

بِسْمِ اللّٰهِ الرَّحْمٰنِ الرَّحِیْمِ



دانشکده ریاضی و رایانه

گروه آمار

پایان نامه تحصیلی برای دریافت درجه کارشناسی ارشد

رشته آمار ریاضی

ماشین بردار پشتیبان و رگرسیون ماشین بردار پشتیبان

با داده های فازی

استاد راهنما:

دکتر علیرضا عربپور

استاد مشاور:

دکتر محمد علی یعقوبی

مؤلف:

محبوبه درانی

شهریورماه ۱۳۹۰

تقدیم به

پشتیبان و یار همیشگی ام مادر مهربانم و همسر عزیزم که مشوقم بوده .

تشکر و قدردانی

خدای مهربان را شاکرم که لطفش بی اندازه و نعمتش گسترده است.

از استاد ارجمند جناب آقای دکتر عربپور که راهنمایی این پایان نامه را بر عهده داشتند، کمال تشکر را دارم و به خاطر تمام زحمات ایشان در سالهای تحصیل سپاسگذارم.

از استاد گرامی جناب آقای دکتر یعقوبی که زحمت مشاوره این پایان نامه را تقبل کردند، کمال تشکر را دارم.

از استادان محترم جناب آقای دکتر امیرزاده و جناب آقای دکتر جمالیزاده که داوری این پایان نامه را بر عهده داشتند، تشکر می کنم.

چکیده

یکی از روش های دسته بندی داده ها ماشین بردار پشتیبان است، همچنین ممکن است داده های موجود مبهم باشند که به عنوان مجموعه های فازی توصیف می شوند. در این پایان نامه ابتدا ماشین بردار پشتیبان معرفی شده است سپس مفهوم مجموعه های فازی را وارد ماشین بردار پشتیبان و رگرسیون ماشین بردار پشتیبان کرده ایم. همچنین ضرایب مدل رگرسیون ماشین بردار پشتیبان فازی برآورد شده است.

در فصل اول این پایان نامه مفاهیم و تعاریف مقدماتی مورد نیاز آورده شده است. در فصل دوم روش ماشین بردار پشتیبان و رگرسیون ماشین بردار پشتیبان معرفی شده است. در فصل سوم مروری بر تاریخچه رگرسیون فازی داریم و به بیان روش های برنامه ریزی ریاضی می پردازیم. در فصل چهارم ماشین بردار پشتیبان فازی و مدل رگرسیون ماشین بردار پشتیبان فازی شرح داده شده است و در نهایت در فصل پنجم نتایج عددی بیان شده است.

کلمات کلیدی : ماشین بردار پشتیبان، رگرسیون ماشین بردار پشتیبان، رگرسیون فازی، ماشین بردار پشتیبان فازی.

فصل اول: مفاهیم و تعاریف مقدماتی

۱-۱ مجموعه های فازی	۲
۱-۱-۱ مقدمه	۲
۱-۱-۲ تعاریف اولیه	۲
۱-۱-۳ اعداد فازی	۳
۱-۱-۴ عملگرهای جبری	۴
۲-۱ رگرسیون کلاسیک	۵
۱-۲-۱ مقدمه	۵
۱-۲-۲ رگرسیون خطی ساده	۵
۱-۲-۳ رگرسیون خطی چند گانه	۹
۳-۱ تحلیل ممیزی	۱۲
۱-۳-۱ جدا سازی و رده بندی برای دو جامعه: روش فیشر	۱۲
۴-۱ تعاریف ریاضی	۱۶

فصل دوم: ماشین بردار پشتیبان و رگرسیون ماشین بردار پشتیبان

- ۲-۱ مقدمه ۲۱
- ۲-۲ ماشین بردار پشتیبان (SVM) ۲۲
- ۲-۳ رگرسیون ماشین بردار پشتیبان ۲۸
- ۲-۴ رگرسیون ماشین بردار پشتیبان با انواع توابع زیان ۳۳
- ۲-۴-۱ رگرسیون خطی با تابع زیان اپسیلون-غیر حساس ۳۳
- ۲-۴-۲ رگرسیون خطی با تابع زیان درجه دو ۳۴
- ۲-۴-۳ رگرسیون خطی با تابع زیان هیوبر ۳۵
- ۲-۴-۴ رگرسیون غیر خطی با تابع زیان اپسیلون-غیر حساس ۳۵

فصل سوم: رگرسیون فازی

- ۳-۱ مقدمه ۳۹
- ۳-۲ تاریخچه ۴۱
- ۳-۳ روش های برنامه ریزی ریاضی ۴۳
- ۳-۳-۱ مدل تاناکا و واتادا ۴۴
- ۳-۳-۲ مدل حجتی و همکاران ۴۶

۳-۳-۳ مدل حسن پور و همکاران ۴۸

فصل چهارم: ماشین بردار پشتیبان فازی و رگرسیون ماشین بردار پشتیبان فازی

۴-۱ مقدمه ۵۲

۴-۲ ماشین بردار پشتیبان فازی (FSVM) ۵۲

۴-۳ رگرسیون خطی SVM فازی ۵۶

۴-۴ رگرسیون خطی SVM فازی با ورودی فازی- خروجی فازی ۶۰

۴-۵ رگرسیون خطی SVM فازی با ورودی غیر فازی- خروجی فازی ۶۴

۴-۶ رگرسیون غیر خطی SVM فازی ۶۸

فصل پنجم: نتایج عددی ۷۲

نتیجه گیری و پیشنهادات ۸۱

ضمیمه: برنامه های R ۸۲

منابع ۱۱۶

فصل اول

تعاریف و مفاهیم مقدماتی

۱- مجموعه های فازی

۱-۱-۱ مقدمه

نظریه ی مجموعه های فازی برای نخستین بار در سال ۱۹۶۵ توسط دکتر لطفی عسکرزاده [۶۵] دانشمندی ایرانی تبار و استاد دانشگاه برکلی به دنیای علم عرضه شد. محققین زیادی از آن تاریخ در توسعه این مفهوم تلاش کرده و آن را در زمینه های مختلف بکار بسته اند.

در این بخش مفاهیم و تعاریف اولیه مجموعه های فازی از قبیل: مفهوم مجموعه های فازی، اعداد فازی و... را بیان خواهیم نمود، که در فصل های بعدی به کار گرفته می شود.

۱-۱-۲ تعاریف اولیه

مفهوم مجموعه و نظریه ی مجموعه ها ابزار های قوی در ریاضیات هستند. در نظریه ی مجموعه - ها، مجموعه به صورت گردایه ای از اشیاء کاملاً مشخص تعریف می شود و عضویت یا عدم عضویت یک شی در مجموعه قطعی است، در بسیاری از مسائل واقعی حدود مجموعه ها کاملاً مشخص نیست برای مجموعه هایی که حدودشان کاملاً مشخص نیست مفهوم مجموعه فازی و عضویت جزئی مطرح می شود [۲].

تعریف ۱-۱-۲-۱: فرض کنید X یک مجموعه ناتهی باشد. هر زیر مجموعه فازی مانند \tilde{A} از X توسط یک تابع $[\cdot, 1]$ به عنوان تابع عضویت مشخص می شود که در آن برای هر x در X ، مقدار $\mu_{\tilde{A}}(x)$ میزان عضویت x در آن زیر مجموعه فازی را نشان می دهد [۶].

به عنوان مثال $\mu_{\tilde{A}}(x) = 0$ یعنی x قطعاً متعلق به A نیست و $\mu_{\tilde{A}}(x) = 1$ یعنی x بدون شک متعلق به A است.

اگر برای هر $x \in X$ ، $\mu_{\tilde{A}}(x) = 0$ یا $\mu_{\tilde{A}}(x) = 1$ در اینصورت \tilde{A} یک مجموعه غیرفازی است و تابع عضویت به تابع نشانگر مجموعه A تبدیل می گردد.

تعریف ۱-۱-۲: تکیه گاه^۱ مجموعه ی فازی \tilde{A} از X که با نماد $S_{\tilde{A}}$ نشان می دهیم به صورت زیر تعریف می شود:

$$S_{\tilde{A}} = \{x \in X | \mu_{\tilde{A}}(x) > 0\}.$$

تعریف ۱-۱-۳: برای هر $\alpha \in [0, 1]$ ، α -برش^۲ یک مجموعه ی فازی \tilde{A} از X که با نماد \tilde{A}_{α} نمایش داده می شود به صورت زیر تعریف می گردد:

$$\tilde{A}_{\alpha} = \begin{cases} \{x \in X | \mu_{\tilde{A}}(x) \geq \alpha\}, & \forall \alpha > 0, \\ \{x \in X | \mu_{\tilde{A}}(x) > 0\}, & \alpha = 0. \end{cases}$$

۱-۱-۳ اعداد فازی

تعریف ۱-۱-۳: کمیت فازی \tilde{A} را عدد فازی $L - R$ گوئیم هر گاه تابع عضویت \tilde{A} به صورت زیر باشد:

$$\mu_{\tilde{A}}(x) = \begin{cases} L\left(\frac{m-x}{\alpha}\right), & x \leq m, (\alpha > 0), \\ R\left(\frac{x-m}{\beta}\right), & x \geq m, (\beta > 0). \end{cases}$$

^۱ support
^۲ α -level

که $L \text{ و } R : [0, \infty) \rightarrow [0, 1]$ توابعی پیوسته، نزولی و روی بازه $[0, 1]$ معکوس پذیر هستند. همچنین به m ، α و β $L(0)=R(0)=1$ و $L(1)=R(1)=0$ ، به توابع L و R توابع مرجع می‌گوییم. برای سادگی عدد فازی \tilde{A} را با نماد $\tilde{A} = (m, \alpha, \beta)_{LR}$ نشان می‌دهیم. اگر $L=R$ و $\alpha=\beta$ ، \tilde{A} را عدد فازی متقارن می‌گوییم و آن را با نماد $\tilde{A} = (m, \alpha)_L$ نشان می‌دهیم.

تعریف ۱-۱-۲: به عدد فازی $L - R$ ، \tilde{A} با توابع مرجع $L(x) = R(x) = \max(0, 1-x)$ عدد فازی مثلثی می‌گوییم و آن را با $\tilde{A} = (m, \alpha, \beta)$ نشان می‌دهیم [۶].

۱-۱-۴ عملگرهای جبری

اگر $\tilde{A}_1 = (m_1, \alpha_1, \beta_1)$ و $\tilde{A}_2 = (m_2, \alpha_2, \beta_2)$ دو عدد فازی مثلثی باشند آنگاه:

$$۱) \quad k \tilde{A}_1 = \begin{cases} (km_1, k\alpha_1, k\beta_1), & k > 0, \\ (km_1, |k|\beta_1, |k|\alpha_1), & k < 0. \end{cases}$$

$$۲) \quad \tilde{A}_1 + \tilde{A}_2 = (m_1 + m_2, \alpha_1 + \alpha_2, \beta_1 + \beta_2),$$

$$۳) \quad \tilde{A}_1 - \tilde{A}_2 = \tilde{A}_1 + (-\tilde{A}_2) = (m_1 - m_2, \alpha_1 + \beta_2, \alpha_2 + \beta_1).$$

در ادامه به مفاهیم مقدماتی رگرسیون کلاسیک می‌پردازیم که در فصل‌های بعدی مورد استفاده قرار می‌گیرد.

۲-۱ رگرسیون کلاسیک

۲-۱-۱ مقدمه

بسیاری از پدیده‌ها و پیشامدها یی که در جهان اتفاق می‌افتند تابع برخی از پیشامدهای دیگر می‌باشند. یکی از هدفهای مهم در تحقیقات علمی، کشف وجود یا عدم وجود رابطه بین پدیده - هاست. یکی از روشهایی که به طور گسترده برای این منظور به کار می‌رود روش تحلیل رگرسیون می‌باشد [۸]. در واقع رگرسیون روشی است برای بررسی رابطه‌ی بین دو یا چند متغیر به طوری که یک متغیر از روی دیگری یا بقیه می‌تواند پیش‌بینی شود.

۲-۲-۱ رگرسیون خطی ساده

در ساده‌ترین حالت در یک مدل رگرسیون فرض می‌کنیم y متغیر وابسته و X متغیر مستقل باشد و این رابطه به صورت زیر باشد:

$$y = \beta_0 + \beta_1 x. \quad (1-1)$$

با توجه به این که نقاط دقیقاً روی خط راست در یک امتداد قرار نمی‌گیرند رابطه‌ی (۱-۱) باید تعدیل شود.

بنابراین مدل (۱-۱) به این صورت اصلاح می‌گردد.

$$y = \beta_0 + \beta_1 x + \varepsilon. \quad (2-1)$$

که ε اختلاف بین مقدار y (متغیر وابسته) و خط $\beta_0 + \beta_1 x$ را بیان می کند و مقدار خطای آماری نامیده می شود. مدل اصلاح شده ی (۱-۲) مدل رگرسیون خطی ساده نامیده می شود.

اعتبار یک مدل رگرسیونی به برقرار بودن شرایط زیر بستگی دارد:

۱- میانگین خطاها برابر صفر و واریانس خطاها مقدار ثابتی باشد.

$$E(\varepsilon_i) = 0, \quad \text{Var}(\varepsilon_i) = \sigma^2, \quad i = 1, 2, \dots, n.$$

۲- خطاها نسبت به هم ناهمبسته باشند.

$$\text{Cov}(\varepsilon_i, \varepsilon_j) = 0, \quad i, j = 1, 2, \dots, n, \quad i \neq j.$$

۳- خطاها دارای توزیع نرمال با میانگین صفر و واریانس ثابت باشند.

$$\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2), \quad i = 1, 2, \dots, n.$$

با استفاده از روش حداقل مربعات خطا می توان پارامترهای مدل یعنی β_0 و β_1 را برآورد کرد [۸].

فرض می کنیم داده ها به صورت زوجی (x_1, y_1) و (x_2, y_2) و ... و (x_n, y_n) وجود دارند،

خطای هر مشاهده به این صورت است:

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

در این روش پارامترها را به گونه ای برآورد می کنیم که مجموع مربعات خطا مینیمم شود.

به عبارتی تابع زیر را مینیمم می کنیم :

$$S(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2.$$

در صورتی که $\hat{\beta}_0$ و $\hat{\beta}_1$ برآورد پارامترهای β_0 و β_1 به روش کمترین مربعات خطا باشند آنگاه

باید داشته باشیم :

$$\frac{\partial S(\beta_0, \beta_1)}{\partial \beta_0} \Big|_{\hat{\beta}_0, \hat{\beta}_1} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0,$$

$$\frac{\partial S(\beta_0, \beta_1)}{\partial \beta_1} \Big|_{\hat{\beta}_0, \hat{\beta}_1} = -2 \sum_{i=1}^n x_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0.$$

با ساده کردن روابط بالا داریم:

$$n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i = \sum_{i=1}^n y_i,$$

$$\hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i y_i.$$

با حل معادلات فوق که معادلات نرمال نامیده می شوند برآورد پارامترها به این صورت به دست

می آیند:

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}, \quad \hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}.$$

$$S_{xy} = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right),$$

$$S_{xx} = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2,$$

$$\bar{y} = \frac{1}{n} \left(\sum_{i=1}^n y_i \right),$$

$$\bar{x} = \frac{1}{n} \left(\sum_{i=1}^n x_i \right).$$

پس از برآورد پارامترها، مدل رگرسیونی به صورت زیر در می آید :

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x.$$

۳-۲-۱ رگرسیون خطی چندگانه

در رگرسیون خطی چندگانه می خواهیم متغیر وابسته Y را براساس متغیرهای مستقل X_1, X_2, \dots, X_k پیش بینی کنیم، برای این منظور مدل خطی زیر پیشنهاد می شود:

$$y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + \varepsilon.$$

که همانند رگرسیون خطی ساده، ε خطای آماری نامیده می شود و اختلاف بین Y و رابطه ی خطی $\beta_0 + \beta_1 X_1 + \dots + \beta_k X_k$ می باشد.

حال بر اساس یک نمونه ی n تایی $(X_{i1}, X_{i2}, \dots, X_{ik}, Y_i)$ $i=1, 2, \dots, n$ پارامترهای مدل یعنی $\beta_0, \beta_1, \dots, \beta_k$ را به روش کمترین مربعات خطا برآورد می کنیم. برای این منظور باید تابع زیر را مینیمم کنیم.

$$S(\beta_0, \beta_1, \dots, \beta_k) = \sum_{i=1}^n (\varepsilon_i)^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 X_{i1} - \dots - \beta_k X_{ik})^2.$$

در صورتی که $\beta_0, \beta_1, \dots, \beta_k$ برآورد کمترین مربعات خطای پارامترها باشند آن گاه باید داشته باشیم:

$$\frac{\partial S(\beta_0, \beta_1, \dots, \beta_k)}{\partial \beta_0} \Big|_{\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{i1} - \dots - \hat{\beta}_k X_{ik}) = 0,$$

$$\frac{\partial S(\beta_0, \beta_1, \dots, \beta_k)}{\partial \beta_1} \Big|_{\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k} = -2 \sum_{i=1}^n X_{i1} (y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{i1} - \dots - \hat{\beta}_k X_{ik}) = 0,$$

⋮

$$\frac{\partial S(\beta_0, \dots, \beta_k)}{\partial \beta_k} \Big|_{\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k} = -2 \sum_{i=1}^n x_{ik} (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \dots - \hat{\beta}_k x_{ik}) = 0.$$

باساده کردن روابط فوق داریم

$$n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_{i1} + \dots + \hat{\beta}_k \sum_{i=1}^n x_{ik} = \sum_{i=1}^n y_i,$$

$$\hat{\beta}_0 \sum_{i=1}^n x_{i1} + \hat{\beta}_1 \sum_{i=1}^n x_{i1}^2 + \dots + \hat{\beta}_k \sum_{i=1}^n x_{i1} x_{ik} = \sum_{i=1}^n x_{i1} y_i,$$

$$\hat{\beta}_0 \sum_{i=1}^n x_{ik} + \hat{\beta}_1 \sum_{i=1}^n x_{i1} x_{ik} + \dots + \hat{\beta}_k \sum_{i=1}^n x_{ik}^2 = \sum_{i=1}^n x_{ik} y_i.$$

با حل $p = k+1$ معادله ی فوق که معادلات نرمال نامیده می شوند $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$ به دست می آیند.

برای سادگی کار وانجام محاسبات مدل رگرسیون چندگانه را به شکل دیگری نیز می توان نوشت.

تعریف می کنیم

$$\mathbf{y}_{n \times 1} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \quad \boldsymbol{\varepsilon}_{n \times 1} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}, \quad \boldsymbol{\beta}_{p \times 1} = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_k \end{pmatrix}, \quad p = k + 1,$$

$$\mathbf{X}_{n \times p} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{nk} \end{pmatrix}.$$

با توجه به تعاریف بیان شده شکل ماتریسی مدل رگرسیون چندگانه را به صورت زیر می توان

نوشت:

$$\mathbf{y}_{n \times 1} = \mathbf{X}_{n \times p} \boldsymbol{\beta}_{p \times 1} + \boldsymbol{\varepsilon}_{n \times 1}, \quad \boldsymbol{\varepsilon}_{n \times 1} \sim N_n(\mathbf{0}_{n \times 1}, \sigma^2 \mathbf{I}).$$

$$\mathbf{y}_{n \times 1} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}) \quad \text{بنابراین}$$

و در نتیجه برآورد کمترین مربعات خطا به صورت زیر به دست می آید.

$$S(\beta_0, \beta_1, \dots, \beta_k) = \boldsymbol{\varepsilon}' \boldsymbol{\varepsilon} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

$$\frac{\partial S(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \Big|_{\hat{\boldsymbol{\beta}}} = \mathbf{0} \mapsto \mathbf{X}' \mathbf{y} = \mathbf{X}' \mathbf{X} \hat{\boldsymbol{\beta}}$$

منظور از \mathbf{X}' ترانزاده ی ماتریس \mathbf{X} می باشد.

اگر $(\mathbf{X}' \mathbf{X})$ معکوس پذیر باشد برآورد کمترین مربعات $\boldsymbol{\beta}$ عبارت است از:

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y}.$$

بنابراین برآورد مدل رگرسیون به صورت $\hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}}$ می باشد که

$$\hat{\mathbf{y}} = \begin{pmatrix} \hat{\beta}_0 + x_{11} \hat{\beta}_1 + \dots + x_{1k} \hat{\beta}_k \\ \hat{\beta}_0 + x_{21} \hat{\beta}_1 + \dots + x_{2k} \hat{\beta}_k \\ \vdots \\ \hat{\beta}_0 + x_{n1} \hat{\beta}_1 + \dots + x_{nk} \hat{\beta}_k \end{pmatrix}, \quad \hat{\boldsymbol{\beta}} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_k \end{pmatrix}.$$

در ادامه به تعریف روش تحلیل ممیزی یکی از روش های دسته بندی و رده بندی می پردازیم که در فصل های بعد مورد استفاده قرار خواهد گرفت.

۱-۳ تحلیل ممیزی

تحلیل ممیزی و رده بندی تکنیک های چند متغیری هستند که با جدا کردن مجموعه های متمایز مشاهدات و تخصیص دادن مشاهدات جدید به دسته های تعریف شده ی قبلی سروکار دارد. در تحلیل ممیزی تاکید روی به دست آوردن قاعده ای است که از آن بتوان برای تخصیص بهینه ی یک شیئی جدید به رده های مشخص استفاده کرد. اصطلاح ممیزی توسط فیشر^۳ [۳۰] در اولین کار پیشرفته ی مسائل جداسازی معرفی گردیده است. در زیر روش ممیزی پیشنهاد شده توسط فیشر بیان می شود.

۱-۳-۱ جداسازی و رده بندی برای دو جامعه (روش فیشر)

در این بخش دو رده (جامعه) را با π_1 و π_2 مشخص می کنیم. مشاهدات به طور معمول بر مبنای اندازه هایی مانند p متغیر تصادفی (X_1, \dots, X_p) ، جداسازی یا رده بندی می شوند. به عنوان مثال جامعه های π_1 و π_2 می تواند دو گونه علف هرز و متغیر X درازای کاسبرگ و گلبرگ، عمق شکاف گلبرگ و ... باشد.

فیشر در ابتدا ترکیبات خطی از X را برای ایجاد Y ها پیشنهاد می کند، زیرا آن ها توابع ساده ای از X هستند و کار کردن با آن ها از نظر ریاضی آسان تر است. اگر μ_{1Y} میانگین Y های به دست آمده از X های متعلق به π_1 و μ_{2Y} میانگین Y های به دست آمده از X های متعلق به π_2 باشد، آن

^۳ Fisher