

چکیده

در این پایان نامه مسئله برآورد پارامترهای رگرسیون در مدل رگرسیون خطی چندگانه هنگامی که هم خطی چندگانه موجود باشد، بررسی می شود. تحت فرض نرمال، سه برآوردگر بیز تجربی شامل برآوردگر کمترین توان های دوم انقباضی، برآوردگر بیز تجربی سلسله مراتبی انقباضی و برآوردگر بیز تجربی تجزیه شده ارائه می شوند. برآوردگرها به ترتیب براساس روش مولفه های اصلی، برآوردگر کمترین توان های دوم دوگانه و با انتخاب پیشین های متفاوت برای پارامترهای مدل رگرسیونی به دست می آیند.

تاکنون ثابت نشده است که برآوردگرهای پیشنهاد شده از نظر همگرایی بهتر از برآوردگرهای کمترین توان های دوم هستند که در این پایان نامه این موضوع ثابت می شود. همچنین نشان داده می شود که این سه برآوردگر هنگامی که هم خطی چندگانه موجود باشد و از طریق شبیه سازی و مطالعات تجربی به دست آیند از نظر کارایی برآوردگرهای مفیدی هستند.

واژه های کلیدی: رگرسیون چندگانه، هم خطی چندگانه، رگرسیون ستیغی، روش بیز تجربی، روش مولفه های اصلی، مینیماکس.

فهرست مندرجات

۱	تعاريف و پيش‌نيازها	۱
۱	۱.۱ مولف‌هاي اصلي جامعه	۱
۲	۲.۱ برآوردگرهاي بيز و مينيماکس	۲
۳	۳.۱ اصول مينيماکس و بيز	۳
۴	۴.۱ توزيع پيشين و توزيع پسین	۴
۶	۵.۱ برآوردگرهاي بيز	۶
۸	۶.۱ برآوردگرهاي مينيماکس	۸
۱۰	رگرسيون	۲

۱۰	مقدمه	۱.۲
۱۰	رگرسیون و ساختن مدل	۲.۲
۱۲	کاربردهای مدل‌های رگرسیون	۳.۲
۱۲	رگرسیون خطی ساده	۴.۲
۱۳	رگرسیون چندگانه خطی	۵.۲
۱۵	برآورد کمترین توان‌های دوم ضرایب رگرسیون چندگانه خطی	۶.۲
۱۸		هم‌خطی چندگانه	۳
۱۸	مقدمه	۱.۳
۱۹	منابع هم‌خطی چندگانه	۲.۳
۲۱	آثار هم‌خطی چندگانه	۳.۳
۲۴	ملاک‌های تشخیص هم‌خطی چندگانه	۴.۳
۲۴	محک ماتریس همبستگی	۱.۴.۳
۲۴	عوامل تورمی واریانس	۲.۴.۳

۲۵	تحلیل سیستم مقادیر ویژه	۳.۴.۳	
۲۷	رگرسیون ستیغی (RR)		۴
۲۷	مقدمه	۱.۴	
۳۲	رابطه برآوردگر ستیغی با برآوردگر بیز	۲.۴	
۳۲	روش‌های انتخاب k	۳.۴	
۳۶	انتخاب متغیر	۴.۴	
۳۶	برآوردگر رگرسیون ستیغی	۵.۴	
۴۲	برآورد بیز پارامترهای رگرسیون	۱.۵.۴	
۴۴	برآورد بیز تجربی c یا λ	۲.۵.۴	
۴۵	برآوردگر مولفه اصلی	۶.۴	
۴۸	اریبی و میانگین توان دوم خطای برآوردگر مولفه اصلی	۱.۶.۴	
۵۰	برآوردگر مولفه اصلی ستیغی	۷.۴	
۵۲	برآوردگرهای رگرسیون ستیغی بیز تجربی تحت هم‌خطی چندگانه		۵
۵۲	مقدمه	۱.۵	

۵۶	برآوردگرهای رگرسینون ستیغی بیز تجربی	۲.۵
۵۷	برآوردگر رگرسینون ستیغی بیز تجربی (EB)	۱.۲.۵
۵۹	برآوردگر رگرسینون ستیغی بیز تجربی سلسله مراتبی (HB)	۲.۲.۵
۶۲	برآوردگر رگرسینون ستیغی بیز تجربی تجزیه شده	۳.۲.۵
۶۳	مینیماکسیتی برآوردگرهای بیز تجربی	۳.۵
۶۳	شرایط کلی برای مینیماکسیتی	۱.۳.۵
۶۵	مینیماکسیتی برآوردگر رگرسینون ستیغی بیز تجربی	۲.۳.۵
۶۶	مینیماکسیتی برآوردگر رگرسینون ستیغی سازوار	۳.۳.۵
۶۸	برآوردگر رگرسینون ستیغی سازوار تعدیل یافته	۴.۳.۵
	مینیماکس برآوردگرهای بیز تجربی تجزیه شده و سلسله	۵.۳.۵
	مراتبی	۶۸
۷۱	نتایج و پیشنهادات	۴.۵
۷۳		برنامه‌های رایانه‌ای	A
۹۰		توضیحات	B
۹۹		واژه‌نامه انگلیسی به فارسی	C

لیست اشکال

فصل ۱

تعاریف و پیش‌نیازها

۱.۱ مولفه‌های اصلی جامعه

اگر بردار تصادفی $X' = [X_1, X_2, \dots, X_p]$ دارای ماتریس کوواریانس Σ و مقادیر ویژه $\lambda_1 \geq \lambda_2 \geq \dots, \lambda_p \geq 0$ باشد. با در نظر گرفتن ترکیبات خطی زیر:

$$Y_1 = l'_1 X = l_{11}X_1 + l_{21}X_2 + \dots + l_{p1}X_p$$

$$Y_2 = l'_2 X = l_{12}X_1 + l_{22}X_2 + \dots + l_{p2}X_p$$

.

.

$$Y_i = l'_i X = l_{1i}X_1 + l_{2i}X_2 + \dots + l_{pi}X_p$$

.

$$Y_p = l'_p X = l_{1p}X_1 + l_{2p}X_2 + \dots + l_{pp}X_p$$

داریم:

$$\text{Var}(Y_i) = l_i' \Sigma l_i \quad i = 1, \dots, p$$

و

$$\text{Cov}(Y_i, Y_k) = l_i' \Sigma l_k \quad i, k = 1, \dots, p$$

مولفه‌های اصلی X_1, X_2, \dots, X_p آن ترکیبات خطی ناهمبسته‌ی Y_1, Y_2, \dots, Y_p هستند که واریانس آن‌ها بیشترین مقدار ممکن را اتخاذ کند.

قضیه ۱-۱: اگر Σ ماتریس کوواریانس بردار تصادفی $X' = [X_1, X_2, \dots, X_p]$ باشد. همچنین اگر Σ دارای p زوج مقدار ویژه - بردار ویژه‌ی $(\lambda_1, e_1), (\lambda_2, e_2), \dots, (\lambda_p, e_p)$ باشد به طوری که $\lambda_1 \geq \lambda_2 \geq \dots, \lambda_p \geq 0$ و اگر مولفه اصلی i ام به صورت:

$$Y_i = e_i' X = e_{1i} X_1 + e_{2i} X_2 + \dots + e_{pi} X_p \quad i = 1, \dots, p$$

تعریف شود، آن‌گاه واریانس i امین مولفه اصلی ماکسیمم خواهد شد. با این انتخاب‌ها داریم:

$$\text{Var}(Y_i) = e_i' \Sigma e_i \quad i = 1, \dots, p$$

$$\text{Cov}(Y_i, Y_k) = e_i' \Sigma e_k \quad i \neq k$$

در صورتی که بعضی از λ_i ها برابر باشند، انتخاب e_i ها و در نتیجه Y_i ها یکتا نخواهند بود (الکساندر مود (۱۹۲۱)).

۲.۱ برآوردهای بیز و مینیماکس

می‌دانیم بهترین برآوردها، برآوردهای بیز است که به طور یکنواخت دارای کمترین مقدار مخاطره باشد. اما متأسفانه موقعیتی که تحت آن بهترین برآوردها وجود داشته باشد، پیش

نمی آید. یعنی برای یک مقدار ثابت θ ، برآوردگری با کمترین مخاطره را می توان یافت، اما این بهترین برآوردگر برای مقادیر مختلف θ تغییر می کند. بنابراین هیچ برآوردگری را نمی توان به عنوان بهترین برآوردگر ارائه داد و به همین دلیل در آمار به یافتن برآوردگرهای بهینه اکتفا می شود. برای این منظور در آمار دو روش متداول است:

۱- قرار دادن محدودیت های مناسب بر روی برآوردگرها

۲- مرتب کردن برآوردگرها تحت یک قانون معین

در این فصل به اختصار روش دوم را مورد بررسی قرار می دهیم. بدین معنی که با اعمال یک قانون خاص، ابتدا برآوردگرها را مرتب کرده سپس برآوردگر مطلوب را انتخاب می کنیم. اصول مینیماکس و بیزدو اصل معمول از این روش اند.

۳.۱ اصول مینیماکس و بیز

در برآورد پارامتر $\gamma(\theta)$ براساس برآوردگر $\delta(X)$ ، دقت برآوردگر با تابع مخاطره:

$$R(\theta, \delta) = E_{\theta}\{L(\theta, \delta(X))\}$$

اندازه گرفته می شود. در مسایل برآوردیابی علاقمند به دستیابی به برآوردگری مانند δ هستیم که برای هر $\theta \in \Theta$ ، $R(\theta, \delta)$ را مینیمم کند. می دانیم که در حالت کلی، مسئله دارای جواب نیست. نتیجه منطقی این موضوع، پیشنهاد ایجاد محدودیت در کلاس برآوردگرها برای دستیابی به برآوردگرهای بهینه است. گفتیم که یک روش برای پیدا کردن برآوردگرهای بهینه، برقراری یک نوع رابطه ترتیبی بین برآوردگرهاست. دوروش اساسی در برقراری رابطه ای ترتیبی بین برآوردگرها، اصول مینیماکس و بیز است. یک رابطه ترتیبی بین برآوردگرها می تواند براساس رخداد بدترین حالت ممکن برای آماردان

در نظر گرفته شود. به عبارت دیگر برآوردگر δ_1 به برآوردگر δ_2 ترجیح داده می شود اگر

$$\sup_{\theta} R(\theta, \delta_1) \leq \sup_{\theta} R(\theta, \delta_2)$$

بر پایه این موضوع، اصل مینیماکس، برآوردگری از δ را بررسی می کند که مقدار $R(\theta, \delta)$ را مینیمم کند. چنین انتخابی از δ منجر به انتخاب برآوردگر مینیماکس δ_m می شود، یعنی برای هر $\delta \in D$ که فضای کل برآوردها است داشته باشیم:

$$\sup_{\theta} R(\theta, \delta_m) \leq \sup_{\theta} R(\theta, \delta)$$

به عبارت دیگر، براساس این اصل، بدترین حالات ممکن از برآوردها به ازای $\theta \in \Theta$ را در نظر گرفته و بین بدترین حالات ممکن، آن δ یی اختیار می شود که کمترین مقدار مخاطره را اختیار کند. واضح است که در این حالت لزومی ندارد برآوردگر مینیماکس یکتا باشد.

اصل بیز، برآوردگری از θ را بررسی می کند که برای یک تابع وزنی نظیر $G(\theta)$ ، مقدار $\int_{\Theta} R(\theta, \delta) dG(\theta)$ حداقل شود. یعنی براساس تابع وزنی G و تابع زیان L ، کلیه مقادیر $\int_{\Theta} R(\theta, \delta) dG(\theta)$ را برای هر $\delta \in D$ ، از کوچک به بزرگ مرتب و برآوردگر δ —ای را انتخاب می کند که مقدار این انتگرال برای آن از همه کمتر باشد که به انتخاب برآوردگر بیز $\delta_B(X)$ می انجامد، یعنی برای هر برآوردگر $\delta \in D$ ،

$$\int_{\Theta} R(\theta, \delta_B) dG(\theta) \leq \int_{\Theta} R(\theta, \delta) dG(\theta).$$

۴.۱ توزیع پیشین و توزیع پسین

توزیع پارامتر قبل از مشاهده هر داده ای، توزیع پیشین آن پارامتر نامیده می شود. توزیع شرطی پارامتر به شرط مشاهده داده ها توزیع پسین نامیده می شود. اگر مقادیر مشاهده

شده داده‌ها را در تابع احتمال شرطی یا تابع چگالی احتمال شرطی داده‌ها به شرط پارامتر، جایگذاری کنیم نتیجه این کار فقط تابعی از پارامتر است که آن را تابع درست‌نمایی می‌نامیم. با توجه به ماهیت اصل بیز، در مسائل استنباط آماری به روش بیز براساس مشاهداتی که از خانواده توزیع‌ها اختیار می‌شود پارامتر θ دارای یک مقدار نامعلوم است. به عبارت دیگر، در حقیقت θ به عنوان یک مقدار متغیر تصادفی W در نظر گرفته می‌شود که مقادیر ممکن آن فضای پارامتر Θ است و دارای توزیع $G(\theta)$ یا تابع احتمال (تابع چگالی احتمال) $g(\theta)$ است و از آن با عنوان توزیع پیشین یا تابع احتمال پیشین (تابع چگالی احتمال پیشین) یاد می‌کنیم. در واقع، توزیع پیشین تبلور کاربر آمار از خلاصه اطلاعات و دانسته‌های او در این باره است که احتمال قرار داشتن θ در چه بخش‌هایی از Θ بیش از همه است. به بیان دیگر قبل از مشاهده و جمع‌آوری هر گونه داده‌ای، اطلاعات و دانسته‌های قبلی کاربر آمار، وی را متقاعد می‌کند بر این باور باشد که شانس قرار گرفتن مقادیر θ در فضای Θ چگونه است و چنین باورهایی را می‌توان در قالب یک تابع توزیع بیان کرد.

تابع مخاطره بیزی در انتخاب δ به عنوان یک برآوردگر نسبت به توزیع پیشین G و تابع زیان L ، که آن را با نماد $r(G, \delta)$ نمایش می‌دهیم، برابر است با:

$$\begin{aligned} r(G, \delta) &= E\{E[L(W, \delta(X))]\} \\ &= E[R(W, \delta)] \\ &= \int_{\Theta} R(\theta, \delta) dG(\theta) \end{aligned} \quad (1.1)$$

و به دنبال برآوردگر δ_B هستیم که مقدار (1.1) را مینیمم کند. در این حالت از برآوردگر $\delta_B(X)$ با عنوان برآوردگر بیز $\gamma(\theta)$ نسبت به توزیع پیشین G تحت تابع زیان L یاد می‌کنیم. فرض می‌کنیم که مدل احتمالی آزمایش بستگی به پارامتر θ دارد که از خانواده توزیع‌های $\{F_\theta : \theta \in \Theta\}$ با خانواده چگالی‌های $\{f_\theta : \theta \in \Theta\}$ است، همچنین فرض کنیم که θ خود مقدار مشاهده شده یک متغیر تصادفی، مثلاً W با توزیع معلوم $G(\theta)$ یا چگالی معلوم $g(\theta)$ است، که از آن با نام توزیع پیشین W یاد می‌کنیم در این حالت، $f_\theta(\cdot)$ را می‌توان به عنوان

چگالی شرطی متغیر X به شرط $W = \theta$ تلقی کرده و چگالی توأم X و W را از رابطه زیر به دست آورد.

$$f_{X,W}(x, \theta) = f_{X|W=\theta}(x)g(\theta) \quad (2.1)$$

توزیع شرطی W به شرط $X = x$ قابل محاسبه است، که از آن با نام توزیع پسین W یاد می‌کنیم. توزیع پسین W دارای چگالی پسین است که با استفاده از رابطه (2.1) و قضیه بیز به صورت زیر قابل محاسبه است.

$$g_{W|X=x}(\theta) = \frac{f_{X,W}(x, \theta)}{f_X(x)}$$

که در آن تابع چگالی احتمال توأم X و W و $f_X(x)$ تابع چگالی احتمال کناری X است، به طوری که

$$f_X(x) = \begin{cases} \sum_{\Theta} f_{X,W}(x, \theta) & \text{گسسته باشد } W \\ \int_{\Theta} f_{X,\theta}(x, \theta) d\theta & \text{پیوسته باشد } W \end{cases}$$

بنابراین تابع چگالی پسین را می‌توان به صورت زیر بازنویسی کرد:

$$\begin{aligned} g_{W|X=x}(\theta) &\equiv g(\theta|x) \\ &= \frac{f_{X,W}(x, \theta)}{\int_{\Theta} f_{X,W}(x, \theta) dG(\theta)}. \end{aligned}$$

۵.۱ برآوردهای بیز

فرض کنید $\mathbf{X} = (X_1, \dots, X_n)$ نمایانگر یک نمونه تصادفی n تایی از خانواده چگالی‌های $\{f_{\theta} : \theta \in \Theta\}$ باشد. اگر توزیع پیشین انتخابی $G(\theta)$ با چگالی $g(\theta)$ و تابع زیان

$L(\theta, \delta)$ باشد، آن گاه تابع مخاطره برآوردگر δ در برآورد پارامتر $\gamma(\theta)$ به صورت زیر خواهد بود:

$$\begin{aligned} R(\theta, \delta) &= E[L(\theta, \delta(\mathbf{X}))] \\ &= \int_{\mathcal{X}} L(\theta, \delta(\mathbf{x})) f_{\theta}(\mathbf{x}) d\mathbf{x} \end{aligned}$$

و تابع مخاطره بیزی برآوردگر δ در برآورد $\gamma(\theta)$ نسبت به توزیع پیشین G و تابع زیان L ، به صورت زیر است (توجه کنید که θ نقش یک متغیر را دارد).

$$\begin{aligned} r(G, \delta) &= E[R(\theta, \delta)] \\ &= \int_{\Theta} R(\theta, \delta) dG(\theta) \\ &= \int_{\Theta} R(\theta, \delta) g(\theta) d\theta \end{aligned}$$

بنابراین به سادگی معلوم می شود که:

$$\begin{aligned} r(G, \delta) &= \int_{\Theta} R(\theta, \delta) g(\theta) d\theta \\ &= \int_{\Theta} \left\{ \int_{\mathcal{X}} L(\theta, \delta(\mathbf{x})) f_{\theta}(\mathbf{x}) d\mathbf{x} \right\} g(\theta) d\theta \\ &= \int_{\Theta} \int_{\mathcal{X}} L(\theta, \delta(\mathbf{x})) f_{\theta}(\mathbf{x}) g(\theta) d\mathbf{x} d\theta \\ &= \int_{\Theta} \int_{\mathcal{X}} L(\theta, \delta(\mathbf{x})) f(\mathbf{x}, \theta) d\mathbf{x} d\theta \\ &= \int_{\Theta} \int_{\mathcal{X}} L(\theta, \delta(\mathbf{x})) g(\theta|\mathbf{x}) f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} d\theta \\ &= \int_{\Theta} \left\{ \int_{\mathcal{X}} L(\theta, \delta(\mathbf{x})) g(\theta|\mathbf{x}) d\mathbf{x} \right\} f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} \end{aligned}$$

انتگرال داخلی در رابطه اخیر، یعنی

$$\int_{\mathcal{X}} L(\theta, \delta(\mathbf{x})) g(\theta|\mathbf{x}) d\mathbf{x}$$

مخاطره پسین نامیده می شود.

برآوردگر بیز پارامتر $\gamma(\theta)$ نسبت به چگالی پیشین g و تابع زیان L ، که آن را برای سهولت با δ_g نمایش می دهیم، آن برآوردگر تعریف می شود که دارای مینیمم تابع مخاطره بیزی باشد،

یعنی

$$r(g, \delta_g) = \min r(g, \delta)$$

توجه داشته باشید که تعریف فوق روش پیدا کردن برآوردگر بیز $\gamma(\theta)$ را ارائه نمی‌کند. برای پیدا کردن برآوردگر بیز پارامتر $\gamma(\theta)$ ، می‌توانیم به صورت زیر عمل کنیم. با توجه به این که به دنبال برآوردگری نظیر $\delta_g(X)$ هستیم که $r(g, \delta)$ را به عنوان تابعی از برآوردگرهای ممکن، یعنی به ازای هر $\delta \in D$ ، مینیمم کند و از طرفی نشان دادیم که:

$$\begin{aligned} r(g, \delta) &= \int_{\Theta} \left\{ \int_{\mathcal{X}} L(\theta, \delta(\mathbf{x})) f_{\theta}(\mathbf{x}) d\mathbf{x} \right\} g(\theta) d\theta \\ &= \int_{\Theta} \left\{ \int_{\mathcal{X}} L(\theta, \delta(\mathbf{x})) g(\theta|\mathbf{x}) d\theta \right\} f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x} \end{aligned}$$

و چون انتگرال دوگانه اخیر غیرمنفی است، انتگرال دوگانه را می‌توان مینیمم کرد، اگر بتوانیم عبارت داخل آکولاد، یعنی مخاطره پسین را برای هر \mathbf{x} داده شده‌ای مینیمم کنیم، که آن را با $\delta_g(\mathbf{x})$ نمایش می‌دهیم. در صورتی که برای هر \mathbf{x} ، $\delta_g(\mathbf{x})$ منحصر به فرد باشد، آنگاه برآوردگر بیز منحصر به فرد خواهد بود. توجه داشته باشید که با تغییر توزیع پیشین و تابع زیان، برآوردگر بیز پارامتر مورد نظر تغییر می‌کند.

۶.۱ برآوردگرهای مینیماکس

همان طور که گفته شد، برآوردگر مینیماکس δ_m عبارت از برآوردگری است که سوپریمم مقدار تابع مخاطره آن کمتر یا مساوی سوپریمم مقدار تابع مخاطره هر برآوردگر دیگر باشد. یعنی δ_m یک برآوردگر مینیماکس است، اگر برای هر $\delta \in D$ ،

$$\sup_{\theta} R(\theta, \delta_m) \leq \sup_{\theta} R(\theta, \delta)$$

استفاده از تعریف درستیابی به برآوردهای مینیماکس چندان ساده به نظر نمی‌رسد. با استفاده از قضیه زیر، به عنوان ابزاری درستیابی به برآوردهای مینیماکس می‌توان استفاده کرد.

قضیه ۱-۲: اگر $\delta_g(x)$ برآوردهای بی‌پارامتر $\gamma(\theta)$ نسبت به چگالی پیشین $g(\theta)$ با تابع مخاطره ثابت باشد، آنگاه برآوردهای δ_g یک برآوردهای مینیماکس است. (اثبات در پیوست B)

فصل ۲

رگرسیون

۱.۲ مقدمه

تحلیل رگرسیونی یک روش آماری برای بررسی و به مدل درآوردن ارتباط بین متغیرهاست. کاربردهای رگرسیون متعدد بوده و معمولاً در زمینه‌های مختلف از جمله مهندسی، فیزیک، اقتصاد، مدیریت، علوم زیستی، پزشکی و علوم اجتماعی به کار می‌رود. در حقیقت تحلیل رگرسیونی شاید از جمله روش‌های آماری با بیشترین و وسیعترین کاربرد باشد.

۲.۲ رگرسیون و ساختن مدل

اولین بار سرفرانسیس گالتون^۱ در سال ۱۸۸۵ اصطلاح رگرسیون را در مقاله خود به کار برد و بعد از آن این موضوع در ابعاد وسیع‌تری بیان شد. تحلیل رگرسیونی در حقیقت یک ابزار آماری برای تشخیص رابطه بین یک یا چند متغیر مستقل با یک متغیر وابسته است.

^۱ Francis Galton

در رگرسیون درباره ارتباط بین متغیرها بحث می‌شود. متغیرها را می‌توان به دو دسته متغیرهای پیشگو^۲ یا مستقل و متغیرهای پاسخ^۳ یا وابسته تقسیم کرد. منظور از متغیرهای پیشگو یا مستقل متغیرهایی هستند که می‌توان آنها را برابر مقادیری گرفت که گرچه ممکن است تحت کنترل نباشند، اما می‌توانند مشاهده شوند و این متغیرها با توجه به نامشان به طور مستقل تغییر می‌کنند. در نتیجه تغییر عمدی این متغیرها اثرهایی در متغیرهای دیگری یعنی متغیرهای پاسخ پدید می‌آورند. وجه تمایز بین متغیرهای پیشگو و پاسخ همیشه روشن نیست و گاهی وابسته به هدف تحلیل‌گر است اما معمولاً در عمل می‌توان نقش متغیرها را تشخیص داد.

در رگرسیون معمولاً فرض بر این است که متغیرهای پیشگو مقید به تغییرات تصادفی نیستند ولی متغیر پاسخ تصادفی است. از نقطه نظر کاربردی این مطلب اغلب صحیح نیست، اما به دلیل این که اگر متغیرهای پیشگو تصادفی باشند، شیوه برآزش بسیار پیچیده می‌شود، برای اجتناب از این پیچیدگی معمولاً فرض می‌شود که تغییرات تصادفی هر یک از متغیرهای پیشگو در مقایسه با دامنه متغیر پیشگوی مشاهده شده آن قدر کوچک است که می‌توان از آن صرف‌نظر کرد.

باید توجه داشت که در رگرسیون رابطه بین متغیرهای پیشگو و پاسخ یک رابطه تابعی ریاضی نیست، بلکه رابطه‌ای آماری است. به این معنی که به ازای هر x مقدار y کاملاً مشخص و دقیق نیست و تعیین مقدار آن با مقداری خطا همراه است. نوع رابطه‌ای که ممکن است بین متغیرهای پیشگو و وابسته برقرار باشد، متنوع است اما معمولاً به دلیل سادگی و نیز به دلیل وجود برخی مبانی نظری در مورد تقریب روابط غیرخطی با روابط خطی، به یافتن رابطه خطی علاقه‌مندی بیشتری نشان داده می‌شود.

در حالت رگرسیون خطی مدل آماری به صورت $Y = X\beta + \varepsilon$ در نظر گرفته می‌شود

Predictive^۲
Response^۳

که در آن X بردار متغیرهای پیشگو، β پارامترها و ε میزان خطا و Y متغیر پاسخ را نشان می‌دهد. پس از تشخیص و پذیرش نوع رابطه اولین مرحله برآورد پارامترهای نامعلوم مدل است. با برآورد کردن پارامترهای نامعلوم، مدل کاملاً مشخص شده و پس از انجام نیکویی برازش می‌توان از آن برای پیش‌بینی استفاده کرد.

۳.۲ کاربردهای مدل‌های رگرسیون

مدل‌های رگرسیونی برای مقاصدی چند مشتمل بر موارد زیر مورد استفاده قرار می‌گیرند.

۱- توصیف داده‌ها

۲- برآورد پارامترها

۳- پیشگویی و برآورد

۴- کنترل

۴.۲ رگرسیون خطی ساده

مدل رگرسیون خطی ساده مدلی است با یک متغیر رگرسیونی مستقل x که با یک متغیر پاسخ y ارتباطی به صورت خط مستقیم دارد. این مدل رگرسیونی خطی ساده عبارت است از:

$$y = \beta_0 + \beta_1 x + \varepsilon \quad (1.2)$$

که β_0 عرض از مبدا و β_1 شیب خط ثابت‌های نامعلوم اند و ε مولفه خطای تصادفی است. فرض می‌شود خطاها دارای میانگین صفر و واریانس σ^2 نامعلوم هستند و بعلاوه مقادیر این خطا ناهمبسته‌اند یعنی این که مقدار یک خطا بستگی به مقدار هر خطای دیگر ندارد.

نکته قابل توجه این که مقادیر متغیر رگرسیونی x توسط محقق و تحلیل‌گر اختیار و با خطای قابل صرف‌نظر کردن اندازه‌گیری می‌شود. در حالی که به ازاء هر مقدار ممکن x ، پاسخ y متغیری تصادفی است بدین معنی که به ازاء هر مقدار ممکن x متغیر تصادفی y دارای یک توزیع احتمال است. میانگین این توزیع عبارتست از:

$$E(y|x) = \beta_0 + \beta_1 x$$

و واریانس آن عبارتست از:

$$V(y|x) = V(\beta_0 + \beta_1 x + \varepsilon) = \sigma^2$$

بنابراین میانگین y تابعی خطی از x است در حالی که واریانس آن به مقدار x بستگی ندارد و بعلاوه چون خطاهای ε ناهمبسته‌اند پاسخ‌ها نیز ناهمبسته خواهند بود. پارامترهای β_0 و β_1 ضرایب رگرسیون نامیده می‌شوند. β_1 (شیب) تغییر در میانگین توزیع y به ازاء یک واحد تغییر در x است. اگر دامنه x شامل $x = 0$ باشد، در این صورت β_0 (عرض از مبدا) میانگین پاسخ y برای $x = 0$ است. اگر دامنه x شامل صفر نباشد، در این صورت β_0 تعبیر عملی (واقعی) ندارد.

۵.۲ رگرسیون چندگانه خطی

مدل رگرسیونی که مشتمل بر بیش از یک متغیر رگرسیونی باشد مدل رگرسیون چندگانه نامیده می‌شود.

فرض می‌کنیم عمر مؤثر یک وسیله برش با سرعت و عمق بریدگی ارتباط دارد. یک مدل رگرسیون چندگانه که می‌تواند این ارتباط را توصیف کند عبارتست از:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon \quad (2.2)$$

که y نشان دهنده عمر مؤثر وسیله، x_1 نشان دهنده سرعت برش و x_2 عمق برش می باشد. این یک مدل رگرسیون خطی چندگانه با دو متغیر رگرسیونی x_1 و x_2 است. کاربرد اصطلاح خطی به لحاظ این است که رابطه (۲.۲) تابعی خطی از پارامترهای نامعلوم β_0 ، β_1 و β_2 می باشد. این مدل یک صفحه را در فضای متغیرهای رگرسیونی x_1 ، x_2 مشخص می کند که پارامتر β_0 ارتفاع عرض از مبدا صفحه رگرسیونی است. اگر دامنه تعریف داده ها شامل $x_1 = x_2 = 0$ باشد در این صورت β_0 میانگین y به شرط $x_1 = x_2 = 0$ خواهد بود. پارامتر β_1 حاکی از تغییر مورد انتظار برای y به ازاء یک واحد تغییر x_1 است، وقتی که x_2 ثابت است. همچنین β_2 برابر است با تغییر مورد انتظار برای y به ازاء یک واحد تغییر x_2 وقتی که x_1 ثابت نگهداشته شود.

در حالت کلی متغیر پاسخ y ممکن است به p متغیر رگرسیونی بستگی داشته باشد، مدل

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \varepsilon \quad (۳.۲)$$

یک مدل رگرسیون چندگانه خطی با p متغیر رگرسیونی نامیده می شود. پارامترهای β_j ، $j = 0, 1, \dots, p$ ضرایب رگرسیون نامیده می شوند. این مدل یک ابرصفحه در فضای p بعدی از متغیرهای رگرسیونی x_j است. پارامتر β_j نشان دهنده تغییرات مورد انتظار متغیر پاسخ به ازاء یک واحد تغییر در x_j است، وقتی که همه متغیرهای رگرسیونی دیگر x_j ($i \neq j$) ثابت نگهداشته شوند. به همین جهت پارامترهای β_j ، $j = 0, 1, \dots, p$ ضرایب جزئی رگرسیون نامیده می شوند.