

دانشگاه پیام نور

دانشکده علوم پایه

پایان نامه

برای دریافت مدرک کارشناسی ارشد

رشته آمار ریاضی

گروه آمار

عنوان پایان نامه:

مطالعه‌ای بر بر آوردیابی استوار در مدل‌های رگرسیونی

مریم افتخاری

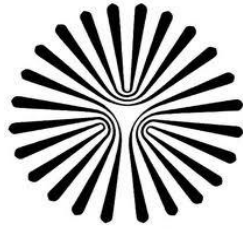
استاد راهنما:

دکتر نرگس عباسی

شهریور ۱۳۹۲



بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه پیام نور

دانشکده علوم پایه

مرکز شیراز

پایان نامه

برای دریافت مدرک کارشناسی ارشد

رشته آمار ریاضی

گروه آمار

عنوان پایان نامه:

**مطالعه‌ای بر برآوردیابی استوار در مدل‌های رگرسیونی**

مریم افتخاری

استاد راهنما:

دکتر نرگس عباسی

شهریور ۱۳۹۲

تاریخ : ۹۲/۰۶/۲۸

شماره : ۱۵/۳۳۲۶

پیوست :



دانشگاه پیام نور شیراز  
باسم تعالی



جمهوری اسلامی ایران  
وزارت علوم، تحقیقات و فناوری  
دانشگاه پیام نور استان فارس

### صور تجلسه دفاع از پایان نامه دوره کارشناسی ارشد

جلسه دفاع از پایان نامه دوره کارشناسی ارشد خانم مریم افتخاری دانشجوی رشته آمار ریاضی به شماره دانشجویی ۹۰۰۰۱۴۲۲۸ با عنوان:

" مطالعه‌ای بر بر آوردیابی استوار در مدل های رگرسیونی "

با حضور هیأت داوران در روز پنجشنبه مورخ ۱۳۹۲/۰۶/۲۸ ساعت ۷:۳۰ صبح در محل ساختمان غدیر دانشگاه پیام نور شیراز برگزار شد و هیأت داوران پس از بررسی، پایان نامه مذکور را شایسته نمره به عدد ۱۹ به حروف فوززده با درجه خوبی تشخیص داد.

ردیف	نام و نام خانوادگی	هیات داوران	رتبه دانشگاهی	دانشگاه	امضاء
۱	دکتر نرگس عباسی	راهنما	دانشیار	پیام نور شیراز	
۲	عبدالرضا بازرگان لاری	داور	استادیار	شیراز	
۳	امیر اکبری	نماینده تحصیلات تکمیلی	مریی	پیام نور شیراز	

رئیس اداره تحصیلات تکمیلی

شیراز- شهرک گلستان بلوار دهخدا  
قبل از نمایشگاه بین المللی  
تلفن : ۰۷۱۱ - ۶۲۲۲۲۵۵  
دورنگار : ۰۷۱۱ - ۶۲۲۲۲۲۹  
صندوق پستی : ۱۳۶۸ - ۷۱۹۵۵  
www.spnu.ac.ir  
Email : admin@spnu.ac.ir

### گواهی اصالت نشر و حقوق مادی و معنوی اثر

اینجانب مریم افتخاری دانشجوی ورودی سال ۱۳۹۰ مقطع کارشناسی ارشد رشته آمار ریاضی گواهی می‌نمایم چنانچه در پایان نامه خود از فکر، ایده و نوشته دیگری بهره گرفته‌ام یا نقل قول مستقیم یا غیر مستقیم منبع و ماخذ آن را نیز در جای مناسب ذکر کرده‌ام. بدیهی است مسئولیت تمامی مطالبی که نقل قول دیگران نباشد بر عهده خویش می‌دانم و جوابگوی آن خواهم بود.

دانشجو تائید می‌نماید که مطالب مندرج در این پایان نامه نتیجه تحقیقات خودش می‌باشد و در صورت استفاده از نتایج دیگران مرجع آن را ذکر نموده است.

نام و نام خانوادگی دانشجو مریم افتخاری

تاریخ و امضاء ۹۲/۶/۲۸

اینجانب مریم افتخاری دانشجوی ورودی سال ۱۳۹۰ مقطع کارشناسی ارشد رشته آمار ریاضی گواهی می‌نمایم چنانچه بر اساس مطالب پایان نامه خود اقدام به انتشار مقاله، کتاب، و ... نمایم ضمن مطلع نمودن استاد راهنما، با نظر ایشان نسبت به نشر مقاله، کتاب، و ... و به صورت مشترک و با ذکر نام استاد راهنما مبادرت نمایم.

نام و نام خانوادگی دانشجو مریم افتخاری

تاریخ و امضاء ۹۲/۶/۲۸

کلیه حقوق مادی مترتب از نتایج مطالعات، آزمایشات و نوآوری ناشی از تحقیق موضوع این پایان نامه مطلع به دانشگاه پیام نور می‌باشد.

شهریور ۹۲

تقدیم به:

پدر، مادر گرامی ام

و به تمام آزاد مردانی که نیک می‌اندیشند و عقل و منطق را پیشه خود نموده و جز رضای الهی و پیشرفت و سعادت جامعه، هدفی ندارند. دانشمندان، بزرگان، و جوانمردانی که جان و مال خود را در حفظ و اعتلای این مرز و بوم فدا نموده و می‌نمایند.

با تشکر و قدردانی فراوان

از استاد گرامی ام سرکار خانم دکتر نرگس عباسی به دلیل یاری‌ها و راهنمایی‌های بی‌چشمداشت ایشان که بسیاری از سختی‌ها را برایم آسان‌تر نمودند، چرا که بدون راهنمایی‌های ایشان تأمین این پایان‌نامه بسیار مشکل می‌نمود.



## چکیده:

رگرسیون استوار یکی از مهمترین روش‌های آنالیز داده‌هایی می‌باشد که آلوده هستند و نقاط پرت دارند. این روش می‌تواند جهت تعیین نقاط پرت استفاده شود و نتایج پایداری را در حضور داده‌های پرت فراهم آورد. در این پایان‌نامه روش رگرسیون استوار معرفی و در نرم افزار SAS/STAT- Version 9 اجرا شده است. این متن روش‌های مهم رگرسیون استوار را معرفی می‌کند. این روش‌ها شامل برآورد  $M$ ، برآورد LTS، برآورد  $S$  و برآورد  $MM$  می‌باشد. روش استوار در آزمون مدل‌های جزئی خطی تعمیم‌یافته بررسی می‌شود به این صورت که یک خانواده از آماره‌های استوار معرفی می‌شود که امکان انتخاب بین یک مدل پارامتری و نیمه‌پارامتری را فراهم می‌آورد. به طور دقیق‌تر تحت یک مدل جزئی خطی تعمیم‌یافته یعنی زمانی که  $(y_i, x_i, t_i)$  مشاهدات مستقل هستند بطوری که  $y_i | (x_i, t_i) \sim F(0, \mu_i)$  با  $Var(y_i | (x_i, t_i)) = V(\mu_i)$  و  $\mu_i = H(\eta(t_i) + x_i^T \beta)$  که در آن  $H = g^{-1}$  یک تابع پیوندی شناخته شده و  $\beta \in \mathbb{R}^p$  یک پارامتر ناشناخته و  $\eta(t_i)$  یک تابع هموار ناشناخته می‌باشد. در این روش برآوردهای حاصل از  $\beta$  دارای توزیع نرمال مجانبی و سازگاری ریشه  $n$  خواهد بود.

## واژگان کلیدی:

مدل‌های جزئی خطی تعمیم‌یافته، آزمون‌سازی استوار، برآوردهای استوار، نقاط نفوذی، داده پرت

## فهرست

۱	مقدمه
۴	فصل اول
۴	مقدمه‌ای بر مدل‌های آماری
۵	۱-۱ مقدمه
۵	۲-۱ رگرسیون خطی ساده
۶	۱-۲-۱ روش کمترین مربعات برآورد پارامترهای مدل
۷	۳-۱ الگوی رگرسیون خطی چندگانه
۸	۱-۳-۱ برآوردگر کمترین مربعات برای بردار پارامترها
۱۰	۴-۱ رگرسیون لجستیک
۱۱	۵-۱ مدل‌های خطی تعمیم یافته
۱۳	۶-۱ مدل‌های جزئی خطی تعمیم یافته
۱۴	۷-۱ برآورد شبه‌ماکسیمم درست‌نمایی
۱۶	فصل دوم
۱۶	رگرسیون استوار و تشخیص نقاط پرت
۱۷	۱-۲ مقدمه
۱۸	۲-۲ برآوردگرهای استوار
۱۹	۱-۲-۲ برآورد $m$ نوع هابر
۲۱	۲-۲-۲ برآورد $lts$
۲۲	۳-۲-۲ برآورد $LMS$
۲۳	۴-۲-۲ برآوردگر $S$
۲۴	۵-۲-۲ برآوردگر $MM$

۲۷	۳-۲ تشخیص مقاوم و کشف داده‌های پرت
۲۷	۱-۳-۲ فاصله استوار
۲۷	۲-۳-۲ داده‌های نافذ
۲۷	۳-۳-۲ داده‌های پرت
۲۸	۴-۲ استنباط استوار
۲۹	۱-۴-۲ آزمون‌های خطی
۳۰	۵-۲ مطالعات شبیه‌سازی
۳۰	۱-۵-۲ آنالیز واریانس استوار
۳۲	۲-۵-۲ مطالعه‌ی رشد
۳۹	۳-۵-۲ نمودارهای گرافیکی
۴۳	فصل سوم
۴۳	آزمون‌سازی در مدل‌های جزئی خطی تعمیم‌یافته با یک رویکرد استوار
۴۴	۱-۳ مقدمه
۴۸	۲-۳ روش برآورد
۵۰	۳-۳ آزمون آماری
۵۳	۴-۳ محاسبات پیشرفته برای برآوردیابی
۵۶	۵-۳ نتیجه‌گیری
۵۷	فصل چهارم
۵۷	مطالعات شبیه‌سازی
۵۸	۱-۴ مقدمه
۵۸	۲-۴ روش کار در مونت کارلو
۶۰	۳-۴ نتایجی از روی مونت کارلو

۶۵	..... ۴-۴ نتیجه گیری
۶۶	..... پیوست
۶۷	..... پیوست الف- خطای اندازه گیری
۷۷	..... پیوست ب- اثبات قضیه (۱-۳)
۸۴	..... منابع

در این پایان‌نامه روش‌های آماری استوار را معرفی می‌شود، روش‌های استوار به روش‌هایی اطلاق می‌شود که حساسیت کمی نسبت به داده‌های پرت داشته و به سادگی تحت تاثیر داده‌های کوچک و بزرگ قرار نگیرند. اولین بار باکس<sup>۱</sup> (۱۹۵۳) از واژه‌ی استوار استفاده کرد. اما چندان مورد توجه قرار نگرفت. آماردانان تا مدت‌ها از بررسی چنین موضوعی گریزان بودند. هویر<sup>۲</sup> (۱۹۶۴) برای اولین بار برآوردگرهای استوار مکان را مورد بررسی قرار داد.

آمار استوار به مفهوم روش‌هایی برای حل مسائل ناشی از اثر داده‌های پرت است. داده‌های پرت اثرات بدی بر روی برآورد پارامتر و همچنین فرضیات مدل دارند. به طور کلی آمار استوار وقتی بکار می‌رود که فرضیات ریاضی داده‌های برقرار نباشد، ولی عموماً داده‌های پرت باعث از بین بردن فرضیات مدل می‌شود که در این حالت از آمار استوار استفاده می‌شود. تقریباً همیشه وقتی صحبت از آمار استوار به میان می‌آید داده‌های پرت به فکر ما خطور می‌کند.

به طور مثال در اکثر برنامه‌های کاربردی، پارامترهای مدل آمیخته بوسیله برآوردکننده‌ی درست‌نمایی ماکسیمم<sup>۳</sup> از طریق الگوریتم امید ماکسیمم‌سازی<sup>۴</sup> برآورد شده است. در حالی که برآورد ماکسیمم درست‌نمایی در برخورد با داده‌های پرت بسیار حساس می‌باشد. هویر (۱۹۸۱) همپل<sup>۵</sup> و همکاران (۱۹۸۶) روسو<sup>۶</sup> و لروی<sup>۷</sup> (۱۹۸۷). برآزش مدل‌های آمیخته را می‌توان مستقیماً برای داده‌ها توسط برآوردکننده‌های استوار محدود استفاده کردند. دلیل آن این است که این برآوردکننده‌های استوار به تناسب مدل پارامتری برای اکثریت داده‌ها طراحی شده‌اند. در حالی که داده‌های باقی مانده که از مدل تبعیت نمی‌کنند به عنوان داده‌های پرت در نظر گرفته می‌شوند. با این حال، در عمل، داده‌ها می‌توانند کاملاً بدون داشتن بخش همگن متشکل از حداقل ۵۰٪ از داده‌ها ناهمگن باشند.

در بسیاری موقعیت‌ها مدل خطی برای بیان رابطه‌ی بین متغیر پاسخ و متغیر کمکی متناظر نارسا است.

---

<sup>۱</sup> Box

<sup>۲</sup> Huber

<sup>۳</sup> MLE: Maximum Likelihood Estimator

<sup>۴</sup> EM: Expectation-Maximum

<sup>۵</sup> Hampel

<sup>۶</sup> Rousseeuw

<sup>۷</sup> Leroy

آزمون فرض برای مدل‌های جزئی خطی تعمیم‌یافته، عمدتاً مبنی بر مقایسه‌ی برآوردگرهای پارامتری هموار شده، متمرکز است. برای نمونه هاردل<sup>۱</sup> و همکارانش (۱۹۹۸) یک آزمون آماری برای انتخاب بین یک مدل خطی و یک مدل نیمه‌پارامتری را مد نظر قرار دادند. طرح پیشنهادی آن‌ها مبنی بر شیوه‌ی تعدیل‌یافته برآورد مطرح شده توسط سورینی<sup>۲</sup> و استنیس‌ولس<sup>۳</sup> (۱۹۹۴) جهت مقابله با درست‌نمایی هموارنشده و هموارشده می‌باشد. یک مطالعه‌ی تطبیقی از روش‌های مختلف توسط مولر<sup>۴</sup> (۲۰۰۱) ارائه گردید در حالی‌که یک رویکرد متفاوت توسط رودریگز کامپاس<sup>۵</sup> و همکاران (۱۹۹۸) بررسی شده بود.

معروف است چنین برآوردهایی در مواجهه با نقاط پرت ناکارآمد هستند و در پی آن آزمون‌های آماری که از روی آن‌ها ساخته می‌شود منجر به نتیجه‌گیری نادرست می‌گردد. روش‌های استوار برای مدل‌های خطی تعمیم‌یافته توسط استفانسکی<sup>۶</sup> و همکاران (۱۹۸۶)، کانش<sup>۷</sup> و همکاران (۱۹۸۹)، بیانکو<sup>۸</sup> و یوهای<sup>۹</sup> (۱۹۹۵)، کانتونی<sup>۱۰</sup> و رانچتی<sup>۱۱</sup> (۲۰۰۱)، کراکس<sup>۱۲</sup> و هاسبروک<sup>۱۳</sup> (۲۰۰۲) و بیانکو و همکاران (۲۰۰۵) مورد توجه قرار گرفته است.

علاوه بر این آزمون‌های استوار برای فرض مقابل، تحت مدل رگرسیون جزئی خطی توسط بیانکو و همکاران (۲۰۰۶) مورد مطالعه قرار گرفت.

در این پایان‌نامه روش‌های آماری استوار برای انواع مدل‌های رگرسیونی معرفی شده و تمرکز بیشتر بر روی مدل‌های جزئی خطی تعمیم‌یافته است. جنبه جدید بودن این پایان‌نامه بررسی روش‌های آماری استوار در مدل‌های جزئی خطی تعمیم‌یافته می‌باشد.

در فصل اول به تعاریف و مفاهیم مربوط به مدل‌های رگرسیون، مدل‌های جزئی خطی تعمیم‌یافته و مدل‌های خطی تعمیم‌یافته پرداخته شده است. در فصل دوم رگرسیون استوار و تشخیص نقاط پرت، در

---

<sup>1</sup>Hardli

<sup>2</sup>Severini

<sup>3</sup>Staniswalis

<sup>4</sup>Müller

<sup>5</sup>Rodríguez Campos

<sup>6</sup>Stefanski

<sup>7</sup>Künsch

<sup>8</sup>Bianco

<sup>9</sup>Yohai

<sup>10</sup>Cantoni

<sup>11</sup>Ronchetti

<sup>12</sup>Croux

<sup>13</sup>Haesbroeck

فصل سوم آزمون‌سازی در مدل‌های جزئی خطی تعمیم‌یافته با یک رویکرد استوار و در فصل چهارم شبیه‌سازی و مقایسه روش استوار با روش کلاسیک و نتیجه‌ی مطالعات مونت‌کارلو<sup>۱</sup> ارائه شده است. و بالاخره در قسمت ضمیمه خطای اندازه‌گیری و اثبات‌ها مورد بررسی قرار گرفته است.

---

<sup>۱</sup>Monte Carlo

## **فصل اول**

### **مقدمه‌ای بر مدل‌های آماری**



## ۱-۱ مقدمه

در آمار رگرسیون یعنی یک نوع رابطه یا تابع ریاضی که بین متغیر وابسته از یک سو و متغیرهای مستقل از سوی دیگر برقرار می‌باشد. تحلیل رگرسیونی یک ابزار آماری مفید برای دیدن رابطه‌ی بین متغیرها به کار می‌رود.

تحلیل رگرسیون فن و تکنیکی آماری برای بررسی و به مدل در آوردن ارتباط بین متغیرهاست. کاربردهای رگرسیون متعدد است و تقریباً در هر زمینه‌ای از جمله مهندسی، فیزیک، اقتصاد، مدیریت، علوم زیستی و بیولوژی و علوم اجتماعی صورت می‌پذیرد. در حقیقت تحلیل رگرسیونی ممکن است فن و تکنیک آماری با بیشترین و وسیعترین کاربرد بین تکنیک‌های آماری باشد.

در این فصل مدل‌های رگرسیونی، مدل‌های جزئی خطی تعمیم‌یافته و مدل‌های خطی تعمیم‌یافته مورد بررسی قرار گرفته است. برای مدل‌های رگرسیون روش کمترین مربعات و برای مدل‌های خطی تعمیم‌یافته روش برآورد درست‌نمایی ماکسیمم آورده شده است که هر دوی این روش‌های برآورد جزء روش‌های کلاسیک می‌باشد. همان‌طور که اشاره شد روش‌های کلاسیک در مقابل داده‌های پرت به نتایج اشتباه منجر می‌شوند. در فصل بعدی روش‌های رگرسیون استوار برای حل این مشکل ارائه شده است.

## ۱-۲ رگرسیون خطی ساده

مدل رگرسیون خطی ساده را برای  $n$  مشاهده می‌توان به صورت زیر نوشت:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, 2, 3, \dots, n. \quad (1-1)$$

یک تخصیص ساده نشان می‌دهد که فقط یک  $x$  برای پیش بینی  $y$  وجود دارد، و مدل خطی یعنی مدل (۱-۱) بر حسب  $\beta_0$  و  $\beta_1$  خطی است. که در واقع داریم  $\mathbb{E}(y_i) = \beta_0 + \beta_1 x_i$  که خطی است. برای مثال، یک مدل مانند  $y_i = \beta_0 + \beta_1 x_i^2 + \varepsilon_i$  نسبت به  $\beta_0$  و  $\beta_1$  خطی است، ولی مدل  $y_i = \beta_0 + e^{\beta_1 x_i} + \varepsilon_i$  خطی نیست.

برای تحلیل مدل (۱-۱)، فرض‌های زیر را در نظر می‌گیریم:

$$(۱) \text{ برای تمام } i = 1, 2, \dots, n, \mathbb{E}(\varepsilon_i) = 0, \text{ یا به عبارت دیگر معادل } \mathbb{E}(y_i) = \beta_0 + \beta_1 x_i.$$

(۲) برای تمام  $i = 1, 2, \dots, n$  ،  $Var(\varepsilon_i) = \sigma^2$  یا به عبارت دیگر  $Var(y_i) = \sigma^2$ .

(۳) برای تمام  $i = 1, 2, \dots, n$  ،  $cov(\varepsilon_i, \varepsilon_j) = 0$  یا به عبارت دیگر  $cov(y_i, y_j) = 0$ .

فرض (۱) بیان می‌کند که مدل (۱-۱) صحیح است، و نتیجه می‌دهد که  $y_i$  فقط به  $x_i$  بستگی دارد و تمام تغییرات دیگر  $y_i$  تصادفی است.

## ۱-۲-۱ روش کمترین مربعات برآورد پارامترهای مدل

با استفاده از یک نمونه تصادفی  $n$  مشاهده  $y_1, y_2, \dots, y_n$  به همراه مقادیر ثابت  $x_1, x_2, \dots, x_n$  می‌توانیم پارامترهای  $\beta_0$  و  $\beta_1$  و  $\sigma^2$  را برآورد کنیم. برای به دست آوردن  $\hat{\beta}_0$  و  $\hat{\beta}_1$  از روش کمترین مربعات استفاده می‌کنیم، که به هیچ‌گونه مفروضات توزیعی نیاز ندارد.

در روش کمترین مربعات، برآوردگرهای  $\hat{\beta}_0$  و  $\hat{\beta}_1$  را به قسمی پیدا می‌کنیم که مجموع مربعات  $y_i - \hat{y}_i$ ، برای  $n$  مشاهده  $y_i$  را از مقادیر پیش‌بینی آنها  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$  مینیمم سازد.

$$\hat{\varepsilon}'\hat{\varepsilon} = \sum_{i=1}^n \hat{\varepsilon}_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) \quad (2-1)$$

توجه کنید که مقدار پیش‌بینی شده  $\hat{y}_i$  مقدار  $E(y_i)$  را برآورد می‌کند نه  $y_i$  را، یعنی  $\hat{\beta}_0 + \hat{\beta}_1 x_i$  مقدار  $\beta_0 + \beta_1 x_i$  را برآورد می‌کند نه  $\beta_0 + \beta_1 x_i + \varepsilon_i$  یک نماد بهتر باید  $E(\hat{y}_i)$  باشد، ولی اغلب به کار برده می‌شود.

برای یافتن مقادیر  $\hat{\beta}_0$  و  $\hat{\beta}_1$  که  $\hat{\varepsilon}'\hat{\varepsilon}$  را در (۲-۱) مینیمم سازد، نسبت به  $\hat{\beta}_0$  و  $\hat{\beta}_1$  مشتق می‌گیریم و حاصل را برابر صفر قرار می‌دهیم:

$$\frac{\partial \hat{\varepsilon}'\hat{\varepsilon}}{\partial \hat{\beta}_0} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \quad (3-1)$$

$$\frac{\partial \hat{\varepsilon}'\hat{\varepsilon}}{\partial \hat{\beta}_1} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i = 0 \quad (4-1)$$

جواب (۳-۱) و (۴-۱) به صورت زیر داده می‌شود:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (5-1)$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (6-1)$$

برای تحقیق این که  $\hat{\beta}_0$  و  $\hat{\beta}_1$  در (5.1) و (6.1) مقدار  $\hat{\epsilon}$  را در (2.1) مینیمم می سازد، می توانیم از مشتق دوم استفاده کنیم یا به طور ساده تر ملاحظه کنیم که  $\hat{\epsilon}$  دارای اکسترمم نیست، در نتیجه مشتق اول مینیمم آن را نشان می دهد.

### ۱-۳ الگوی رگرسیون خطی چندگانه

پاسخ  $y$  اغلب تحت تاثیر بیش از یک متغیر مستقل است. به عنوان مثال، محصول یک بذر ممکن است به مقدار نیتروژن، پتاسیم و فسفات در کود مورد استفاده بستگی داشته باشد. این متغیرها توسط آزمایشگر کنترل می شوند، ولی محصول ممکن است به متغیرهای کنترل نشده ای مانند آن ها که با آب و هوا در ارتباط اند نیز وابسته باشد.

یک الگوی خطی متغیر پاسخ  $y$  را به چندین متغیر مستقل به صورت زیر مربوط می سازد.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \epsilon_0 \quad (7-1)$$

پارامترهای  $\beta_0, \beta_1, \dots, \beta_k$  را ضرایب رگرسیون می نامند.  $\epsilon_0$  تغییرات تصادفی  $y$  را که به وسیله  $x$  ها بیان نمی شوند نشان می دهد. این متغیر تصادفی ممکن است قسمتی به علت سایر متغیرهایی باشد که بر  $y$  اثر می گذارند ولی مجهولند یا در نظر گرفته نشده اند. الگوی (7-1) بر حسب  $\beta$  ها خطی است؛ لازم نیست بر حسب  $x$  ها خطی باشد. پس الگوهایی به صورت زیر:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_2 + \beta_4 x_2^2 + \epsilon$$

نیز جزء الگوهای خطی به حساب می آیند.

برای برآورد  $\beta$  ها نمونه ای به حجم  $n$  مشاهده روی  $y$  ها و  $x$  های مربوطه را به کار خواهیم برد. این مدل برای مشاهده نام عبارت است از:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + \epsilon_i, \quad i = 1, 2, \dots, n. \quad (8-1)$$

فرض های مربوط به  $\epsilon_i$  یا  $y_i$  مشابه آن هایی است که برای رگرسیون ساده داده شد.

(۱) برای  $i = 1, 2, \dots, n$   $E(\epsilon_i) = 0$  یا به عبارت معادل

$$E(y_i) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik}.$$

(۲) برای  $i = 1, 2, \dots, n$ ،  $Var(\varepsilon_i) = \sigma^2$  یا به عبارت معادل  $Var(y_i) = \sigma^2$

(۳) برای تمام  $i \neq j$ ،  $cov(\varepsilon_i, \varepsilon_j) = 0$  یا به عبارت معادل  $cov(y_i, y_j) = 0$

فرض (۱) بیان می‌کند که به طور متوسط مدل صحیح است یعنی تمام  $x$ های مربوطه منظور شده‌اند و مدل خطی است. (۲) بیان می‌کند که واریانس  $\gamma$  ثابت است و در نتیجه به  $x$ ها بستگی ندارد. فرض (۳) بیان می‌کند که  $\gamma$ ها با یکدیگر ناهم‌بسته‌اند، معمولاً در یک نمونه تصادفی برقرار است (در یک سری زمانی یا وقتی اندازه‌های مکرر از یک طرح یا حیوان گرفته می‌شود، مشاهدات نوعاً وابسته می‌شوند). بعداً فرض نرمال بودن را اضافه خواهیم کرد که تحت آن  $\gamma$ ها علاوه بر ناهم‌بسته بودن مستقل نیز خواهند بود.

وقتی تمام سه فرض برقرار است، برآوردگر کمترین مربعات  $\beta$ ها دارای چندین خاصیت خوب است. اگر یک فرض یا بیشتر برقرار نباشد، برآوردگرها ممکن است ضعیف باشند. تحت فرض نرمال بودن برآوردگرهای درست‌نمایی ماکسیمم دارای خواص عالی هستند.

هر یک از سه شرط ممکن است برای داده‌های حقیقی برقرار نباشد. چندین راه حل برای امتحان فرض‌ها پیشنهاد شده است.

### ۱-۳-۱ برآوردگر کمترین مربعات برای بردار پارامترها

در این بخش از روش کمترین مربعات در برآورد  $\beta$ ها در مدل  $x$  ثابت بحث می‌کنیم. هیچ فرضی درباره توزیع  $\gamma$  برای به دست آوردن برآوردگرها ضروری نیست. برای پارامترهای  $\beta_0, \beta_1, \dots, \beta_k$  برآوردگرهایی را پیدا می‌کنیم که مجموع مربعات انحراف‌های  $n$  مقدار مشاهده شده  $\gamma$  از مقادیر پیش‌بینی شده  $\hat{\gamma}$  را مینیمم سازد. می‌دانیم که  $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$  عبارت (۹-۱) را می‌نیمم می‌سازد.

$$\sum_{i=1}^n \hat{\varepsilon}_i^2 = \sum_{i=0}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \beta_2 x_{i2} - \dots - \hat{\beta}_k x_{ik})^2 \quad (9-1)$$

توجه کنید که مقدار پیش‌بینی شده  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \beta_2 x_{i2} + \dots + \hat{\beta}_k x_{ik}$  برآورد  $E(y_i)$  است؛ نه  $y_i$ ؛ یک نماد بهتر  $E(\hat{y}_i)$  است، ولی عموماً از  $\hat{y}_i$  استفاده می‌شود.

برای به دست آوردن برآوردگرهای کمترین مربعات، لازم نیست که معادله‌ی پیش‌بینی

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \beta_2 x_{i2} + \dots + \hat{\beta}_k x_{ik}$$