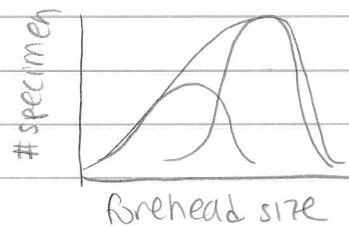# Mixture Models

Introduced by Karl Pearson in 1894

Naples crab:



Are there two species? Is the distribution a mixture of two Gaussians?

Recall: A Gaussian has pdf

$$N(\mu, \sigma^2, x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

A mixture has pdf

$$F(x) = w N(\mu_1, \sigma_1^2, x) + (1-w) N(\mu_2, \sigma_2^2, x)$$

where $0 \leq w \leq 1$ is the mixing weight.

Interpretation: Flip a biased coin to determine which component sample comes from

Other applications include modeling height, velocities in gasses etc

Pearson invented the method of moments
to attack the learning problem

def: Let $M_r = \mathbb{E}_{x \in F(x)}[x^r]$

Fact: $M_r$ is a polynomial in the unknown
parameters, i.e. $(w, \mu_1, \sigma_1^2, \mu_2, \sigma_2^2)$

In particular

① $M_1 = w\mu_1 + (1-w)\mu_2$

② $M_2 = w(\mu_1^2 + \sigma_1^2) + (1-w)(\mu_2^2 + \sigma_2^2)$

③ $M_3 = w(\mu_1^3 + 3\mu_1\sigma_1^2) + (1-w)(\mu_2^3 + 3\mu_2\sigma_2^2)$

etc

Let $\tilde{M}_r = \frac{1}{|S|} \sum_{i \in S} x_i^r$ denote empirical avgs

          ↑
        samples

Sixth Moment Test

- Given samples $S$, compute $\tilde{M}_r$ for $r = 1$ to $6$

- Solve for simultaneous roots of

$$\left\{ M_r(w, \mu_1, \sigma_1^2, \mu_2, \sigma_2^2) = \tilde{M}_r \right\}_{r = 1 \text{ to } 5}$$

- Among all valid solns, choose the one that

is closest in sixth moment

Main Questions:

① Is it stable to sampling noise?

② Do the first six moments uniquely determine the parameters?

Milestones

Pearson (1894): method of moments
(no guarantees)

Fisher (1912-1922): maximum likelihood estimator

$$\hat{\Theta}_{MLE} = \arg\max p(x_1, \ldots, x_s; \theta)$$

Consistent and asymptotically efficient, usually computationally hard

Teicher (1961): identifiability

i.e. if we knew the density function exactly, it determines the parameters

Proof [sketch] The component with the largest variance dominates the behavior of $F(x)$ in the tails.

Find its mean, variance and mixing
weight, subtract it off from $F(x)$
and proceed. ▨

Intuitively, this requires tons of samples

Dempster, Laird, Rubin (1977): expectation
$\qquad\qquad\qquad\qquad\qquad$ maximization

① initial guess $(\hat{w}, \hat{\mu}_1, \hat{\sigma}_1^2, \hat{\mu}_2, \hat{\sigma}_2^2)$

② Iterate:

$\quad$ Cluster: For each $x \in S$, calculate posterior

$$P_x = \frac{\hat{w} \, \mathcal{N}(\hat{\mu}_1, \hat{\sigma}_1^2, x)}{\hat{w} \, \mathcal{N}(\hat{\mu}_1, \hat{\sigma}_1^2, x) + (1-\hat{w}) \, \mathcal{N}(\hat{\mu}_2, \hat{\sigma}_2^2, x)}$$

$\quad$ update parameters

$$\hat{w} \leftarrow \frac{\sum\limits_{x \in S} P_x}{|S|} \quad ; \quad \hat{\mu}_1 \leftarrow \frac{\sum\limits_{x \in S} P_x \, x}{\hat{w}}$$

$$\hat{\sigma}_1^2 \leftarrow \frac{\sum\limits_{x \in S} P_x (x - \hat{\mu}_1)^2}{\hat{w}}$$

$\quad$ and similarly for $\hat{\mu}_2, \hat{\sigma}_2^2$

This is a heuristic to maximize likelihood, but often gets stuck

—

## Learning via Clustering

Dasgupta gave the first provable guarantees

Claim: If we can accurately cluster the samples into which component generated them, can estimate the parameters

But how do you cluster?

Let's start with some intuitive and counter-intuitive properties of high-dimensional Gaussians

Recall: $\mathcal{N}(\mu, \Sigma, x) = \dfrac{e^{-\frac{(x-\mu)^T \Sigma^{-1}(x-\mu)}{2}}}{(2\pi)^{d/2} \det(\Sigma)^{1/2}}$

Fact #1: $\mathcal{N}(\mu, \Sigma, x)$ is maximized at $x = \mu$

Fact #2: For $x \sim \mathcal{N}(\mu, \sigma^2 I, x)$

$$\mathbb{P}\left[\left|\|x-\mu\|^2 - \sigma^2 d\right| \geq C\sigma^2 \sqrt{d \log d}\right] \leq d^{-\frac{C^2}{4}}$$

How can these facts simultaneously be true?

The growth of the volume of the ball counteracts the decay in the pdf as we move away from $\mu$

Sketch of Fact #2: First we note if

$$x \sim N(\mu, \sigma^2) \text{ then } bx+a \sim N(b\mu+a, b^2\sigma^2)$$

Now consider                    $i^{th}$ coordinate of random $x$

$$\sum_{i=1}^{d} z_i^2 \text{ where } z_i \triangleq \frac{(x_i - \mu_i)}{\sigma_i}$$

Then $\sum_{i=1}^{d} z_i^2 = \frac{\|x-\mu\|^2}{\sigma^2}$

Now each $z_i \sim N(0,1)$ and $\sum_{i=1}^{d} z_i^2$ is called a $\chi^2$-distribution

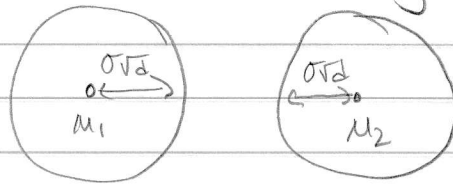It has an explicit expression for its pdf, but all we need is

$$\frac{\sum_{i=1}^{d} z_i^2 - d}{2\sqrt{d}} \xrightarrow{\lim d \to \infty} N(0,1)$$

Hence $\sum_{i=1}^{d} z_i^2 \to N(d, 4d)$

$$\Rightarrow \sum_{i=1}^{d} (x_i - \mu_i)^2 = \|x-\mu\|^2 \to N(\sigma^2 d, 4\sigma^4 d)$$

Finally $\mathbb{P}\left[ \left| \|x-\mu\|^2 - \sigma^2 d \right| > c \, \sigma^2 \sqrt{d \ln d} \right]$

$$\lesssim e^{-\frac{c^2 \sigma^4 d \ln d}{4 \sigma^2 d}} = d^{-\frac{c^2}{4}} \quad \boxed{2}$$
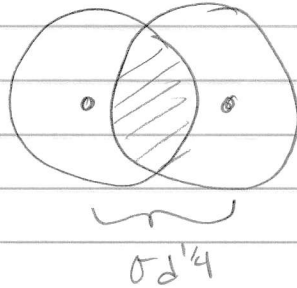
Now back to clustering: If $\|M_1 - M_2\| \gg \sigma\sqrt{d}$



We should be able to cluster [Dasgupta]

[Arora, Kannan]

Proposition: If $\|M_1 - M_2\| \gg d^{1/4}\sigma \ln^{1/2}d$ then whp all samples from first component are closer to each other than to any sample from second component and vice-versa

How can this be? Pictorially



$\sigma d^{1/4}$

The measure of the overlap region is negligible

Proof: Consider $a, a' \sim N(M_1, \sigma^2 I)$ and $b \sim N(M_2, \sigma^2 I)$. Then whp the vectors

$$a - M_1, \quad a' - M_1, \quad M_1 - M_2 \text{ and } b - M_2$$

are nearly orthogonal. This follows b/c three of them are random, so pairwise inner-products are small

Now we can compute

$$\|a-a'\|^2 = \|a - \mu_1 + \mu_1 - a'\|^2$$

$$= \underbrace{\|a-\mu_1\|^2}_{①} + \underbrace{\|\mu_1-a'\|^2}_{②} + \underbrace{2\langle a-\mu_1, \mu_1-a'\rangle}_{③}$$

Now ① and ② are each $\sigma^2 d \pm c\sigma^2\sqrt{d \ln d}$

and ③ is $\dfrac{\sigma^2 d}{\sqrt{d}}$ → negligible

Thus $\|a-a'\|^2 = 2\sigma^2 d \pm 2c\sigma^2\sqrt{d \ln d}$

Similarly we have

$$\|a-b\|^2 = \|a-\mu_1 + \mu_1 - \mu_2 + \mu_2 - b\|^2$$

$$= \|a-\mu_1\|^2 + \|\mu_1-\mu_2\|^2 + \|\mu_2-b\|^2$$

$$\pm \text{ lower order terms}$$

$$= 2\sigma^2 d \pm 2c\sigma^2\sqrt{d \ln d} + \sigma^2\sqrt{d}\ln d$$

Hence we have

$$\|a-b\|^2 \geq \|a-a'\|^2 + \sigma^2\sqrt{d}\ln d \quad \text{whp}$$

$\boxminus$

This gives a polynomial running time / sample complexity algorithm at separation $\gtrsim d^{1/4}$

Is this the best we can do?

def. the total variation distance btwn
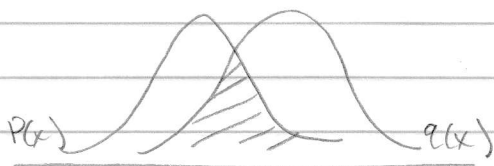pdfs $p(x)$ and $q(x)$ is

$$d_{TV} \triangleq \frac{1}{2} \int |p(x) - q(x)| \, dx$$

Fact: Clustering $w(1)$ samples, requires (with $w = \frac{1}{2}$)

$$d_{TV}(N(\mu_1, \Sigma_1), N(\mu_2, \Sigma_2)) \geq 1 - o(1)$$

Proof: we can couple samples from the two
distributions, e.g.



$p(x)$          $q(x)$

Throw darts and output samples if they
are below the pdf.

first dart below $p(x) \sim p(x)$
first dart below $q(x) \sim q(x)$

Now to sample from the mixture

① throw a dart

② flip a ~~biased~~ coin. on heads, output
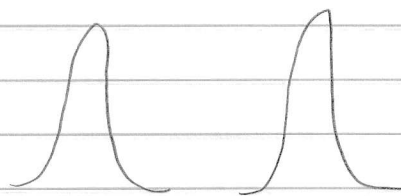the sample if it's below $p(x)$. on tails

output the sample if it's below $q(x)$

Notice that if the dart lands in the overlap region, which has area $1 - d_{TV}$, then it could have come from either. ☒

For what separation is $d_{TV} = 1 - o(1)$?

This holds even when $\|\mu_1 - \mu_2\| \gg \sigma \sqrt{\ln d}$, and this is tight

e.g. if we could project on the line connecting $\mu_1$ and $\mu_2$ we'd get



Main Question: How can we find the right directions to project on?

Following [Vempala, Wang], let

$$M = \mathbb{E}_{x \sim F(x)}[x x^T]$$

Lemma: Let $u_1, \ldots, u_k$ be the top $k$ singular vectors of $M$. If $F(x)$ is a mixture of $k$ spherical Gaussians with means $\mu_1, \ldots, \mu_k$ then
linearly independent

$$\text{span}(u_1 \ldots u_k) = \text{span}(M_1 \ldots M_k)$$

Proof: We can write $X = C + Z$ where

$$C = \begin{cases} M_1 & \text{w/ prob } w_1 \\ \vdots \\ M_k & \text{w/ prob } w_k \end{cases} \quad \text{and } Z \sim N(0, \sigma^2 I)$$

Since $C \perp\!\!\!\perp Z$ we have

$$\mathbb{E}[XX^T] = \underbrace{\mathbb{E}[CC^T]}_{\sum_{i=1}^{k} w_i M_i M_i^T} + \underbrace{\mathbb{E}[ZZ^T]}_{\sigma^2 I}$$

The variational characterization of singular values tells us

$$\sigma_{k+1}(M) = \min_{\dim(V) = k} \max_{u \perp V} \frac{u^T M u}{u^T u}$$

Hence $\sigma_{k+1}(M) = \sigma_{k+2}(M) = \ldots = \sigma^2$

Thus all but the top $k$ singular vectors must be $\perp \text{span}(M_1, \ldots, M_k)$. ☑

So if we estimate $M$ well enough, we can reduce to a $k$-dimensional problem

$$\frac{\text{separation}}{d^{1/4} \sigma \sqrt{\log d}} \implies \frac{\text{separation}}{k^{1/4} \sigma \sqrt{\log d}}$$

why is this $d$ and not $k$? we still
need to cluster $d^{??}$ samples