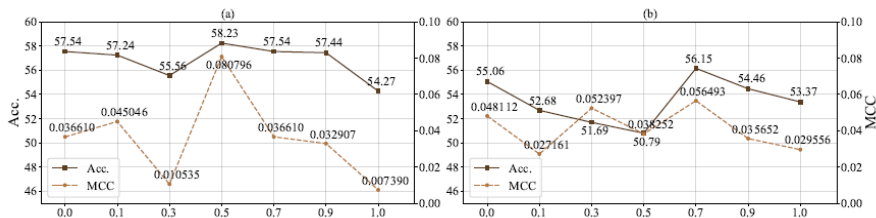


Figura: (a) Desempeño de HedgeFundAnalyst con diferentes valores de α . (b) Desempeño de DiscriminativeAnalyst con diferentes valores de α .

Comentarios sobre la implementación con datos de México

Para realizar la implementación del modelo StockNet con datos de México primero sería necesario conectar con la API de Twitter que permita descargar los tweets de las empresas mexicanas de las que se desee analizar el movimiento de los precios de sus acciones.

Habría que recolectar la información de los precios de las acciones de interés, que se podría realizar a partir de una plataforma como *mx.investing.com*.

Se tendrían que obtener (o crear) embeddings preentrenados de tweets en español.

Una posible problemática sería la falta de conexión de la población general con el comportamiento de la Bolsa Mexicana de Valores, es decir, existe la posibilidad de que no existan suficientes tweets que permitan alimentar el modelo.

También se tendría que realizar el preprocesamiento de los tweets en español, o bien, considerar como una posibilidad el incluir también tweets en inglés.

Actualmente no existen artículos de investigación que cumplan con las características anteriores.

- Se demostró la efectividad de los modelos generativos profundos para la predicción del movimiento de acciones a partir de información de redes sociales al introducir StockNet.
- Se probó el modelo en un nuevo dataset y se mostró que tiene un mejor desempeño que otros baselines robustos, incluyendo implementaciones de trabajos previos.
- Se comentó que las adecuaciones principales para implementar el modelo StockNet con datos de México sería primero generar la base de datos compuesta por los tweets y los precios históricos, así como la creación de embeddings preentrenados de tweets en español.

- Ding, X., Zhang, Y., Liu, T., & Duan, J. (2014, October). Using structured events to predict stock price movement: An empirical investigation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)* (pp. 1415-1425).
- Ding, X., Zhang, Y., Liu, T., & Duan, J. (2015, June). Deep learning for event-driven stock prediction. In *Twenty-fourth international joint conference on artificial intelligence*.
- Hu, Z., Liu, W., Bian, J., Liu, X., & Liu, T. Y. (2018, February). Listening to chaotic whispers: A deep learning framework for news-oriented stock trend prediction. In *Proceedings of the eleventh ACM international conference on web search and data mining* (pp. 261-269).
- Luss, R., & d'Aspremont, A. (2015). Predicting abnormal returns from news using text classification. *Quantitative Finance*, 15(6), 999-1012.

- Le, Q., & Mikolov, T. (2014, January). Distributed representations of sentences and documents. In *International conference on machine learning* (pp. 1188-1196).
- Oliveira, N., Cortez, P., & Areal, N. (2013, June). Some experiments on modeling stock market behavior using investor sentiment analysis and posting volume from Twitter. In *Proceedings of the 3rd International Conference on Web Intelligence, Mining and Semantics* (pp. 1-8).
- Schumaker, R. P., & Chen, H. (2009). Textual analysis of stock market prediction using breaking financial news: The AZFin text system. *ACM Transactions on Information Systems (TOIS)*, 27(2), 1-19.
- Xu, Y., & Cohen, S. B. (2018, July). Stock movement prediction from tweets and historical prices. *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (pp. 1970-1979).