

① Validation و training برای hyper parameter test یک بار جدا train برای لینی

② (الف) y_1 و $V \downarrow$, $b \uparrow$ ← under fit $b \uparrow$, $V \uparrow$ ←
 (ب) y_2 : $b \downarrow$, $V \downarrow$ ← good fit $b \downarrow$, $V \downarrow$ ←
 (ج) y_3 : $b \downarrow$, $V \uparrow$ ← over fit $b \downarrow$, $V \uparrow$ ←

③ (الف) E بعد از مدتی تغییر نمی کند دلیل رفتار KNN
 (ب) B ابتدا over fit و سپس under fit
 (ج) D شکست اورفیت، رفت سریع
 (د) B دلیل over fit

Converge و local minima } B (ب)
 تقریباً مانند training } B (ب)
 زیاده over fit } A (ت)
 بعد از مدتی over fit } B (ت)



باسمه تعالی

دانشگاه صنعتی شریف

دانشکده مهندسی کامپیوتر

مقدمه‌ای بر یادگیری ماشین

استاد درس: دکتر زارچی

آزمون میانترم - سوال دوم

خوشه خوشه

در این مسئله می‌خواهیم به بررسی الگوریتم خوشه‌بندی k -means بپردازیم. فرض کنید:

$$X = x_1, x_2, \dots, x_n$$

داده‌های ما باشد و γ یک ماتریس Indicator باشد به این صورت که $\gamma_{ij} = 1$ اگر x_i متعلق به خوشه j باشد و در غیر این صورت برابر 0 است. فرض کنید $\mu_1, \mu_2, \dots, \mu_k$ میانگین خوشه‌ها باشند. اعوجاج J برای داده‌ها به صورت زیر محاسبه می‌شود:

$$J(\gamma, \mu_1, \dots, \mu_k) = n \sum_{j=1}^k \sum_{i=1}^n \gamma_{ij} \|x_i - \mu_j\|^2$$

همچنین k را به عنوان مجموعه خوشه‌ها در نظر بگیرید $C = 1, \dots, k$.

۱

آیا k -means نسبت به انتخاب نقاط اولیه حساس است، یعنی پاسخ آن بر اساس مجموعه‌ی نقاط اولیه تغییر می‌کند؟ اگر بله یک مثال ارائه کنید و اگر خیر، اثبات کنید. (۱۰ نمره)

پاسخ: بله (۳ نمره) مثال (۷ نمره)

۲

نشان دهید که الگوریتم k -means در زمان متناهی قدم به پایان می‌رسد. (راهنمایی: نشان دهید J ، تعداد محدودی حالت دارد). (۲۰ نمره)

پاسخ:

نشان دادن اینکه مقدار اعوجاج بعد از هرگام یا ثابت میماند یا کاهش پیدا میکند (۸ نمره) و اگر ثابت ماند الگوریتم به پایان می‌رسد (در گام‌های بعد هم ثابت میماند) (۳ نمره) اشاره به اینکه فضای حالات محدود است و به دست آوردن این مقدار (۵ نمره) و نتیجه‌گیری نهایی (۴ نمره)

۳

اگر ابعاد داده نسبت به تعداد نمونه‌ها خیلی زیاد باشد و عملاً نمونه‌ها در یک فضای بزرگ پراکنده باشند، برای بهبود خوشه‌بندی از چه روشی استفاده می‌کنید؟

پاسخ:

اشاره به الگوریتم‌های مناسب مانند PCA و DBSCAN و ... (۱۰ نمره)

۴

نشان دهید که کمینه J ، یک تابع غیر افزایشی بر حسب k یا همان تعداد خوشه‌ها است. در این صورت آیا انتخاب مقدار هایپرپارامتر K بر اساس کمینه کردن مقدار J ، ایده‌ی خوبی است؟ اگر نه، چه ایده‌ی بهتری دارید؟ (۲۵ نمره)

پاسخ:

با استفاده از استقرا، فرض کنید $L(k)$ کمینه تابع اعوجاج J برای k خوشه است. هدف نشان دادن این است که با افزایش تعداد خوشه‌ها از k به $k+1$ ، مقدار $L(k)$ کاهش یا ثابت می‌ماند: $L(k+1) \leq L(k)$.

خوشه‌بندی بهینه برای k خوشه، کمینه $L(k)$ است:

$$L(k) = \min_{\gamma, \mu_1, \dots, \mu_k} J(\gamma, \mu_1, \dots, \mu_k)$$

با افزایش تعداد خوشه‌ها به $k+1$ ، می‌توان یکی از خوشه‌های قبلی را به دو خوشه جداگانه تقسیم کرد. این تقسیم‌بندی جدید می‌تواند باعث کاهش مقدار تابع اعوجاج J شود یا حداقل مقدار آن ثابت بماند، زیرا تقسیم یک خوشه به دو خوشه می‌تواند انحراف درون‌خوشه‌ای را کاهش دهد. راه

دیگر اثبات: استفاده از انحراف درون خوشه‌ای و بین خوشه‌ای (۱۵ نمره)
 شرح اینکه چرا انتخاب هایپرپارامتر K بر اساس کمینه کردن J ایده‌ی خوبی نیست (۴ نمره) و ایده‌ی بهتر (۶ نمره)

۵

فرض کنید \hat{x} میانگین داده‌های نمونه باشد. مقادیر زیر را در نظر بگیرید:

$$T(X) = \frac{\sum_{i=1}^n \|x_i - \hat{x}\|^2}{n}$$

$$W_j(X) = \frac{\sum_{i=1}^n \gamma_{ij} \|x_i - \mu_j\|}{\sum_{i=1}^n \gamma_{ij}}$$

$$B(X) = \sum_{j=1}^k \frac{\sum_{i=1}^n \gamma_{ij}}{n} \|\mu_j - \hat{x}\|^2$$

در اینجا $T(X)$ نشان دهنده انحراف کلی، $W_j(X)$ انحراف درون خوشه‌ای و $B(X)$ انحراف بین خوشه‌ای است. رابطه‌ی بین این ۳ مقدار به چه صورت است؟ نشان دهید که k -means می‌تواند به عنوان کمینه‌کننده میانگین وزن‌دار مقادیر درون خوشه‌ای و به طور تقریبی بیشینه کردن انحراف بین خوشه‌ای دیده شود. ۳۵ نمره

پاسخ:

رابطه‌ی بین این ۳ مقدار به این صورت است:

$$\begin{aligned} nT(X) &= \sum_{i=1}^n \|x_i - \hat{x}\|^2 \\ &= \sum_{i=1}^n \|x_i - \mu_j + \mu_j - \hat{x}\|^2 = \sum_{i=1}^n \|x_i - \mu_j\|^2 + \sum_{i=1}^n \|\mu_j - \hat{x}\|^2 + \sum_{i=1}^n (x_i - \mu_j)^T (\mu_j - \hat{x}) \\ &= \sum_{j=1}^k \sum_{i=1}^n \gamma_{ij} \|x_i - \mu_j\|^2 + \sum_{j=1}^k \sum_{i=1}^n \gamma_{ij} \|\mu_j - \hat{x}\|^2 + \sum_{j=1}^k \sum_{i=1}^n \gamma_{ij} (x_i - \mu_j)^T (\mu_j - \hat{x}) \end{aligned}$$

۸ نمره

حال داریم:

۱.

$$\sum_{j=1}^k \sum_{i=1}^n \gamma_{ij} \|x_i - \mu_j\|^2 = \sum_{j=1}^k n_j W_j(X) \quad \text{Where } n_j = \sum_{i=1}^n \gamma_{ij}$$

۵ نمره

۲.

$$\sum_{j=1}^k \sum_{i=1}^n \gamma_{ij} \|\mu_j - \hat{x}\|^2 = \sum_{j=1}^k n_j \|\mu_j - \hat{x}\|^2 = nB(X)$$

۵ نمره

۳.

$$\mu_j = \frac{\sum_{i=1}^n \gamma_{ij} x_i}{\sum_{i=1}^n \gamma_{ij}} \Rightarrow \sum_{i=1}^n \gamma_{ij} (x_i - \mu_j)^T (\mu_j - \hat{x}) = 0$$

۵ نمره

پس داریم:

$$nT(X) = \sum_{j=1}^k n_j W_j(X) + nB(X)$$

۴ نمره

از آنجا که $T(X)$ ثابت است، الگوریتم k -Means $W(X)$ را کمینه و در نتیجه $B(X)$ را بیشینه می‌کند. ۸ نمره

پاسخنامه و ریزنمرات بخش خلاف شیب

سوال ۱:

$$w^{t+1} = w^t - \eta \nabla_{w_i} F(w) \quad (۵ \text{ نمره})$$

$$F(w) = l(x^{(j)}, y^{(j)}, w) = 0 \rightarrow \nabla_w F(w) = \nabla_w \left[y^{(j)} \left(\sum_{i=1}^d w_i x_i^{(j)} \right) - \ln(1 + \exp(\sum_{i=1}^d w_i x_i^{(j)})) \right]$$
$$\nabla_{w_j} F(w) = y^{(j)} x_i^{(j)} - \frac{x_i^{(j)} \exp(\sum_{i=1}^d w_i x_i^{(j)})}{1 + \exp(\sum_{i=1}^d w_i x_i^{(j)})}$$

(۱۵ نمره)

سوال ۲:

در صورت استفاده از داده ساختارهای متراکم میانگین پیچیدگی زمانی از حاصلضرب ابعاد پیروی می‌کند $O(NX)$ (۱۰ نمره)

در صورت استفاده از داده ساختارهای تنک میانگین پیچیدگی زمانی از تعداد داده‌های غیرصفر پیروی می‌کند. $O(d)$ (۱۰ نمره)

سوال ۳:

$$w_i^{(t+1)} = w_i^t - \eta \nabla_{w_i} F(w) \quad (۵ \text{ نمره})$$

$$\frac{\partial F(w)}{\partial w_i} = y^{(j)} x_i^{(j)} - \frac{x_i^{(j)} \cdot \exp(\sum_{i=1}^d w_i x_i^{(j)})}{1 + \exp(\sum_{i=1}^d w_i x_i^{(j)})} - \lambda w_i$$

(۱۵ نمره)

سوال ۴:

تنها تفاوت این است که از یک ضریب ثابت جدید استفاده کردیم که تاثیری در پیچیدگی ندارد. (۱۰ نمره)

بنابراین میانگین پیچیدگی زمانی برابر است با $O(NX)$. (۱۰ نمره)

سوال ۵:

اشاره به صفر بودن جمله دوم معادله آپدیت وزن‌ها (۵ نمره)

$$w_i^{(t+1)} = w_i^t - \eta \lambda w_i^t = (1 - \eta \lambda) w_i^t$$

$$w_i^{(t+k)} = (1 - \eta \lambda) w_i^{(t+k-1)} = (1 - \eta \lambda)^k w_i^t$$

نوشتن دو معادله بالا (۵ نمره)

سوال ۶:

باتوجه به زمان نیاز برای به روزرسانی وزن‌ها و خلاقیت الگوریتم از ۱۰ نمره داده شود.

①

+	-	-
+	+	-
+	+	+

②

+	-	-
+	+	-
+	+	+

③

+	-	-
+	+	-
+	+	+

④

+	-	-
+	+	-
+	+	+

⑤

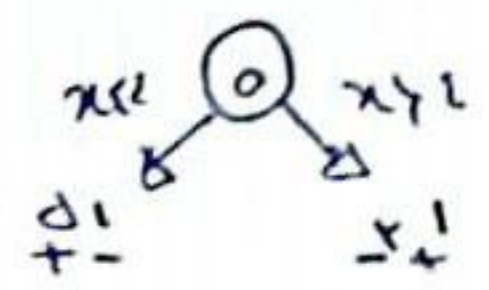
+	-	-
+	+	-
+	+	+

⑥

+	-	-
+	+	-
+	+	+

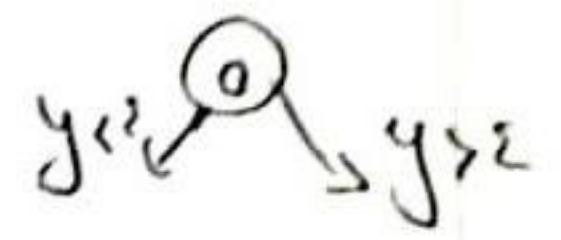
①

+	-	+	-
+	+		-
+	+		+



④

+	-	-
+	+	-
+	+	+

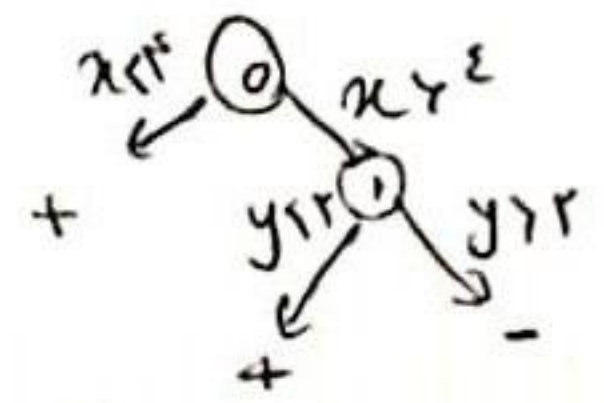


②

+	-	-
+	+	-
+	+	+

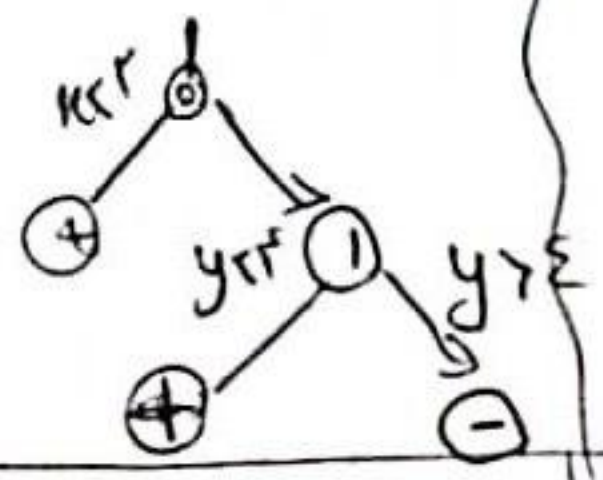
③

+	-	-
+	+	-
+	+	+



⑤

+	-	-
+	+	-
+	+	+



⑥

+	-	-
+	+	-
+	+	+

$$\frac{\partial L(\omega)}{\partial \omega} = -2X^T y + 2X^T X \omega + 2\Gamma^T \Gamma \omega = 0$$

جواب:
$$L(\omega) = y^T y + \omega^T X^T X \omega - 2\omega^T X^T y + \omega^T \Gamma^T \Gamma \omega$$

وائر مستقیم عبارت بالا می شود

$$\omega^* = (X^T X + \Gamma^T \Gamma)^{-1} X^T y$$

یعنی برای $\hat{\beta}$ دارد.

برای اینم سوابق به دست می آید که $\hat{\beta}$ بهترین حالت

$$A^T A + \lambda I$$

است عملیات

یعنی بردارهای

$$A^T A + \lambda I$$

است عملیات

X می باشد.

$$l(w) = E \left[\|y - (D \odot X)w\|^2 \right]$$

$$D \sim \text{Bern}(p)$$

→ هر دایره با احتمال p اینست

$$= \cancel{E} l(w) = E \|y - pXw\|^2$$

$$= E \left[(y^T y - 2w^T (D \odot X)^T y + w^T (D \odot X)^T (D \odot X) w) \right]$$

$$= y^T y - 2 E w^T \underbrace{E (D \odot X)^T y}_p + \cancel{w^T} E (\|D \odot X\|^2) w$$

$$= \|y - pX\hat{w}\|^2 + p(1-p) \|\hat{\Gamma} \hat{w}\|^2$$

$$L(\hat{w}) = \|y - X\hat{w}\|_2^2 + \lambda \|\hat{w}\|_2^2$$

$$L(w) = \|y - Xw\|_2^2 + \|\Gamma w\|_2^2$$

minimize $\Rightarrow \hat{w} = \frac{(\Gamma^T w)}{\sqrt{\lambda}} \Rightarrow \Gamma w = \sqrt{\lambda} \hat{w}$

$$\|\Gamma w\|_2^2 = \lambda \|\hat{w}\|_2^2$$

درون Γ $\Rightarrow \hat{w} = \sqrt{\lambda} T^{-1} \hat{w} \Rightarrow \|y - Xw\|_2^2 = \|\cancel{\Gamma} \hat{w}\|_2^2$

$$= \|y - \underbrace{\sqrt{\lambda} X T^{-1}}_{\hat{X}} \hat{w}\|_2^2 = \|y - \hat{X} \hat{w}\|_2^2$$

$$\Rightarrow L(\hat{w}) = \|y - \hat{X} \hat{w}\|_2^2 + \lambda \|\hat{w}\|_2^2$$

1. اشاره به تغییرات لازم در String های داده شده لازم است. (0.5 نمره)
جواب یکنایی وجود ندارد. یک مثال می تواند استفاده از k-means و یا رگرسیون باشد. (2.5 نمره)

2. کشتن Centralize مهم است (1 نمره)
میزان تأثیرگذاری را برای ویژگی های یکسان کنیم یا به عبارتی نرمال سازی ویژگی ها را داشته باشیم (2 نمره)
مطابق به توضیحات و دلایل برای سوار بالا و یا جواب ملی دیگر (1 نمره)

3. $C = \begin{bmatrix} 3 & -4 \\ -4 & 3 \end{bmatrix}$ \rightarrow حرف اولیه: یافتن مقادیر ویژه

$$\Rightarrow \det(C - \lambda I) = 0 \rightarrow (3 - \lambda)^2 - 16 = 0 \rightarrow \lambda^2 - 6\lambda - 7 = 0$$

$$(1 + 1)(1 - 7) = 0$$

مقادیر ویژه: $\lambda_1 = 7, \lambda_2 = -1$ (2 نمره)

$$Cv_1 = \lambda_1 v_1 \rightarrow v_1 = \begin{bmatrix} \sqrt{2}/2 \\ -\sqrt{2}/2 \end{bmatrix}$$

$$Cv_2 = \lambda_2 v_2 \rightarrow v_2 = \begin{bmatrix} \sqrt{2}/2 \\ \sqrt{2}/2 \end{bmatrix} \quad (2 \text{ نمره})$$

4. ابتدا حایه ماتریس کواریانس را بیابید. سه مقادیر ویژه را به دست آورد $(\lambda_1, \lambda_2, \dots, \lambda_n)$ با ترتیب نزولی $\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_n$

چاقچه به قدر اطلاعاتی که نیاز داریم اکثر $k\%$ جاست $\left(\frac{\sum_{i=1}^m \lambda_i}{\sum_{i=1}^n \lambda_i} \right) \cdot 100 \geq k$ $m \leftarrow \min m$ به صورت نیاز است (3 نمره).

$$\arg \max_{\|v\|_2^2 = 1} \text{Var}(Xv) \rightarrow \text{Var}(Xv) = v^T X^T X v$$

طبق کلاثر : $L = v^T X^T X v - \lambda (\|v\|_2^2 - 1)$

$$\rightarrow \frac{\partial L}{\partial v} = 0 = 2 \underbrace{X^T X v}_{\substack{n \sum \\ \text{Cov mat}}} - \lambda v \rightarrow n \sum v = \lambda v \quad (2 \text{ غره})$$

$$\Rightarrow \text{Var}(Xv) = v^T X^T X v = n v^T \left(\sum v \right) = \lambda v^T v = 1$$

بزرگترین مقدار و بزرگ جاترین واریانس دارد (2 غره)

6. طبق خلاصت و جالب بودن سؤال مطرح شده. (2 غره)