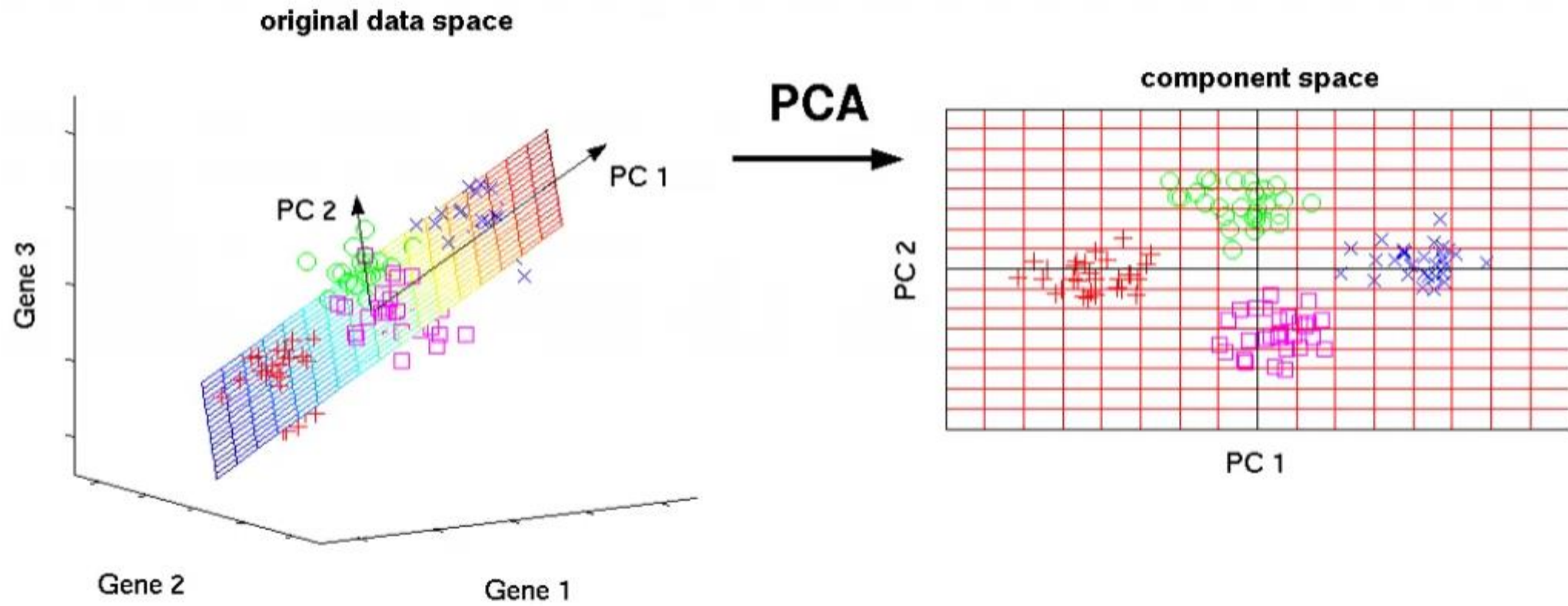


# ضروریات آمار برای علم داده

تحلیل مولفه اصلی

تحلیل عاملی



## تحلیل مؤلفه اصلی

یکی از روش‌های تحلیل چند متغیره است

هدف اصلی آن تقلیل بعد مسئله مورد مطالعه است

از کاربردهای مهم تحلیل مؤلفه اصلی در رگرسیون است

با استفاده از این روش تعداد زیادی متغیر توضیحی هم‌بسته با تعداد محدودی متغیر توضیحی جدید که مولفه اصلی نامیده می‌شوند و ناهمبسته هستند جایگزین می‌شود؛ به این ترتیب نه تنها بعد مسئله کاهش می‌یابد بلکه مسئله هم‌خطی نیز حل می‌شود.

## تحلیل مؤلفه اصلی

- ❖ انجام بهترین بصری سازی از داده های دارای ابعاد بالا
- ❖ برای غلبه بر افزونگی متغیرهای موجود
- ❖ روی مجموعه داده های دارای ویژگی های عددی قابل اعمال است
- ❖ مولفه های این متغیرها حاصل ترکیب خطی نرمال شده متغیرهای پیش بین اصلی هستند.
- ❖ هدف این مولفه ها حفظ بیشترین اطلاعات ممکن با واریانس های بالا است
- ❖ اولین مولفه اصلی بالاترین واریانس را داراست و پس از آن مولفه اصلی دوم دارای بیشترین مقدار واریانس است و این موضوع برای مولفه های اصلی سوم و دیگر مولفه های اصلی نیز صادق است.
- ❖ مولفه ها باید ناهمبسته باشند (جهت های آنها متعامد است).
- ❖ نرمال سازی داده ها هنگامی که متغیرها دارای واحدهای (یکاهای) گوناگونی هستند، فوق العاده مهم است

# Principal component Analysis

$$z = \frac{\text{value} - \text{mean}}{\text{standard deviation}}$$

$$\begin{bmatrix} \text{Cov}(x, x) & \text{Cov}(x, y) & \text{Cov}(x, z) \\ \text{Cov}(y, x) & \text{Cov}(y, y) & \text{Cov}(y, z) \\ \text{Cov}(z, x) & \text{Cov}(z, y) & \text{Cov}(z, z) \end{bmatrix}$$

Compute the eigenvectors and eigenvalues of the covariance matrix

$$\text{FinalDataSet} = \text{FeatureVector}^T * \text{StandardizedOriginalDataSet}^T$$

# مفاهيم جبري

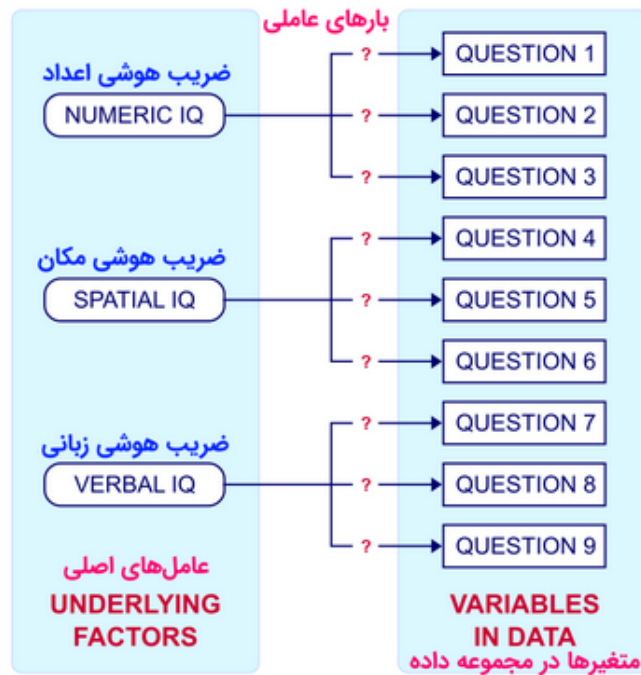
$$M \text{ positive semi-definite} \iff z^T M z \geq 0 \text{ for all } z \in \mathbb{R}^n \setminus \mathbf{0}$$

$$\det(A - \lambda I) = 0$$

$$A = \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} \quad |A| = a(ei - fh) - b(di - fg) + c(dh - eg)$$

$$(A - \lambda I) \mathbf{v} = \mathbf{0}$$

# تحلیل عاملی



هدف اصلی تحلیل عاملی در صورت امکان بیان روابط کوواریانس میان بسیاری از متغیرها براساس چند کمیت تصادفی غیرقابل مشاهده است که عامل‌ها نامیده می‌شوند.

فرض کنید متغیرها را بتوان براساس همبستگی‌شان دسته‌بندی کرد یعنی تمام متغیرهای یک گروه خاص در میان خودشان همبستگی بالایی دارند ولی با متغیرهای یک گروه دیگر همبستگی نسبتاً کمی دارند؛ می‌توان گفت هرگروه از متغیرها یک ترکیب یا عامل مورد نظر را نشان می‌دهند که نشان‌دهنده همبستگی‌های مشاهده شده است.

## تحليل عاملی

$$x_{i,m} - \mu_i = l_{i,1}f_{1,m} + \cdots + l_{i,k}f_{k,m} + \varepsilon_{i,m}$$

whereby

- $x_{i,m}$  is the value of the  $i$ th observation of the  $m$ th individual,
- $\mu_i$  is the observation mean for the  $i$ th observation,
- $l_{i,j}$  is the loading for the  $i$ th observation of the  $j$ th factor,
- $f_{j,m}$  is the value of the  $j$ th factor of the  $m$ th individual, and
- $\varepsilon_{i,m}$  is the  $(i, m)$ th *unobserved stochastic error term* with mean zero and finite variance.