

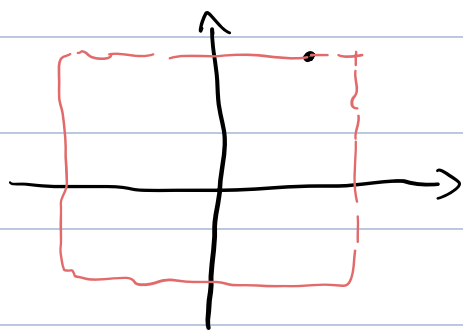
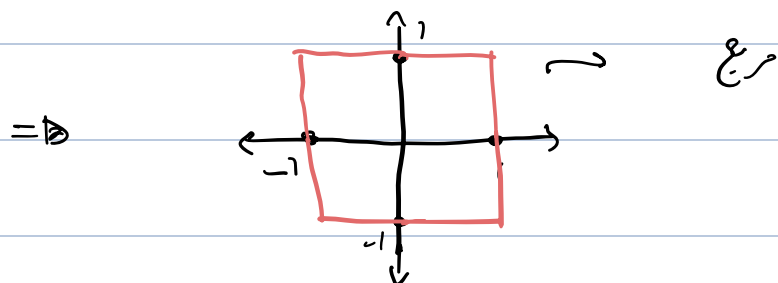
سوال یک

الف)

برای مثال بازه نقاط $[-1, 1] \times [-1, 1]$ را در نظر بگیرید.

مرکزهای نرم ها را رسم می کنیم.

برای (ایمپ) $\rightarrow 1 = \infty$ مرکز $\|z\|_2$

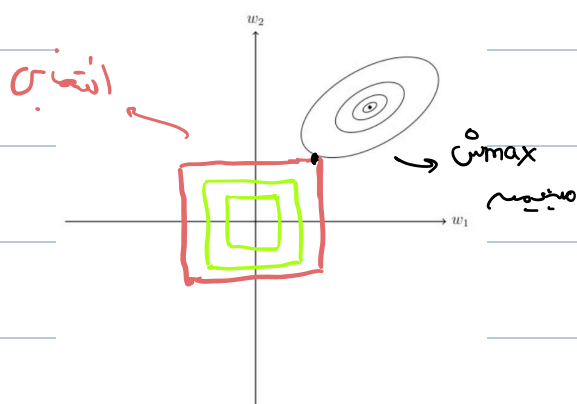


برای مجموعه این از وزن ها اگر نقطه را بگیریم.

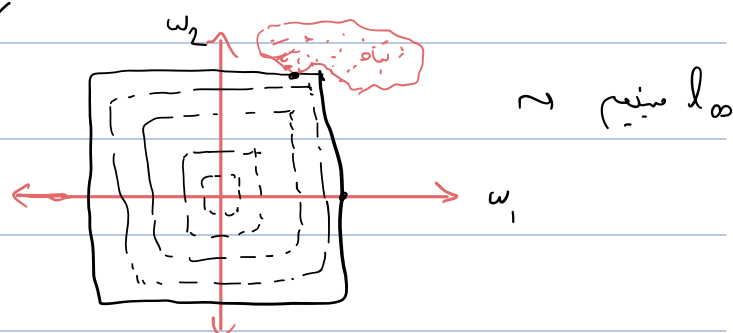
حال مربع شامل این را (مرفض و فاصل $= \infty$ می شه)

در نظر بگیریم سعی در کوچک کردن این مربع است.

پس نقطه انتخابی اگر از مرکز و مربع شروع کنیم و مربع را بزرگ کنیم نقطه انتخابی برضرد است \rightarrow که کوچکترین مربع میرا این نقطه است.



نگاه سوال:



(ب)

با استفاده از این نرم بین $\lambda \|w\|_\infty$ و ماکسیم اعصاب w هزینه ϕ میزنیم که بهین سبب ماکسیم وزن ها خیلی بزرگ نشود.

$$\min_w \ell(w; x, y) - \lambda \|w\|_\infty \quad \text{مقادیر}$$

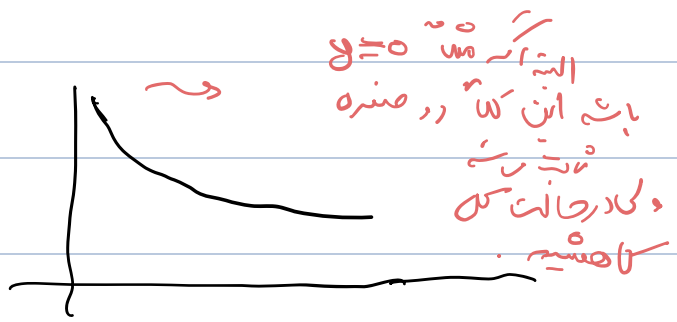
$$\Leftrightarrow \min_w \ell(w; x, y)$$

$$\text{s.t. } \|w\|_\infty \leq C$$

یعنی روی ماکس وزن ها محدودیت قرار می دهیم.

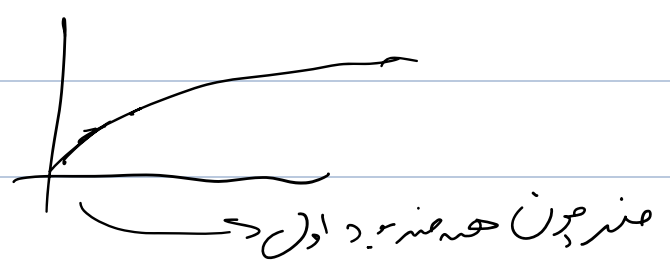
در شرایطی که نخواهیم وزن ها محدودیت اندازه داشته باشند از یک لیمیت خاص کمتر باشند کاربرد دارند.

الف | انبساط محدودیت : بدون هدف قدیمی پس احتیاطاً (محدود) دریاچه ها
 مایه های زیادی دارد، به مرور با افزایش S ، رنگ و اثرش در دریاچه کمتر
 کمتر شود پس مایه های کاهش می یابد اما ممکن است دچار تورفتگی شود
 که در این بخش که فضای ترین مد نظر مهم نیست.
 تا جایی که دریاچه به نقطه فوت شدن برسد پس از آن پس S
 اهمیت ندارد و فقط می شود



پس 4 کاهش می یابد.

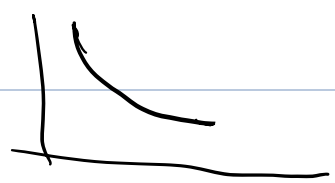
ب | با افزایش محدودیت S و H فضای فرغی مان را گسترش می دهیم



مدتی پس واریش افزایش می یابد.

پس 3

ج | انبساط با افزایش حوضه محدود حاصله در وی با گسترش H



کاهش می یابد. (حتی که H و S فاصله داشته باشند)
 با افزایش S بهریت نزدیک و دیر
 می شود = کاهش
 پس 4

سوال ۳

فکر کنم بخش الف و ب سوال جابجا شد ولی جفتش هست :-}}

Data Augmentation

یکی از کاربردهای روتین در ویژن است. وقتی دیتاستمون محدودیت هایی داره مثلا از نظر تعدادی یا جهت و ... برای اینکه مدل لرن شده تحت تاثیر اون محدودیت نشه می توانیم اگمنت شده های عکس ها رو هم به دیتاست اضافه کنیم. برای مثال تبدیل های افاین مختلف (روتیشن های مختلف هر عکس برای مثال) روی عکس ها زده یا کراپ شده بخش هایی یا ... را به دیتاست می افزاییم. شهود: می دانیم برای مثال در یک تسک کلسیفیکیشن مثل تشخیص سگ: / توی عکس توقع داریم زاویه عکس تاثیری در سگ بودن نداشته باشد بدین ترتیب این رو به عنوان یک inductive bias با افزودن روتیشن های مختلف عکس به دیتاست به مدل اضافه میکنیم و توقع داریم مدل روی زاویه و ... تاثیر نگیره که اورفیت بشه و بعد روی دیتای جدید که لزوما همون زاویه نباشه غلط نشه. کلا هم با اضافه کردن این دیتاها موجب میشیم نسبت به اون نویز یا حالا پترن روباست تر رفتار کنه. از مثال روش های توی ویژن میشه افزودن دیتای اگمنتش شده تغییرهای crop, flip, rotation, translation, brightness, contrast, color, saturan اشاره کرد. کلا توی خیلی تسکا نیاز به چیزایی عین سلف سوپروایزد هستش مثل خیلی مدلا در کانترستو لرنینگ. دیتامون اساسا لیبیل نداره میایم اگمنت شده های یه عکس رو توی دیتا قرار میدیم و سعی میکنیم امبد شدشونو نزدیک کنیم. اینجا اصلا لیبیل نداریم واسه همین نیاز به دیتای این طوری داریم. کارهای دیگه ای هم میشه کرد مثلا بیایم ادورسریال اتک هایی به دیتا بزنی خود این اتک خورده ها رو هم به دیتا اضافه کنیم برای ترین مدل. یا اگه دیتا کم داریم از مدلای جنریتو مثل گن و ... استفاده کنیم که دیتا بیشتر بشود. البته لزوما همیشه دیتا اگمنتیشن موجب بهبود اورفیت نمیشه. حال دیتا اگمنتیشن های دیگه ای رو بررسی میکنیم. یکی از روش ها توی ام ال undersampling و oversampling هستش. یکی از روش ها رندوم سمپلینگه. حالا یا اورسمپلینگ میکنیم که میایم برای کلاس ماینورمون دیتاهای تکراری استفاده میکنیم که این احتمال اورفیت رو میتونه زیاد کنه به خاطر دیتاهای تکراری. ازونور میتونیم اندرسمپلینگ رندوم کنیم از کلاس میجرمون که این میتونه موجب از دست رفتن اینفورمیشن مهم بشه. روش های دیگه ای هم داریم مثل smote. اینجا میایم توی کلاس ماینور یه دیتارو میگیریم. k همسایه نزدیکش رو هم میگیریم. (k تعیین میکنیم معمولا ۵ مثلا) حالا یکی ازینا رو رندوم انتخاب میکنیم و روی پاره خط بین این دو تا (کلا دارم توی فضای این فیچر بحث میکنم) یه نقطه جدید میگیریم و به دیتا اضافه می کنیم. اما این متوهم اگه واریانس دیتای ماینورمون زیاد باشه و شباهت با میجره زیاد مشکل زاعه. همچنین کلا احتمال تولید دیتای نویزی و مشکلات اورفیت رو همچنان داره. روش های دیگه اِپسمپل مثل adasyn که بیسش همون smote هستش و یا اندرسمپلینگ مثل totem links هست که میاد اورلپ کلاسی میجره کم میکنه با حذف دیتاهایی که توتم لینکن و ...

یکم ساده لوحانه نگاه کنیم داریم که

$$\text{rademacher complexity: } \hat{\mathcal{R}}_N(\mathcal{H}) = \mathbb{E}_{\sigma} \left[\sup_{h \in \mathcal{H}} \left| \frac{1}{N} \sum_{i=1}^N \sigma_i h(x_i) \right| \right]$$

$$\mathcal{R}_N(\mathcal{H}) = \mathbb{E} \left[\hat{\mathcal{R}}_N(\mathcal{H}) \right]$$

حالا ارور جنرالیزیشنو داریم که :

$$\hat{\epsilon}_N(h) \leq \mathcal{R}_N(\mathcal{H}) + \mathcal{O} \left(\sqrt{\frac{\ln 1/\delta}{N}} \right)$$

با احتمال $1 - \delta$

خب با این نگاه با افزایش N که تعداد سمپلامونه داره ارور جنرالیزیشنمون رو کم میکنه و خب یعنی داره از اورفیت جلوگیری میکنه ولی مشکلی که هست حداقل متوذهای که ظاهر بیرونی رو حفظ میکنند برا ادم یعنی مثلا توی کلسیفیکیشن کلاسمون همونه استفاده زیاد از پرایار نالجمون درمورد دامنه دیتامون و ویژگی های تصویر مثلا داره خب iid نیستند و همیشه به چشم دیتای جدید صاف بهشون نگاه کرد که خب نگاه تئوری لرنینگ طور بالامون اساسا همونطور که گفتیم ساده لوحانه هستش. بررسی دقیق ریاضیش از توان من خارجه هعب. یه پیپریم دیدم بررسی کرده بود اساسا اگه میتونید رگولاریزیشن استفاده کنید و قصد کاهش اورفیت تاثیر خیلی بهتر و قابل کنترل تر (از لحاظ تئوری) از دیتا اگمنتیشنه.

PCA and Dimension Reduction

پی سی ای و روشای کاهش ابعاد لزوما موجب کاهش اورفیت نمیشن. درسته بعضا میتونه کاهش بده ولی افزایش ممکنه بده. وقتی یه مدل فیت میکنیم داریم کوواریانس بین فیچرها و اون تارگت رو می یابیم و توی هدفمون کواریانس خود فیچرها با هم نیستش و دنبال تاثیر تغییرات فیچر روی تارگتیم. حالا pca عملا به کواریانس با اون تارگته کاری نداره داره میاد واریانس فیچر رو بررسی میکنه و ازون حالا ایگن وکتورا رو میگیریم و ... جهت بیشترین واریانس فضای فیچرمون لزوما دلیلی نداره ربط به کواریانس بین فیچر و تارگت داشته باشه. ازینکه pca برنیم فیچرهای مهمون از دست بره تاثیرش. یه مثال :

H = height of cloud-base
Ts = surface temperature
Td = surface dewpoint

حالا داریم که چون کواریانس دو تا فیچر مثبتیه

$$\text{var}(Ts+Td) = \text{var}(Ts) + \text{var}(Td) + 2\text{Cov}(Ts, Td) > \text{var}(Ts) + \text{var}(Td) - 2\text{Cov}(Ts, Td) > \text{var}(Ts-Td)$$
 ولی H مستقیما به رطوبت ربط داره که از Ts-Td میاد. ولی اگه مثلا اینجا ما واریانس بیشتر و نگه داریم دیگه Ts+Td به ما نتیجه خوبی توی پیشبینی H نمیده.

نکته: کلا این که واریانس کما نویز نیستند و اهمیت دارند واسه همین مدلای دیمنشن ریداکشنمون علی الخصوص pca لزوما اینجا خوب عمل نمیکنن و میتونن موجب اورفیت بشن توی این صفحه هم بررسی شده:
<https://stats.stackexchange.com/questions/87198/low-variance-components-in-pca-are-they-really-just-noise-is-there-any-way-to/87231#87231>

بررسی تئوری: به هر حال تعداد دیمنشنا داره کم میشه واسه همین کامپلکسیتی مدل کم شده این داره واریانس مدل رو کم میکنه که باعث کاهش اورفیت میتونه بشه ولی همیشه تاثیرش روی بایاس هم صرف نظر کرد ممکنه کلا اینفورمیشن خویمون رو از دست بدیم. پس میتونه باعث کاهش اورفیت بشه ولی لزوما همیشه ممکنه زیادهش کنه. همچنین PCA و کلا متوذهای کاهش بعد میتونه موجب حذف نویز هم مقداری بشه که باعث کاهش اورفیت میشه. ولی ازونور میتونه اینفورمیشن مهم از دست بده یا ریلیشن بین فیچرهای مختلف رو مقداری بهم بریزه که درنتیجه اورفیت زیاد شه. اینها برای کلا روشای کاهش بعد صدق میکنه چه خطی عین pca چه عین t-sne و ...

$$-l = \sum_{i=1}^n \log P(y_i | x_i, w_0, w_1, w_2) - C w_j^2$$

$$\Rightarrow l = - \sum_{i=1}^n \log P(y_i | x_i, w_0, w_1, w_2) + C w_j^2 \quad \text{minimize}$$

* برای حل می‌توانیم مستقیماً به $\nabla l = 0$ برویم

$$\rightarrow l = - \sum_{i=1}^n (y_i \log P(1 | x_i, w_0, w_1, w_2) + (1-y_i) \log P(0 | x_i, w_0, w_1, w_2)) + C w_j^2$$

$$= \sum_{i=1}^n (-y_i \log \overset{\text{sigmoid}}{\sigma(w^T x_i)} - (1-y_i) \log (1 - \sigma(w^T x_i))) + C w_j^2$$

مشتق w بردار w است

$$w_k = w_k - \alpha \left(\sum_{i=1}^n (\sigma(w^T x_i) - y_i) x_j^i + 2 C w_j I(j=k) \right)$$

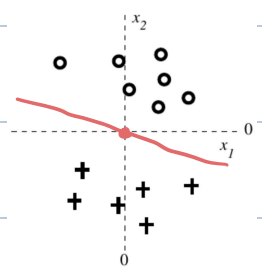
مشتق $k = j$ باشد اگر نه 0

طبق لاسر مشاهده شد هنگامی که بزرگ باشد ضریب کاهش w زیاد و مدل w را به صفر نزدیک می‌کند.

حالت 1) $j=0$:

در این مدل w به صفر میل می‌کند. پس $w^T X = w_2 x_2 + w_1 x_1 + 0$
یعنی خط بین کلاس‌ها فقط بر مبنای x_1 و x_2 می‌گذرد.

حال چون خط با w هم‌راهِ w داریم به همین جهت شده (موتور است) در این بهینه می‌رسد.
یعنی خطای کمترین:

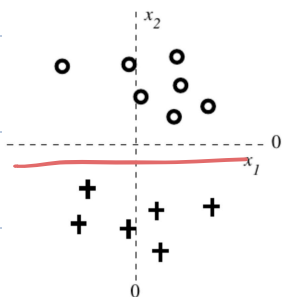


$$w^T X \approx w_2 x_2 + w_0$$

حالت دوم: $w_1 = 0$ و $w_2 \neq 0$

به ندرت مسئله از ویژگی x_1 می‌شود.

در فضای x_1 و x_2 پس خطی عمود بر محور x_2 است. خط شایسته بر حسب x_2 است. در مسئله خط با این ویژگی با شایسته‌ها داریم پس همچنان چون اردر کم است خط شایسته بهتر است.



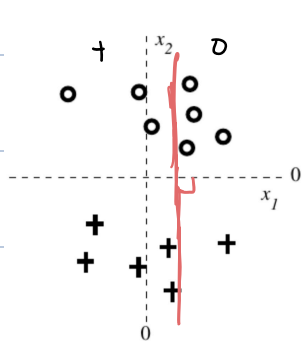
$$w^T X \approx w_1 x_1 + w_0 \quad w_2 = 0$$

مثلاً حالت دوم یکی این بار نسبت به x_2 نیست.

پس خطی عمود بر محور x_1 است و شایسته بر حسب x_1 است.

اما چنین خطی در فضا وجود ندارد که اردر 0 بدهد پس خطی تشریح در این مرحله بهتر می‌شود.

چون فضا مقایسه فقط روی x_1 است. انبار شده در x_1 تغییر شونده بهترین خط می‌باشد x_1 را گرفتیم.



(چنین صود فقط در یک)

بهترین حالت را با دست حالتها مختلف

ببینیم. اگر خط تفرز بگیریم

$$\frac{5+4}{5+4+3+1} = \frac{9}{13}$$

اکثریت پس در دیتا ترین داریم.

در یک خط بیان می‌شود. افزایش یافت از دیتا قبلی. چون x_1 تنها برای جداسازی اطلاعات کافی ندارد.

(5)

$$p(y|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2} (y-\mu)^2\right\}$$

الف)

حالت

سادس نسبی :

$$p(y|\eta) = b(y) \exp\{\eta^T T(y) - a(\eta)\}$$

$$\frac{1}{\sqrt{2\pi}} \frac{1}{\sigma} \exp\left\{-\frac{1}{2\sigma^2} (y^2 - 2y\mu + \mu^2)\right\}$$

$$\frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2\sigma^2} y^2 + \frac{y\mu}{\sigma^2}\right\} \exp\left\{-\frac{\mu^2}{2\sigma^2} - \ln(\sigma)\right\}$$

$$T(y) = \begin{pmatrix} y \\ y^2 \end{pmatrix} \quad b(y) = \frac{1}{\sqrt{2\pi}}$$

$$\eta = \begin{pmatrix} \frac{\mu}{\sigma^2} \\ -\frac{1}{2\sigma^2} \end{pmatrix} \quad a(\eta) = \frac{\mu^2}{2\sigma^2} + \ln(\sigma) = -\frac{\eta_1^2}{4\eta_2} - \frac{1}{2} \ln(-2\eta_2)$$

حالت سادس نسبی :

σ غیر متغیر \Leftarrow فرمت نهاده شده است \Leftarrow درستی

$$\Rightarrow p(y|\mu) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} y^2\right) \exp\left(\mu y - \frac{1}{2} \mu^2\right)$$

$$\Rightarrow \eta = \mu, \quad T(y) = y, \quad a(\eta) = \frac{\eta^2}{2} = \frac{\eta^2}{2} \quad b(y) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right)$$

(ب) ابتدا لم های را اثبات می کنیم.

این ها از خواص اولیه می نامیم.

برابر دهیم $c(y) = e^{b(y)}$

ادبی: $\int P_\eta(y) dy = 1$ چون توزیع

$$\nabla_\eta \int \exp(\eta \cdot T(y) - \alpha(\eta) + c(y)) (T(y) - \nabla_\eta \alpha(\eta)) dy = 0$$

$$\Rightarrow \int \exp(\eta \cdot T(y) - \alpha(\eta) + c(y)) T(y) dy = \int \exp(\eta \cdot T(y) - \alpha(\eta) + c(y)) \nabla_\eta \alpha(\eta) dy$$

$$= \nabla_\eta \alpha(\eta) \underbrace{\int \exp(\eta \cdot T(y) - \alpha(\eta) + c(y)) dy}_1 = \nabla_\eta \alpha(\eta)$$

$$\Rightarrow \nabla_\eta \alpha(\eta) = \int \exp(\eta \cdot T(y) - \alpha(\eta) + c(y)) T(y) dy =$$

$$\int P_\eta(y) T(y) dy = E(T(y)) \quad y \sim P_\eta(y)$$

ادبی اگر می بیند

دومی:

$$\frac{\partial \alpha}{\partial \eta_i} = \int \exp(\eta \cdot T(y) - \alpha(\eta) + c(y)) T_i(y) dy$$

یک بر سر متوجه می شویم

$$\begin{aligned} \frac{\partial^2 \alpha}{\partial \eta_i \partial \eta_j} &= \int \exp(\eta \cdot T(y) - \alpha(\eta) + c(y)) \left(T_j(y) - \frac{\partial \alpha(\eta)}{\partial \eta_j} \right) T_i(y) dy \\ &= \int \exp(\eta \cdot T(y) - \alpha(\eta) + c(y)) T_i(y) T_j(y) dy - \frac{\partial \alpha(\eta)}{\partial \eta_j} \int \exp(\eta \cdot T(y) - \alpha(\eta)) T_i(y) dy \\ &= E(T_i(y) T_j(y)) - E(T_i(y)) E(T_j(y)) = \text{Cov}(T_i(y), T_j(y)) \end{aligned}$$

$$\Rightarrow D_{\eta}^2 \alpha = \text{Cov}(T(y))$$

ایضا مندرجہ ذیل اثبات میں ہم براہِ خود ملاحظہ فرمائیں:

$$(۱۳) \quad \frac{\partial A}{\partial \eta^T} = E(T(X)) \quad \text{درجہ اول کے}$$

$$\frac{\partial A}{\partial \eta^T} = \frac{\int T(x) e^{\eta^T T(x)} h(x) v(dx)}{\int e^{\eta^T T(x) - A(\eta)} h(x) v(dx)} = E(T(X))$$

$$(۱۴) \quad \frac{\partial^2 A}{\partial \eta \partial \eta^T} = \text{Var}(T(X))$$

$$\frac{\partial^2 A}{\partial \eta \partial \eta^T} = \int T(x) (T(x) - \frac{\partial A}{\partial \eta^T})^T e^{\eta^T T(x) - A(\eta)} h(x) v(dx) =$$

$$\int T(x) (T(x) - E(T(X)))^T e^{\eta^T T(x) - A(\eta)} h(x) v(dx) = E(T(X) T(X)^T) -$$

$$E(T(X)) E(T(X))^T = \text{Var } T(X)$$

له ۳) داریم در حالت گاوسیانه

$$\frac{\partial a(\eta)}{\partial \eta_1} = \frac{-2\eta_1}{4\eta_2} = \frac{-\eta_1}{2\eta_2} = \frac{-\frac{\mu}{\sigma^2}}{-\frac{1}{\sigma^2}} = \mu$$

$$\frac{\partial^2 a(\eta)}{\partial \eta_1 \partial \eta_2} = \frac{\partial}{\partial \eta_1} \left(\frac{-\eta_1}{2\eta_2} \right) = -\frac{1}{2\eta_2} = -\frac{1}{2 \times \left(\frac{1}{2\sigma^2} \right)} = -\sigma^2 \quad (4)$$

حل واحد شده سردرآهون شده است. لهما در اول اون درآهون شده است.

$$\ell(\theta) = - \sum_{i=1}^m \log b(y_i) + x_i^T \theta T(y_i) - a(\theta^T x_i)$$

$$\frac{\partial \ell}{\partial \theta} = - \sum_{i=1}^m T(y_i) x_i - \frac{\partial a}{\partial \theta} (\theta^T x_i) x_i$$

مبداً ثابت کردم دو لم می آید - سرچش رو

$$\frac{\partial \ell}{\partial \theta} = - \sum_{i=1}^m (T(y_i) - \frac{\partial a}{\partial \theta}(\eta)) x_i$$

$$\frac{\partial^2 \ell}{\partial \theta_j \partial \theta_k} = - \frac{\partial}{\partial \theta_j} \sum_{i=1}^m (T(y_i) - \frac{\partial a}{\partial \theta_k}(\eta)) x_k^i =$$

$$- \sum_{i=1}^m \left(- \frac{\partial a}{\partial \theta_j \partial \theta_k}(\eta) \frac{\partial \eta}{\partial \theta_j} x_k^i + 0 \right) = \sum_{i=1}^m \frac{\partial a}{\partial \theta_j \partial \theta_k}(\eta) x_j^i x_k^i = H_{cov}$$

زیر محاسبه شده: $\sigma, \nabla a = E(\psi)$ (مبداً ثابت کردم)

$\frac{\partial a}{\partial \theta_j} = \mu$

$\Rightarrow \frac{\partial a}{\partial \theta_j \partial \theta_k} = \frac{\partial \mu}{\partial \theta_k}$

$\Rightarrow H = X^T W X$

$W = \begin{bmatrix} \frac{\partial \mu}{\partial \theta_1} & \dots \\ \vdots & \ddots \\ \vdots & \dots \end{bmatrix}$

cov

$$\textcircled{1} \quad \nabla a(\eta) = E(T(y))$$

کتاب سر خوانی که از قبل داریم :

$$\textcircled{2} \quad \nabla \nabla^T a(\eta) = \text{Cov}(T(y))$$

حالت کامل² کلی³
پارامترهای ترسیده، که حالت خاص است
هم ثابت کردم

$$l = \sum_{i=1}^m [\log b(y^i) + \eta^T T(y^i) - a(\eta)]$$

$$\nabla_{\theta_j} l = \sum_{i=1}^m x^i (T_j(y^i) - E(T_j(y^i) | \eta))$$

① از

$$\eta = \theta^T x$$

$$T = (T_1(y), T_2(y), \dots)^T$$

$$\nabla_{\theta_j} \nabla_{\theta_k}^T l = \sum_{i=1}^m x^i x^{iT} \left(-\text{Cov}[T_{j,k}(y^i) | \eta] + \delta_{jk} \eta''_j (T_k(y^i) - E[T_k(y^i) | \eta]) \right)$$

②

در سوال ما چه می باشد

* حالت کلی ساده تر: در حالت $\eta = T(y)$ داریم که ماتریس:

$$f(E(y)) = \eta$$

هر چند خاص و مقصوره

$$\nabla_{\theta} l = X^T (y - f^{-1}(X\theta))$$

$$\nabla \nabla^T l = -X^T W X$$

ماتریس واریانس-کواریانس

$$W_{ij} = \text{var}[y^i | \eta]$$

$$l = - \sum_{i=1}^m \log b(y^i) + \eta^T T(y^i) - a(\eta)$$

$$\nabla_{\theta_j} l = - \sum_{i=1}^m x^{iT_j(y^i)} - \frac{\partial}{\partial \theta_j} a(\eta)$$

↓

$$\theta: r \times 2$$

$$X: r \times n$$

$$\eta^T T(y^i) = x^{iT} \theta T(y^i) = (x_1^i, \dots, x_r^i) \begin{pmatrix} \theta_1^1 & \theta_1^2 \\ \vdots & \vdots \\ \theta_r^1 & \theta_r^2 \end{pmatrix} \begin{pmatrix} T_1 \\ T_2 \end{pmatrix}$$

$$= (x_1^i \theta_1^1, \dots, x_r^i \theta_r^1, x_1^i \theta_1^2, \dots, x_r^i \theta_r^2) \begin{pmatrix} T_1 \\ T_2 \end{pmatrix}$$

$$= x_1^i \theta_1^1 T_1 + \dots + x_r^i \theta_r^1 T_1 + x_1^i \theta_1^2 T_2 + \dots + x_r^i \theta_r^2 T_2$$

$$\Rightarrow \frac{\partial \eta^T T(y^i)}{\partial \theta_j} = x^i T_j$$

$$\Rightarrow \nabla_{\theta_j} l = - \sum_{i=1}^m \left(x^i T_j(y^i) - \frac{\partial a(\eta)}{\partial \eta} \cdot \frac{\partial \eta}{\partial \theta_j} \right)$$

$$= - \sum_{i=1}^m x^i (T_j(y^i) - E(T_j(y^i)))$$

$$\frac{\partial l}{\partial \theta_1} = - \sum_{i=1}^m x^i (y^i - \mu)$$

$$\hookrightarrow E(T_1(y^i)) = E(y^i)$$

در اینجا:

$$\frac{\partial l}{\partial \theta_2} = - \sum_{i=1}^m x^i (y^{i2} - E(y^{i2}))$$

و اینجاست

حال دوباره مشتق بگیریم :

$$\frac{\partial l}{\partial \theta_j \partial \theta_k} = + \sum_{i=1}^m x^i x^{iT} \text{cov}(T_j(y^i), T_k(y^i)) \rightarrow H$$

↑

$$\frac{\partial a}{\partial \theta_j \partial \theta_k} = \frac{\partial}{\partial \theta_j} \left(\frac{\partial a(\eta)}{\partial \eta_k} \cdot \frac{\partial \eta_k}{\partial \theta_k} \right) = \frac{\partial}{\partial \theta_j} \left(\frac{\partial a(\eta)}{\partial \eta_k} \right) \frac{\partial \eta_k}{\partial \theta_k} + \dots$$

$$= \underbrace{\frac{\partial a(\eta)}{\partial \eta_j \partial \eta_k}}_{\text{cov}} \frac{\partial \eta_j}{\partial \theta_j} \frac{\partial \eta_k}{\partial \theta_k} \rightarrow \text{cov}(T_j(y^i), T_k(y^i)) x^i \cdot x^{iT}$$

به سبب اینکه هر دو را داریم $\rightarrow (-1) \times (-1) = 1$

$$\Rightarrow H = X^T S X$$

$$S = \text{cov}[T(y)]$$

همیشه که قبلاً هر صفت بردار.

حال S را می توان بر حسب η ها نوشت.

$$E[(T_1(y) - E(T_1(y)))(T_2(y) - E(T_2(y)))]$$

$$= E[(y - \mu)(y^2 - E(y^2))] = E[(y - \mu)(y^2 - \mu^2 - \sigma^2)]$$

↙

$$E(y^2) = E(y)^2 + E((y-E(y))^2) \Rightarrow$$

$$\Rightarrow E(y^3 - y\mu^2 - y\sigma^2 - \mu y^2 + \mu^3 + \mu\sigma^2)$$

$$= E(y^3) - \cancel{\mu^3} - \cancel{\mu\sigma^2} - \mu(\sigma^2 + \mu^2) + \cancel{\mu^3} + \cancel{\mu\sigma^2} = E(y^3) - \mu(\sigma^2 + \mu^2)$$

$$\text{Var}(T_1(y)) = E((y - E(y))^2) = \sigma^2$$

$$\text{Var}(T_2(y)) = E((y^2 - E(y^2))^2) = E((y^2 - \mu^2 - \sigma^2)^2) =$$

$$E(y^4) + \cancel{\mu^4} + \cancel{\sigma^4} - 2\mu^2(\sigma^2 + \mu^2) - 2\sigma^2(\mu^2 + \sigma^2) + 2\mu^2\sigma^2 =$$

$$E(y^4) - \mu^4 - \sigma^4 - 2\mu^2\sigma^2 = E(y^4) - (\mu^2 + \sigma^2)^2$$

حالت " μ و σ اینها بر حسب η_1 و η_2 (رابطه) $X^T X$ هم که از η ها داریم.

$$\eta = \begin{pmatrix} \frac{\mu}{\sigma^2} \\ \frac{1}{2\sigma^2} \end{pmatrix} \Rightarrow \frac{\eta_1}{-2\eta_2} = \mu$$

$$\frac{1}{2\eta_2} = \sigma^2$$

بر حسب η_1 و η_2
 معادله

بر حسب η_1 و η_2 ، X ها $X^T \text{CON}(T) X$ در مورد η