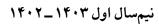
آمار و احتمال مهندسی





دانشکدهی مهندسی کامپیوتر مکتر نجفی

تمرین سری ششم موعد تحویل: ۱۶ دی

- مهلت ارسال پاسخ تا ساعت ۲۳:۵۹ روز مشخص شده است.
- همکاری و همفکری شما در انجام تمرین مانعی ندارد اما پاسخ ارسالی هر کس حتما باید توسط خود او نوشته شده باشد.
- در صورت هم فکری و یا استفاده از هر منابع خارج درسی، نام هم فکران و آدرس منابع مورد استفاده برای حل سوال مورد نظر را ذکر کنید.
 - لطفا تصویری واضح از پاسخ سوالات نظری بارگذاری کنید. در غیر این صورت پاسخ شما تصحیح نخواهد شد.

مسئلهی ۱. (۱۴ نمره)

مادربزرگ مهدی که چایخور قهاری است، ادعا میکند که می تواند تشخیص دهد که ابتدا چای در فنجان ریخته شده یا آب جوش! اما مهدی که به تازگی در کلاس آمار و احتمال مهندسی با آزمون فرض آشنا شده تصمیم می گیرد که ادعای مادربزرگ را بررسی کند. به همین منظور او Λ فنجان چای آماده میکند و در Λ تای آنها ابتدا چای و سپس آب جوش، و در Λ تای دیگر ابتدا آب جوش و سپس چای را اضافه میکند. حال او Λ فنجان به طور تصادفی روبروی مادربزرگش قرار می دهد تا او Λ فنجانی که ابتدا در آنها چای ریخته شده است را تشخیص دهد.

آ) برای آزمایش طراحی شده فرض صفر و جایگزین را بنویسید.

ب) آزمون فیشر را اجرا کرده و مقادیر p-value را بنویسید.

حل. آ) فرض صفر: مادربزرگ نمی تواند فنجانهایی که ابتدا در آنها چای و سپس آب جوش ریخته شده است را تشخیص دهد.

فرض جایگزین: مادربزرگ می تواند فنجانهایی که ابتدا در آنها چای و سپس آب جوش ریخته شده است را تشخیص دهد.

ب) امکان دارد مادربزرگ ۰،۲،۴،۶ و یا ۸ فنجان را به درستی تشخیص دهد. احتمال هر یک از این حالات را بدست می آوریم:

$$p_{\bullet} = \frac{\binom{\mathfrak{f}}{\bullet}\binom{\mathfrak{f}}{\bullet}}{\binom{\Lambda}{\mathfrak{f}}} = \frac{1}{\mathsf{V}^{\bullet}}, p_{\mathsf{T}} = \frac{\binom{\mathfrak{f}}{\bullet}\binom{\mathfrak{f}}{\bullet}}{\binom{\Lambda}{\mathfrak{f}}} = \frac{1}{\mathsf{V}^{\bullet}}, p_{\mathsf{F}} = \frac{\binom{\mathfrak{f}}{\bullet}\binom{\mathfrak{f}}{\bullet}}{\binom{\Lambda}{\mathfrak{f}}} = \frac{\mathsf{V}^{\mathsf{F}}}{\mathsf{V}^{\bullet}}, p_{\mathsf{F}} = \frac{\binom{\mathfrak{f}}{\bullet}\binom{\mathfrak{f}}{\bullet}}{\binom{\Lambda}{\mathfrak{f}}} = \frac{1}{\mathsf{V}^{\bullet}}, p_{\mathsf{A}} = \frac{1}{\mathsf{V}^{\bullet}}, p_{\mathsf$$

حال مقادير p-value را بدست مي آوريم:

p-value,
$$= p \cdot + p_{\uparrow} + p_{f} + p_{f} + p_{h} = 1$$

p-value, $= p_{\uparrow} + p_{f} + p_{f} + p_{h} = \frac{6}{7}$
p-value, $= p_{f} + p_{f} + p_{h} = \frac{67}{7}$
p-value, $= p_{f} + p_{h} = \frac{17}{7}$
p-value, $= p_{h} = \frac{1}{7}$

مسئلهی ۲. (۱۵ نمره)

اداره هواشناسی شهر تهران، ۴ دستگاه سنجش آلودگی هوا را در یک منطقه قرار داده است. فرض کنید شاخص آلودگی هوا در این منطقه ثابت است اما این دستگاهها دقیق نیستند و شاخص آلودگی را با کمی نویز گزارش میدهند. در یک روز نسبتا آلوده، مقادیر گزارش شده توسط این ۴ دستگاه به شرح زیر است:

آ) با کمک دادههای جمع آوری شده، یک بازه اطمینان ۹۹ درصدی برای شاخص آلودگی هوا در این منطقه ارائه دهید. ب) اگر شاخص آلودگی هوا بیش از ۱۵۰ باشد، هوا در شرایط ناسالم برای تمامی گروهها قرار میگیرد و مدارس و دانشگاهها تعطیل می شوند. برخی از صاحب نظران معتقند که از آنجایی که میانگین شاخص آلودگی هوا بیش از ۱۵۰ بوده، هوا ناسالم است و مدارس و دانشگاهها می بایست تعطیل شوند. از طرفی گروهی دیگر از صاحب نظران بر این باورند که میانگین بدست آمده به طور محسوس و معناداری از ۱۵۰ بیشتر نبوده و مدارس و دانشگاهها نباید تعطیل شوند. با توجه به اینکه p-value قابل قبول برای تعطیلی مدارس حداکثر ۱/۰ است، بررسی کنید که آیا مدارس تعطیل شوند یا خیر.

 ϕ) برای کاهش خطای نوع اول و دوم باید مقدار p-value قابل قبول را افزایش دهیم یا کاهش چرا

حل.

آ) از آنجایی که توزیع عدد گزارش شده توسط هر دستگاه از توزیع نرمال با واریانس نامشخص میآید، میتوانیم توزیع میانگین را t Student's در نظر بگیریم. با توجه به این نکته برای بازه اطمینان خواهیم داشت:

$$-t_{/ \cdots \Diamond} \leqslant \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \leqslant t_{/ \cdots \Diamond} \to \bar{x} - \frac{s}{\sqrt{n}} t_{/ \cdots \Diamond} \leqslant \mu \leqslant \bar{x} + \frac{s}{\sqrt{n}} t_{/ \cdots \Diamond}$$

با توجه به اعداد گزارش شده و از آنجایی که m-1=n-1 داریم:

$$\bar{x} = \frac{\text{12T} + \text{14A} + \text{12T} + \text{12T}}{\text{4}} = \text{121/2}, \quad s = \sqrt{\frac{\text{1/2}^{\text{4}} + \text{4/2}^{\text{4}} + \text{4/2}^{\text{4}}}{\text{4}}} \approx \text{4/4A} \cdot \text{4}$$

چون t_{\prime} ۹۹٪ را بدست آورد: t_{\prime} است، می توان بازه اطمینان /۹۹ را بدست آورد:

$$1\Delta 1/\Delta \pm 9/9\Delta$$

t-value و فرض مقابل ۱۵۰ استفاده می نیم. مقدار $\mu = 10۰$ با فرض صفر ۱۵۰ با زآزمون تک نمونه ای $\mu = 10۰$ با فرض صفر ۱۵۰ با فرض مقابل ۱۵۰ با نیم محاسبه می کنیم:

$$rac{ar{x} - \mu}{rac{s}{\sqrt{n}}} pprox \ {
m N/FTA}$$

پس نمی توانیم فرض صفر را رد کنیم و در نتیجه مدارس و دانشگاهها تعطیل نمی شوند.

پ) خطای نوع اول برابر با سطح اهمیت آزمون است بنابراین برای کاهش خطای نوع اول باید مقدار عددی سطح اهمیت را کاهش دهیم. خطای نوع دوم هنگامی اتفاق میافتد که فرض صفر غلط باشد و ما به اشتباه نتوانیم آن را رد کنیم. افزایش مقدار عددی سطح اهمیت رد کردن فرض صفر را ساده کرده و احتمال آن را بیشتر میکند (در حالی که تاثیری روی درست یا غلط بودن فرض صفر در عالم واقعیت ندارد!) بنابراین افزایش مقدار عددی سطح اهمیت باعث کاهش احتمال خطای نوع دوم می شود.

مسئلهی ۳. (۱۴ نمره)

متغیر تصادفی پیوسته X با تابع چگالی احتمالاتی زیر را در نظر بگیرید:

$$f_X(x) = \begin{cases} \mathbf{Y}x^{\mathbf{Y}} & x \in [\mathbf{\cdot}, \mathbf{1}] \\ \mathbf{\cdot} & o.w. \end{cases}$$

اگر X = Y باشد، تخمین MAP متغیر تصادفی X به شرط Y = X را بدست آورید.

حل. با توجه به تعریف توزیع هندسی برای $y=1,7,\cdots$ داریم:

$$P_{Y|X}(y|x) = x(1-x)^{y-1}$$

پس با جایگزینی $y = \mathbf{r}$ خواهیم داشت:

$$P_{Y|X}(\Upsilon|x) = x(\Upsilon-x)^{\Upsilon}$$

اکنون باید به دنبال $x \in [*, 1]$ بگردیم که مقدار زیر را بیشینه کند:

$$P_{Y|X}(y|x)f_X(x) = x(1-x)^{\Upsilon} \cdot \Upsilon x^{\Upsilon} = \Upsilon x^{\Upsilon}(1-x)^{\Upsilon}$$

برای یافتن بیشینه این تابع از آن برحسب x مشتق گرفته و حاصل را برابر با صفر قرار می دهیم.

$$\frac{d}{dx} \left[\mathbf{\tilde{r}} x^{\mathbf{\tilde{r}}} (\mathbf{1} - x)^{\mathbf{\tilde{r}}} \right] = \mathbf{\tilde{q}} x^{\mathbf{\tilde{r}}} (\mathbf{1} - x)^{\mathbf{\tilde{r}}} - \mathbf{\tilde{r}} x^{\mathbf{\tilde{r}}} (\mathbf{1} - x) = \mathbf{\tilde{r}} \to \begin{cases} x = \mathbf{\tilde{r}} \\ x = \frac{\mathbf{\tilde{r}}}{\mathbf{\tilde{o}}} \end{cases} \max_{x = \mathbf{\tilde{r}}}$$

با بررسی مقادیر فوق تخمینگر MAP را بدست می آوریم و داریم:

$$\hat{x}_{MAP} = \frac{\Upsilon}{\Delta}$$

 \triangleright

مسئلهی ۴. (۲۰ نمره)

مسئله ی رگرسیون خطی ساده را در نظر بگیرید که در آن ورودی های $\{(x_i,y_i)\}_{i=1}^n$ داده شده است. مقادیر $\{x_i\}_{i=1}^n$ یقینی و مشخص هستند. اما به ازای هر x_i مقدار y_i از طریق رابطه ی زیر به دست می آید:

$$y_i = \beta_i + \beta_1 x_i + \epsilon_i$$

که $\epsilon_i \sim \mathcal{N}(\, {f \cdot}\,, \sigma^{\, {f r}})$ که $\epsilon_i \sim \mathcal{N}(\, {f \cdot}\,, \sigma^{\, {f r}})$ که $\epsilon_i \sim \mathcal{N}(\, {f \cdot}\,, \sigma^{\, {f r}})$

آ) اثبات کنید که تخمین بیشینه درستنمایی دو پارامتر β و β معادل انتخاب مقادیری برای β و β است که میانگین مربعات خطا را کمینه میکند.

ب) اثبات کنید که تخمینهای بدست آمده در بخش پیشین نااُریب بوده و از توزیعهای زیر پیروی میکنند:

$$\hat{\beta}_1 \sim \mathcal{N}\left(\beta_1, \frac{\sigma^{\Upsilon}}{\sum_i (x_i - \bar{x})^{\Upsilon}}\right), \quad \hat{\beta}_* \sim \mathcal{N}\left(\beta_*, \frac{\sigma^{\Upsilon} \sum_i x_i^Y}{n \sum_i (x_i - \bar{x})^{\Upsilon}}\right)$$

 ψ) بررسی کنید که آیا تخمینگر بیشینه درستنمایی عضوی از خانواده ی خطی تخمینگرهای زیر است یا نه؟ اگر هست رابطه ی γ_i را برحسب دادههای ورودی بدست آورید.

$$\tilde{eta}_{1} = rac{\sum \gamma_{i} y_{i}}{\sum \gamma_{i} x_{i}}$$
 such that $\sum_{i} \gamma_{i} = \cdot$

ت) اثبات کنید هر تخمینگری که عضو خانواده فوق است نااُریب میباشد.

ث) اثبات کنید به ازای هر انتخابی از مقادیر γ_i در خانواده فوق داریم $Var(\hat{eta}_1) \leqslant Var(\hat{eta}_1) \leqslant Var(\hat{eta}_1)$ نتیجه بدست آمده را توضیح دهید.

حل.

(Ĩ

$$f_{L} = f(y_{1}, \dots, y_{n} | \beta_{\bullet}, \beta_{1}) = \prod_{i=1}^{n} p(y_{i} | \beta_{\bullet}, \beta_{1}) = \prod_{i=1}^{n} \frac{1}{\sqrt{\sqrt{\gamma}\pi\sigma}} \exp\left(-\frac{(y_{i} - \beta_{\bullet} - \beta_{1}x_{i})^{\gamma}}{\gamma\sigma^{\gamma}}\right)$$

$$\rightarrow f_{L} = (\gamma\pi)^{\frac{-n}{\gamma}} \sigma^{-n} \exp\left(-\frac{\sum_{i} (y_{i} - \beta_{\bullet} - \beta_{1}x_{i})^{\gamma}}{\gamma\sigma^{\gamma}}\right)$$

$$\mathcal{L} = \log(f_{L}) = c - \frac{1}{\gamma\sigma^{\gamma}} \sum_{i} (y_{i} - \beta_{\bullet} - \beta_{1}x_{i})^{\gamma}$$

$$\frac{\partial \mathcal{L}}{\partial \beta_{\bullet}} \propto \sum_{i} (y_{i} - \beta_{\bullet} - \beta_{1}x_{i}) = \cdot \rightarrow \beta_{\bullet} = \bar{y} - \beta_{1}\bar{x}$$

$$\frac{\partial \mathcal{L}}{\partial \beta_{1}} \propto \sum_{i} (y_{i} - \beta_{\bullet} - \beta_{1}x_{i})x_{i} = \cdot \rightarrow \beta_{1} = \frac{\sum_{i} (x_{i} - \bar{x})(y_{i} - \bar{y})}{\sum_{i} (x_{i} - \bar{x})^{\gamma}}$$

ب) رابطه را برای $\hat{\beta}_1$ اثبات میکنیم و باقی روند برای $\hat{\beta}_1$ به همین صورت میباشد.

$$\hat{\beta}_{1} = \frac{\sum_{i} (x_{i} - \bar{x}) y_{i}}{\sum (x_{i} - \bar{x})^{\Upsilon}} = \sum_{i} \frac{(x_{i} - \bar{x})}{\sum (x_{i} - \bar{x})^{\Upsilon}} y_{i}$$

از آنجایی که y_i خود یک توزیع گاوسی دارد، $\hat{\beta}_1$ نیز جمع تعدادی جملهی گاوسی میشود و در نتیجه توزیع آن نیز گاوسی میشود. پس داریم:

$$\mathbb{E}[\hat{\beta}_{\mathbf{1}}] = \sum_{i} \frac{(x_{i} - \bar{x})}{\sum (x_{i} - \bar{x})^{\mathsf{T}}} (\beta_{\mathbf{1}} + \beta_{\mathbf{1}} x_{i}) = \beta_{\mathbf{1}}$$

$$Var(\hat{\beta_1}) = \sum_i \left(\frac{(x_i - \bar{x})}{\sum (x_i - \bar{x})^{\mathsf{Y}}}\right)^{\mathsf{Y}} Var(y_i) = \frac{\sigma^{\mathsf{Y}}}{\sum (x_i - \bar{x})}$$

پ) بدیهیست با جایگذاری $\gamma_i = x_i - \bar{x}$ به همیان تخمین کمینهی مربعات خطا می رسیم.

ت)

$$\mathbb{E}[\tilde{\beta_{\mathsf{I}}}] = \sum_{i} \frac{\gamma_{i}}{\sum \gamma_{i} x_{i}} \mathbb{E}[y_{i}] = \sum_{i} \frac{\gamma_{i}}{\sum \gamma_{i} x_{i}} (\beta_{\mathsf{I}} + \beta_{\mathsf{I}} x_{i}) = {}^{\mathsf{I}} + \frac{\sum \gamma_{i} x_{i}}{\gamma_{i} x_{i}} \beta_{\mathsf{I}} = \beta_{\mathsf{I}}$$

ث)

$$\frac{\gamma_i}{\sum \gamma_i x_i} := d_i \to Var(\tilde{\beta}_1) = \sigma^{\Upsilon} \sum d_i^{\Upsilon}$$

$$\sigma^{\Upsilon} \sum d_i^{\Upsilon} \sum (x_i - \bar{x})^{\Upsilon} \geqslant \sigma^{\Upsilon} \Big(\sum d_i (x_i - \bar{x}) \Big)^{\Upsilon} = \sigma^{\Upsilon}$$

$$\to \sigma^{\Upsilon} \sum d_i^{\Upsilon} = Var(\tilde{\beta}_1) \geqslant \frac{\sigma^{\Upsilon}}{\sum (x_i - \bar{x})^{\Upsilon}} = Var(\hat{\beta}_1)$$

در این خانواده از تخمینگرها همه نااُریب هستند و تخمین بیشینه درستنمایی در میان آنها از کمترین واریانس برخوردار است پس میتوان ادعا کرد که این تخمین بهترین است.

مسئلهی ۵. (۱۵ نمره)

 $\epsilon \sim N(ullet, \sigma^{\Upsilon})$ یک مدل رگرسیون خطی به شکل $Y = \beta \cdot + \beta_1 X_1 + \beta_7 X_7 + \epsilon$ در شرایطی که خطای آن از توزیع نرمال پیروی میکند داریم.

نابید. β ۰, β ۰, β ۱ ضرایب MLE

حل.

$$L(\beta, \beta_{1}, \beta_{1}, \sigma^{7}) = (\mathbf{Y}\pi\sigma^{7})^{-n/7} \exp\left(-\sum_{i=1}^{n} \frac{(y_{i} - \beta, -\beta_{1}x_{i1} - \beta_{1}x_{i1})^{7}}{\mathbf{Y}\sigma^{7}}\right)$$

$$l(\beta, \beta_{1}, \beta_{1}, \sigma^{7}) = -\frac{n}{\mathbf{Y}} \log(\mathbf{Y}\pi\sigma^{7}) - \sum_{i=1}^{n} \frac{(y_{i} - \beta, -\beta_{1}x_{i1} - \beta_{1}x_{i1})^{7}}{\mathbf{Y}\sigma^{7}}$$

$$\frac{\partial l}{\partial \beta_{1}} = \sum_{i=1}^{n} (y_{i} - \beta, -\beta_{1}x_{i1} - \beta_{1}x_{i1}) = \bullet$$

$$\frac{\partial l}{\partial \beta_{1}} = \sum_{i=1}^{n} (y_{i} - \beta, -\beta_{1}x_{i1} - \beta_{1}x_{i1})x_{i1} = \bullet$$

$$\frac{\partial l}{\partial \beta_{2}} = \sum_{i=1}^{n} (y_{i} - \beta, -\beta_{1}x_{i1} - \beta_{1}x_{i1})x_{i1} = \bullet$$

$$\hat{\beta} \cdot = \bar{y} - \hat{\beta}_{1}\bar{x}\mathbf{1} - \hat{\beta}\mathbf{1}\bar{x}\mathbf{1}$$

$$\hat{\beta} \cdot = \frac{\sum_{i=1}^{n} (x_{i1}\mathbf{1} - \bar{x}\mathbf{1})(y_{i} - \bar{y})}{\sum_{i=1}^{n} (x_{i1}\mathbf{1} - \bar{x}\mathbf{1})(y_{i} - \bar{y})}$$

$$\hat{\beta} \cdot = \frac{\sum_{i=1}^{n} (x_{i1}\mathbf{1} - \bar{x}\mathbf{1})(y_{i} - \bar{y})}{\sum_{i=1}^{n} (x_{i1}\mathbf{1} - \bar{x}\mathbf{1})(y_{i} - \bar{y})}$$

$$\hat{\beta} \cdot = \frac{\sum_{i=1}^{n} (x_{i1}\mathbf{1} - \bar{x}\mathbf{1})(y_{i} - \bar{y})}{\sum_{i=1}^{n} (x_{i1}\mathbf{1} - \bar{x}\mathbf{1})(y_{i} - \bar{y})}$$

 \triangleright

مسئلهی ۶. (۱۰ نمره)

فرض کنید میانگین معدل همه ی دانشجویان فارغالتحصیل شده از دانشگاه صنعتی شریف در سال ۲۰۰۵ برابر ۳/۰۵ (بر اساس معیار GPA) بوده است. سازمان ثبت احوال قصد دارد سوابق ۲۰۰ دانشجوی فارغ التحصیل در سال ۲۰۲۲ را بررسی کند تا ببیند میانگین GPA تغییر کرده است یا خیر. در این نمونه ی ۱۰۰ تایی، میانگین معدل برابر ۲/۱۵ بوده و انحراف معیار برابر ۱/۵ می باشد.

الف. فرصیه های صفر و جایگزین را برای این تحقیق بیان کنید.

ب. قبول/رد فرضیه H. را در مقابل H_1 جایگزین را در سطح اهمیت $\alpha = 1.0$ بررسی کنید. فرض کنید حجم نمونه به اندازه کافی بزرگ بوده و توزیع میانگین نمونه از توزیع نرمال پیروی می کند. (منظور از سطح اهمیت، Significance Level

حل.

 $H_1: \mu \neq \Upsilon/ \cdot \Delta$: فرض جایگزین $H_1: \mu = \Upsilon/ \cdot \Delta$ الف) فرض صفر

<u>(</u>ب

$$t_{oss} = \sqrt{n} imes (ar{x} - \mu_{ullet})/\sigma = \sqrt{1 \cdot \cdot \cdot} imes (\mathrm{Y/VD} - \mathrm{Y/\cdot D})/\mathrm{V/D} = -\mathrm{Sol}$$

بنابراين

$$p-value = P(Z > |t_{oss}|) = \mathsf{Y}P(Z > \mathsf{P}) \approx \bullet$$

چون $p-value < \alpha$ پس فرض صفر رد می شود.

 \triangleright

مسئلهی ۷. (۱۲ نمره)

داده $\{w[i], x[i]\}_{i=1}^{N-1}$ مشاهده شده است. میدانیم $\{w[i], x[i]\}_{i=1}^{N-1}$ که $\{w[i], x[i]\}_{i=1}^{N-1}$ داده احتمال پیشین زیر است.

$$p(A) = \begin{cases} \lambda \exp(-\lambda A) & A \geqslant \bullet \\ \bullet & A < \bullet \end{cases}$$

A از w[i] از w[i] همچنین مقادیر w[i] از w[i] از w[i] میباشند. همچنین مقادیر w[i] از w[i] مستقل هستند. برآوردگر w[i] را روی w[i] بیابید.

حل. ابتدا میدانیم

$$p(w[n]) = \frac{1}{\sqrt{\mathbf{Y}\pi\sigma^{\mathbf{Y}}}} exp(-\frac{w^{\mathbf{Y}}[n]}{\mathbf{Y}\sigma^{\mathbf{Y}}})$$

در نتیجه با توجه به اینکه x[n] = A + w[n]، خواهیم داشت:

$$p(x[n]|A) = \frac{1}{\sqrt{\mathbf{Y}\pi\sigma^{\mathbf{Y}}}}exp(-\frac{(x[n]-A)^{\mathbf{Y}}}{\mathbf{Y}\sigma^{\mathbf{Y}}})$$

برای محاسبهی Liklihood خواهیم داشت

$$p(x|A) = \prod_{i=1}^{N-1} p(x[n]|A) = \frac{1}{\sqrt{\Upsilon \pi \sigma^{\Upsilon}}} exp(-\frac{1}{\Upsilon \sigma^{\Upsilon}} \sum_{n=1}^{N-1} (x[n] - A)^{\Upsilon})$$

اگر تخمین MAP متغیر A را \hat{A}_{MAP} بنامیم، خواهیم داشت:

$$\hat{A}_{MAP} = \underset{A}{\operatorname{arg\,max}} \ p(A|x) = \underset{A}{\operatorname{arg\,max}} \ p(x|A)p(A) = \underset{A}{\operatorname{arg\,max}} \ (\log p(x|A) + \log p(A))$$

اگر توزیع A و x|A را در این رابطه قرار دهیم، خواهیم داشت:

$$\hat{A}_{MAP} = \underset{A}{\operatorname{arg\,max}} \ -\frac{1}{\mathbf{Y}\sigma^{\mathbf{Y}}} \sum_{i=1}^{N-1} (x[n] - A)^{\mathbf{Y}} - \lambda A$$

اگر تابع هدف را f(A) نامیده و مشتق آن را برابر • قرار دهیم، خواهیم داشت:

$$\frac{df(A)}{dA} = \frac{1}{\sigma^{\gamma}} \sum_{n=1}^{N-\gamma} (x[n] - A) - \lambda = \cdot \leftrightarrow A = \bar{x} - \frac{\lambda \sigma^{\gamma}}{N}$$

که $\bar{x} = \sum_{n=1}^{N-1} (x[n])$ که انتجابی که

$$\frac{d^{\mathsf{Y}}f(A)}{dA^{\mathsf{Y}}} = -\frac{N}{\sigma^{\mathsf{Y}}} - \lambda < \bullet$$

مراحل بالا f(A) را بیشینه میکند.

در نهایت با در نظر گرفتن اینکه احتمال رخداد مقادیر منفی برای A مساوی \bullet است، خواهیم داشت:

$$\hat{A}_{MAP} = max(\bullet, \bar{x} - \frac{\lambda \sigma^{\Upsilon}}{N})$$

 \triangleright

موفق باشيد :)