

رگرسیون بیزی و بیت کوین

علیرضا کشاورز (a. keshavarz@khu. ac. ir)

پروژه یادگیری ماشین-ترم اول ۱۴۰۱

مقاله: <http://arxiv.org/pdf/1410.1231.pdf>

دانشکده ریاضی و کامپیوتر-دانشگاه خوارزمی

معرفی

این مقاله تحت عنوان رگرسیون بیزی و بیت کوین روش رگرسیون بیزی و اثربخشی آن برای پیش‌بینی تغییرات قیمت بیت‌کوین (یک ارز دیجیتال) را مورد بحث قرار می‌دهد.

رگرسیون بیزی به استفاده از داده‌های تجربی برای انجام استنتاج بیزی می‌پردازد. در بسیاری کاربردها از رگرسیون به عنوان "مدل منبع نهفته" (latent source model) استفاده می‌شود. اولین بار، رگرسیون بیزی برای مدل منبع نهفته به منظور دسته‌بندی باینری (binary classification) معرفی و مورد بحث قرار گرفت به این صورت که کارایی تئوری روش به خوبی کارایی تجربی آن برای ایجاد دسته‌بندی باینری، بررسی شد.

در این مقاله، از آن برای پیش‌بینی مقدار واقعی قیمت بیت‌کوین استفاده شده است. بر اساس این روش پیش‌بینی قیمت، یک استراتژی ساده برای معامله بیت‌کوین طراحی شده است. این استراتژی می‌تواند سرمایه‌گذاری را در مدت کمتر از ۶۰ روز تقریباً دو برابر کند، در صورتی که بر اساس ردیابی داده‌های واقعی اجرا شود.

I. رگرسیون بیزی

طرح مساله رگرسیون:

به ما n داده برچسب دار آموزشی داده شده است

(x_i, y_i) for $1 \leq i \leq n$. هدف مساله استفاده از داده‌های آموزشی برای پیش‌بینی برچسب (y) برای داده تست (x) است.

رویکرد حل مساله:

K منبع نهفته متمایز $s_1, \dots, s_K \in \mathbb{R}^d$ ، یک توزیع نهفته روی $1, \dots, K$ با احتمالات $\{\mu_1, \dots, \mu_K\}$ و K توزیع احتمال پنهان روی \mathbb{R} که با P_1, \dots, P_K نشان داده شده است، وجود دارد. هر نقطه داده‌ی برچسب دار (x, y) به صورت زیر تولید می‌شود:

$$T \in \{1, \dots, K\}$$

$$P(T = k) = \mu_k \quad \text{for } 1 \leq k \leq K$$

$$x = s_T + \varepsilon$$

که در آن ε یک متغیر تصادفی مستقل d بعدی است که نویز را نشان می‌دهد (آن را گوسی فرض می‌کنیم با میانگین بردار $0 = (0, \dots, 0) \in \mathbb{R}^d$ و ماتریس کوواریانس شناسایی y که از \mathbb{R} مطابق با توزیع P_T نمونه برداری می‌شود).

با توجه به این مدل، برای پیش‌بینی برچسب y مرتبط با مشاهده x ، می‌توانیم از توزیع شرطی y با توجه به x به صورت زیر استفاده کنیم:

$$\begin{aligned} P(y|x) &= \sum_{k=1}^T P(y|x, T = k) P(T = k|x) \\ &\propto \sum_{k=1}^T P(y|x, T = k) P(x|T = k) P(T = k) \\ &= \sum_{k=1}^T P_k(y) P(\varepsilon = (x - s_k)) \mu_k \\ &= \sum_{k=1}^T P_k(y) \exp\left(-\frac{1}{2} \|x - s_k\|_2^2\right) \mu_k \end{aligned}$$

(نویسنده، هم‌ارزی در رابطه بالا را به سادگی از روابط قانون بیز نتیجه گرفته است

$$(\text{posterior} \propto \text{likelihood} \times \text{prior})$$

$$Z(x) = \sum_{i=1}^n \exp(-\frac{1}{4} \|x - x_i\|_2^2)$$

و $y \in \mathbb{R}^n$ با i امین عنصر y_i باشد، آنگاه

$$\hat{y} \equiv E_{emp} [y|x]$$

$$\hat{y} = X(x) y$$

در این مقاله، از رابطه قبل (\hat{y}) برای پیش‌بینی تغییرات آینده قیمت بیت‌کوین استفاده شده است.

هدف کلی محقق در این مقاله، علاقه به درک این موضوع است که آیا اطلاعاتی در داده‌های تاریخی مربوط به بیت‌کوین وجود دارد که بتواند به پیش‌بینی تغییرات قیمت در آینده در بیت‌کوین کمک کند و بنابراین به توسعه استراتژی کمی سودآور با استفاده از بیت‌کوین کمک کند. همانطور که قبل از این هم بررسی کردیم، از رگرسیون بیزی الهام گرفته از مدل منبع نهفته، استفاده خواهد شد.

II. تجارت بیت‌کوین

ارتباط مدل منبع پنهان:

یکی از رویکردهای رایج در صنعت مالی، تحلیل تکنیکال است، که فرض می‌کند حرکات قیمت از مجموعه‌ای از الگوها پیروی می‌کند و می‌توان از حرکات قیمت گذشته برای پیش‌بینی بازده‌های آتی تا حدی استفاده کرد. مطالعات نشان دادند که برخی از الگوهای هندسی توسعه یافته تجربی، مانند heads – and – shoulders، triangle و double – top – and – bottom می‌توانند برای پیش‌بینی تغییرات قیمت در آینده استفاده شوند.

مدل منبع نهفته دقیقاً در تلاش است تا وجود چنین الگوهای زیربنایی را که منجر به تغییرات قیمت می‌شود، مدل‌سازی کند. تلاش برای توسعه الگوها با کمک یک متخصص انسانی یا تلاش برای شناسایی دقیق الگوها در داده‌ها، می‌تواند چالش برانگیز و تا حدی ذهنی باشد. در عوض، استفاده از رویکرد رگرسیون بیزی همانطور که قبل

بنابراین، تحت مدل منبع نهفته، مسئله رگرسیون به یک مسئله استنتاج بیزی بسیار ساده تبدیل می‌شود. با این حال، مشکل، عدم آگاهی از پارامترهای "مخفی" مدل منبع است. به طور خاص، عدم آگاهی از K منابع (s_1, \dots, s_k) ، احتمالات (μ_1, \dots, μ_k) و توزیع‌های احتمال P_1, \dots, P_k .

برای غلبه بر این چالش، الگوریتم ساده زیر را پیشنهاد شده است:

از داده‌های تجربی برای تخمین توزیع شرطی y با توجه به x در روابط بالا می‌توان استفاده کرد. به طور خاص، با توجه به n نقطه داده (x_i, y_i) ، $1 \leq i \leq n$ احتمال شرطی تجربی برابر خواهد بود با:

$$P_{emp}(y|x) = \frac{\sum_{k=1}^T P(y = y_i) \exp(-\frac{1}{4} \|x - x_i\|_2^2)}{\sum_{i=1}^n \exp(-\frac{1}{4} \|x - x_i\|_2^2)}$$

تخمین تجربی پیشنهادی در رابطه بالا بدین صورت است که در دسته‌بندی باینری، y مقادیر $\{0, 1\}$ را می‌گیرد و بر این اساس قانون دسته‌بندی زیر را پیشنهاد می‌کند:

$$\text{compute ratio: } \frac{P_{emp}(y = 1|x)}{P_{emp}(y = 0|x)} = \frac{\sum_{i=1}^n P(y_i = 1) \exp(-\frac{1}{4} \|x - x_i\|_2^2)}{\sum_{i=1}^n P(y_i = 0) \exp(-\frac{1}{4} \|x - x_i\|_2^2)}$$

اگر نسبت $1 < y$ ، در غیر این صورت

$y = 0$ را اعلام کنید. به طور کلی، برای تخمین امید

شرطی y ، برای مشاهده x ، رابطه زیر پیشنهاد شده است:

$$E_{emp}[y|x] = \frac{\sum_{i=1}^n y_i \exp(-\frac{1}{4} \|x - x_i\|_2^2)}{\sum_{i=1}^n \exp(-\frac{1}{4} \|x - x_i\|_2^2)}$$

تخمین در رابطه بالا را می‌توان به طور معادل به عنوان یک تخمین زننده "خطی" مشاهده کرد:

اگر بردار $X(x) \in \mathbb{R}^n$ به صورت زیر

$$X(x)_i = \exp\left(-\frac{1}{4} \|x - x_i\|_2^2\right) / Z(x)$$

از این هم گفته شد، به ما امکان می دهد از وجود الگوها برای پیش بینی بهتر بدون یافتن صریح آنها استفاده کنیم.

داده:

محقق در مقاله اصلی، برای انجام آزمایش‌ها، از داده‌های مربوط به قیمت و دفتر سفارش به دست آمده از Okcoin.com که یکی از بزرگترین صرافی‌های فعال در چین است، استفاده کرده است. داده‌های استفاده شده مربوط به دوره زمانی بین فوریه ۲۰۱۴ تا ژوئیه ۲۰۱۴ است. کل نقاط داده ی خام بیش از ۲۰۰ میلیون بوده است. اطلاعات دفترچه سفارش، شامل ۶۰ بهترین قیمت است که در یک زمان معین مایل به خرید یا فروش است. (نقاط داده در فاصله زمانی هر دو ثانیه به دست آمده است)

به منظور سهولت محاسباتی، یک سری زمانی جدید با فاصله زمانی ۱۰ ثانیه ساخته شده است. هر یک از نقاط داده ی خام به نزدیکترین نقطه در ۱۰ ثانیه آینده نقشه برداری شده است. با افزایش بازه زمانی، "خطای" جزئی در دقت مدل به وجود خواهد آمد، و از آنجایی که استراتژی معاملاتی ما در مقیاس زمانی بزرگتر عمل می کند، این خطا ناچیز خواهد بود.

استراتژی تجارت و پیش بینی تغییر قیمت:

روش اصلی برای میانگین تغییر قیمت (Δp) در بازه ی ۱۰ ثانیه‌ای، رگرسیون بیزی است.

روش بیان شده در مقاله بدین صورت است که از یک سری زمانی تاریخ دار از تغییرات بیت کوین استفاده می کند. این سری زمانی از تجمیع بازه های زمانی ۱۰ ثانیه ای در چند ماه بدست آمده است. با تجمیع این بازه های زمانی یک سری زمانی (یک بردار) بسیار بزرگ بدست می آید. از این سری زمانی استفاده می شود و از آن سه زیر مجموعه داده سری زمانی با سه طول مختلف تولید می کنیم:

S_1 با طول زمانی ۳۰ دقیقه، S_2 با طول زمانی ۶۰ دقیقه و S_3 با طول زمانی ۱۲۰ دقیقه.

اکنون در یک نقطه زمانی معین، برای پیش‌بینی تغییر آتی Δp ، از داده‌های تاریخی سه طول زیر استفاده می‌کنیم: ۳۰ دقیقه قبل، ۶۰ دقیقه قبل و ۱۲۰ دقیقه قبل که به ترتیب با x^1, x^2, x^3 نشان داده می شوند.

پیش بینی میانگین تغییر قیمت با استفاده از رگرسیون بیزی انجام می شود. در همین راستا، رگرسیون بیزی از x^j با نمونه های تاریخی S^j برای Δp^j ($1 \leq j \leq 3$) و

$r = (v_{bid} - v_{ask}) / (v_{bid} + v_{ask})$ استفاده می کند (v_{bid} کل حجمی است که افراد مایل به خرید در ۶۰ سفارش برتر هستند و v_{ask} کل حجمی است که افراد مایل به فروش در ۶۰ سفارش برتر بر اساس داده های دفتر سفارش فعلی هستند). تخمین نهایی Δp به صورت تولید می شود:

$$\Delta p = w_0 + \sum_{j=1}^3 w_j \Delta p^j + w_4 r$$

$w = (w_0, \dots, w_4)$ نرخ های یادگیری هستند.

اکنون در مورد یافتن S_j ، $1 \leq j \leq 3$ و یادگیری.

برای یافتن S_j ، $1 \leq j \leq 3$ و نرخ های یادگیری w ، مدت زمان را باید به سه دوره تقریباً مساوی تقسیم کرد. از اولین دوره زمانی برای یافتن الگوهای S_j استفاده می شود. دوره دوم برای یادگیری پارامترهای w و دوره سوم آخر برای ارزیابی عملکرد الگوریتم استفاده می شود. یادگیری w به سادگی با یافتن بهترین تناسب خطی روی همه گزینه ها با توجه به انتخاب S_j ، $1 \leq j \leq 3$ انجام می شود.

انتخاب S_j ، برای این کار، تمام سری های زمانی ممکن با طول مناسب را باید انتخاب شوند (برای تاثیر بیشتر

بردارهای با ابعاد ۱۸۰، ۳۶۰ و ۷۲۰ به ترتیب برای S_1 ، (S_2, S_3) . هر کدام از این x_i ها برچسب مربوط به آن (y_i) با مشاهده میانگین تغییر قیمت در بازه زمانی ۱۰ ثانیه

پس از پایان مدت زمان x_i محاسبه می‌شود. این مخزن داده بسیار بزرگ است. برای تسهیل محاسبات روی ماشین تکی با رم ۱۲۸G و ۳۲ هسته، الگوها در ۱۰۰ خوشه با استفاده از الگوریتم $k - \text{means}$ خوشه‌بندی شده است. از بین این‌ها، ۲۰ خوشه موثر انتخاب و الگوهای نماینده از این خوشه‌ها برداشته شده است.

پس از محاسبه Δp باید استراتژی معاملاتی که بسیار ساده است را توضیح داد: در هر زمان، ما موقعیت $+1$ بیت کوین، 0 بیت کوین یا -1 بیت کوین را حفظ می‌کنیم. در هر نمونه زمانی، میانگین حرکت قیمت را در بازه ۱۰ ثانیه پیش‌بینی می‌کنیم. اگر $\Delta p > t$ (یک آستانه است) و همچنین موقعیت فعلی بیت‌کوین کمتر یا مساوی 0 باشد، یک بیت‌کوین می‌خریم. اگر $\Delta p < -t$ و همچنین موقعیت فعلی بیشتر یا مساوی 0 باشد، یک بیت‌کوین می‌فروشیم. و در غیر این صورت هیچ کاری انجام نمی‌دهیم.

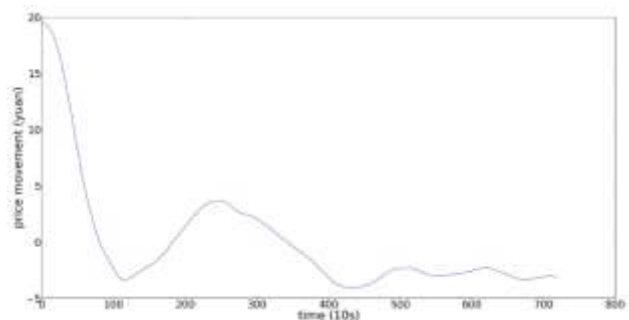
بر اساس شواهد، مشخص است که سود کل در ۳۳۶۲ یوان با ۲۸۷۲ معامله به طور کلی، با میانگین سرمایه گذاری ۳۷۸۱ یوان، به اوج خود خواهد رسید.

شکل زیر عملکرد بهترین استراتژی را در طول زمان نشان می‌دهد. قابل ذکر است، زمانی که نوسانات بازار بالا است، استراتژی در بخش میانی بهتر عمل می‌کند.

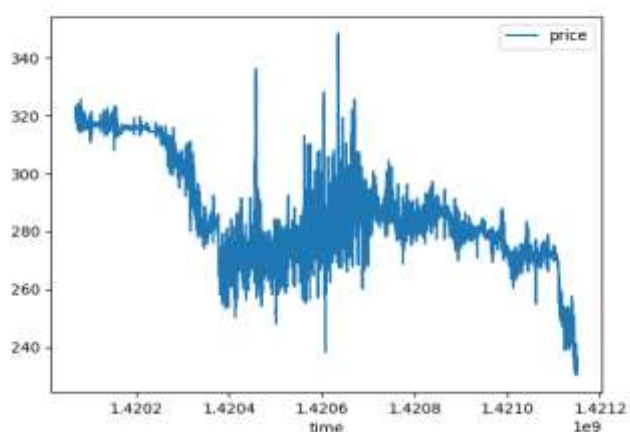
علاوه بر این، استراتژی همچنان سودآور است حتی زمانی که قیمت در آخرین بخش از دوره آزمایش کاهش می‌یابد.



IV. مقایسه نتایج



نمودار موجود در مقاله که در آن نقاط داده‌های بدست آمده بر حسب میانگین تغییر قیمت بیت‌کوین، با استفاده از دو پارامتر زمان و قیمت رسم شده است.



III. نتیجه گیری

در مقاله اصلی، استراتژی معاملاتی شرح داده شده بر روی یک سوم کل داده‌ها در مدت زمان ۶ می ۲۰۱۴ تا ۲۴ ژوئن ۲۰۱۴ به صورت علنی شبیه‌سازی شده است تا عملکرد استراتژی مورد بررسی قرار بگیرد. داده‌های آموزشی استفاده شده تماماً تاریخی هستند (یعنی قبل از ۶ می ۲۰۱۴ جمع‌آوری شده‌اند).

در این بررسی از آستانه‌های t مختلف استفاده شده است تا تغییرات عملکرد استراتژی مشاهده شود. بررسی‌ها نشان داده است که آستانه‌های مختلف عملکردهای مختلفی به همراه دارند. به طور مشخص، با افزایش آستانه، تعداد معاملات کاهش می‌یابد و میانگین زمان نگهداری افزایش می‌یابد. در عین حال، میانگین سود هر معامله نیز افزایش پیدا می‌کند.

نمودار بدست آمده پس از اجرای کد و گرفتن جواب که نشاندهنده میزان تغییرات قیمت داده ها بر حسب زمان و قیمت بیت کوین است. لازم به ذکر نمودار دوم در ابعاد بزرگتری رسم شده است اما از مقایسه نمودارها با یکدیگر می توان متوجه شد که نقاط داده ها پس از محاسبه میانگین تغییر قیمت و با افزایش زمان، رفته رفته قیمتشان کم شده است (می توان این نتیجه را از هر دو نمودار دریافت).