

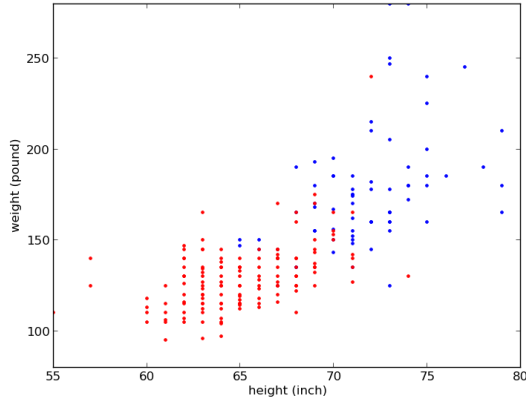
## Karisimler ve Idare Edilmeyen Kumeleme (Unsupervised Clustering)

Gaussian (normal) dagilimi tek tepesi olan (unimodal) bir dagilimdir. Bu demektir ki eger birden fazla tepe noktası olan bir veriyi modellemek istiyorsak, degisik yaklasimlar kullanmamiz gerekecektir.

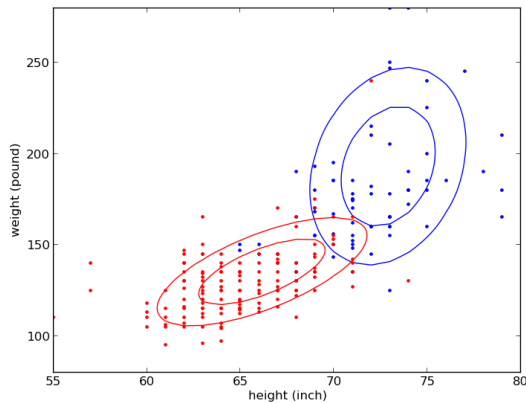
Birden fazla Gaussian'i "karistirmek (mixing)" bu tur bir yaklasim olabilir. Karistirmek, karisim icindeki her Gaussian'dan gelen sonuclari toplamaktir, yani kelimenin tam anlamıyla her veri noktasini teker teker karisimdaki tum dagilimlara gecip sonuclari toplamaktir. Eger cok boyutlu normal dagilimleri topluyorsak, formül:

$$p(x) = \sum_z \pi_z N(x|\mu_z, \Sigma_z)$$

$\pi_z$  karistirma oranlaridir (mixing proportions). Iki Gaussian oldugunu dusunelim, oranlar 0.2, 0.8 olabilir mesela (toplam her zaman 1 olmalidir). Karisim oranlarına degisik bir bakis acisini simdi isleyecegiz. Ornek olarak alttaki grafige bakalim.



Bu grafik kadınlar ve erkeklerin boy ve kilolarini iceren bir veri setinden geliyor, veri setinde erkekler ve kadınlara ait olan olcumler isaretlenmis, biz de bu isaretleri kullanarak kadınlari kirmizi erkekleri mavi ile grafikledik. Bu isaretler verilmiş olsun ya da olmasın, eger bu veriye bir dagilim uydurmak (fit) istersek, bir karisim kullanmamiz gereklidir. Karisimi olusturan Gaussian'lar mesela su sekilde olabilir



Nihai olasılık değeri  $p(x)$ 'i hesaplamak için her noktayı her iki Gaussian'a teker teker geçiriz, ve sonuçları karışım oranları ile carparak toplarız. Kesisme olmayan bölgelerdeki noktaları düşünürsek, o noktaların olasılık değeri zaten ağırlıkla tek bir Gaussian'dan geliyor olacak (diğer Gaussian o bölge için sifıra yakın bir değer verecek, bu değer toplamda bir fark yaratmayacaktır). Kesisim olan bölgelerdeki noktalar ise ağırlıklara göre carpilip toplanacak. Ağırlığı fazla olan Gaussian bu bölgeler için anlamlı, demek ki o bölgede o Gaussian daha yoğun noktalara sahip (bu yüzden o bölgede daha fazla olasılık veriyor olması gerekir).

