

Veri Analizi - Ders 1

Gaussian Kontrolu

Diyelim ki Gaussian dagilimina sahip oldugunu dusundugumuz $\{x_i\}$ verilerimiz var. Bu verilerin Gaussian dagilimina uyup uymadigini nasil kontrol edecegiz? Normal bir dagilimin her veri noktası için şöyle temsil edebiliriz,

$$y_i = \Phi\left(\frac{x_i - \mu}{\sigma}\right)$$

Burada Φ standart Gaussian'ı temsil ediyor (detaylar için *Istatistik Ders 1*) ve CDF fonksiyonuna tekabül ediyor. CDF fonksiyonunun aynı zamanda ceyregi (quantile) hesapladığı söylenir, aslında CDF son derece detaylı bir olasılık değeri verir fakat evet, dolaylı yoldan noktanın hangi ceyrek içine dustugu de görülecektir.

Simdi bir numara yapalım, iki tarafa ters Gaussian formülünü uygulayalım, yani Φ^{-1} .

$$\Phi^{-1}(y_i) = \Phi^{-1}\left(\Phi\left(\frac{x_i - \mu}{\sigma}\right)\right)$$

$$\Phi^{-1}(y_i) = \frac{x_i - \mu}{\sigma}$$

$$x_i = \Phi^{-1}(y_i)\sigma + \mu$$

Bu demektir ki elimizdeki verileri $\Phi^{-1}(y_i)$ bazında grafiklersek, bu noktalar eğimi σ , başlangıcı (intercept) μ olan bir düz çizgi olmalıdır. Eğer kabaca noktalar düz çizgi oluşturmuyorsa, verimizin Gaussian dağılıma sahip olmadığına karar verebiliriz.

Ustte tarif edilen grafik, olasılık grafiği (probability plot) olarak bilinir.

Ters Gaussian teorik fonksiyonunu burada vermeyeceğiz, Scipy `scipy.stats.invgauss` hesaplar için kullanılabilir. Fakat y_i 'nin kendisi nereden geliyor? Eğer y_i , CDF'in bir sonucu ise, pur veriye bakarak bir CDF değeri de hesaplayabilmemiz gerekir. Bunu yapmak için bir başka numara lazım.

1. Eldeki sayıları artan şekilde sıralayın
2. Her veri noktasına bir derece (rank) atayın (sıralama sonrası hangi seviyede olduğu yeterli, 1'den başlayarak).
3. Ceyrek değeri y_i bu sıra / $n + 1$, n eldeki verinin büyüklüğü.

Bu teknik niye isliyor? x 'in CDF'i $x_i < x$ şartına uyan x_i 'lerin oranı değil midir? Yani bir sıralama söz konusu ve üstteki teknik te bu sıralamayı biz elle yapmış olduk, ve bu sıralamadan gereken bilgiyi aldık.

Ozet İstatistikleri

Genellikle istatistik kitapları hemen ortalama (mean), medyan (median) ve bağlantılı özet istatistiklerinden (summary statistics) bahsederek işe girerler. Bu istatistikleri dikkatle kullanmak gerekir, çünkü her türlü veri, her yerde geçerli değildir. Mesela ortalama sadece tek merkezi bir tepesi olan (unimodal) dağılımlar için geçerlidir.

Eger bu temel varsayim gecerli degilse, ortalama kullanarak yapilan hesaplar bizi yanlis yollara goturur. Bu uyaridan sonra ortalama ve standart sapmayi (standart deviation) gorelim.

$$m = \frac{1}{n} \sum_i x_i$$

Standart sapma veri noktalarin “ortalamadan farkinin ortalamasini” verir. Tabii bazen noktalar ortalamanin altinda, bazen ustunde olacaktir, bizi bu negatiflik, pozitiflik ilgilendirmez, biz sadece farkla alakaliyiz. O yuzden her sapmanin karesini aliriz, bunlari toplayip nokta sayisina boluruz .

$$s^2 = \frac{1}{n} \sum_i (x_i - m)^2$$

Eger m tanimini ustte yerine koyarsak,

$$\begin{aligned} &= \frac{1}{n} \sum_i x_i^2 + \frac{1}{n} \sum_i m^2 - \frac{2}{n} \sum_i x_i m \\ &= \frac{1}{n} \sum_i x_i^2 + \frac{m^2 n}{n} - \frac{2mn}{n} m \\ &= \frac{1}{n} \sum_i x_i^2 + m^2 - 2m^2 \\ &= \frac{1}{n} \sum_i x_i^2 - m^2 \end{aligned}$$

Bu olcuye varyans (variance) denir ve teorik olarak ortalamadan daha onemli oldugu soylenebilir. Fakat dagilimin yayilma olcusu olarak biz bu olcuyu oldugu gibi degil, onun karesini kullanacagiz (ki standart sapma buna deniyor aslinda). Niye? Cunku o zaman veri noktalarinin ve yayilma olcusunun birimleri birbiri ile ayni olacak. Eger veri setimiz bir alisveris sepetindeki malzemelerin lira cinsinden degerleri olseydi, varyans bize sonucu “karekok lira” olarak verecekti ve bunun pek anlami olmayacakti.