

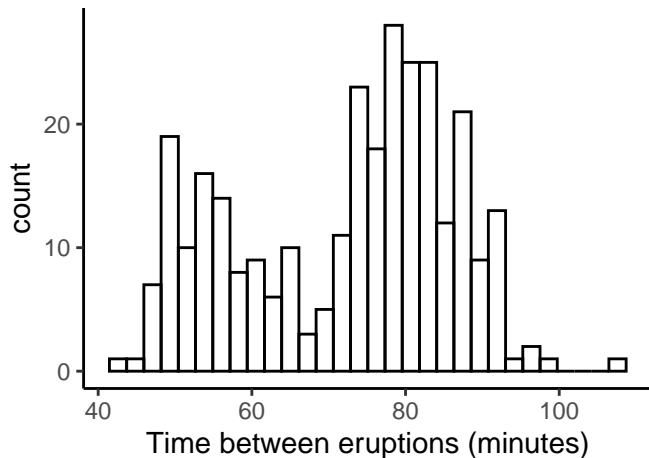
Lecture 31: Gaussian mixture models with multivariate data

Ciaran Evans

Course logistics

- ▶ HW 8 released, due Monday 4/28
- ▶ This week: wrap up Gaussian mixtures and EM algorithm
- ▶ Course evals: Wednesday during class
- ▶ Monday 4/28: wrap-up work day

Previously



Previously: Gaussian mixture model

- ▶ Observe data X_1, \dots, X_n
- ▶ Assume each observation i comes from one of k groups. Let $Z_i \in \{1, \dots, k\}$ denote the group assignment
 - ▶ The group Z is an unobserved (**latent**) variable

Model:

- ▶ $P(Z_i = j) = \lambda_j$
- ▶ $X_i | (Z_i = j) \sim N(\mu_j, \sigma_j^2)$

Posterior probabilities and parameter estimation

- ▶ **If** we know the parameters λ , μ , σ , we can calculate posterior probabilities:

$$P(Z_i = j|X_i) = \frac{\lambda_j f(X_i|Z_i = j)}{\lambda_1 f(X_i|Z_i = 1) + \dots + \lambda_k f(X_i|Z_i = k)}$$

- ▶ **If** we know the posterior probabilities, we can estimate the model parameters λ , μ , and σ :

$$\hat{\lambda}_j = \frac{1}{n} \sum_{i=1}^n P(Z_i = j|X_i)$$

$$\hat{\mu}_j = \frac{\sum_{i=1}^n X_i P(Z_i = j|X_i)}{\sum_{i=1}^n P(Z_i = j|X_i)}$$

$$\hat{\sigma}_j = \sqrt{\frac{\sum_{i=1}^n (X_i - \hat{\mu}_j)^2 P(Z_i = j|X_i)}{\sum_{i=1}^n P(Z_i = j|X_i)}}$$

Putting everything together

Model: $P(Z_i = j) = \lambda_j$ $X_i | (Z_i = j) \sim N(\mu_j, \sigma_j^2)$

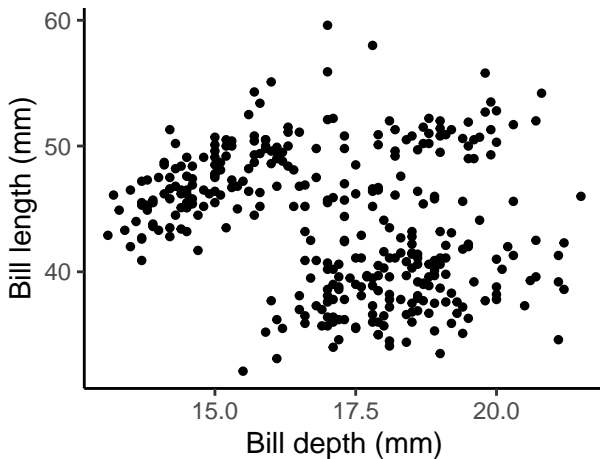
Parameters: $\lambda = (\lambda_1, \dots, \lambda_k)$, $\mu = (\mu_1, \dots, \mu_k)$, $\sigma = (\sigma_1, \dots, \sigma_k)$

Estimation:

1. Initialize parameter guesses $\lambda^{(0)}$, $\mu^{(0)}$, $\sigma^{(0)}$
2. Given current parameter estimates, compute $P^{(0)}(Z_i = j | X_i)$ for all i, j
3. Given current posterior probabilities $P^{(0)}(Z_i = j | X_i)$, update parameter estimates to $\lambda^{(1)}$, $\mu^{(1)}$, $\sigma^{(1)}$
4. Iterate: repeat steps 2–3 until convergence

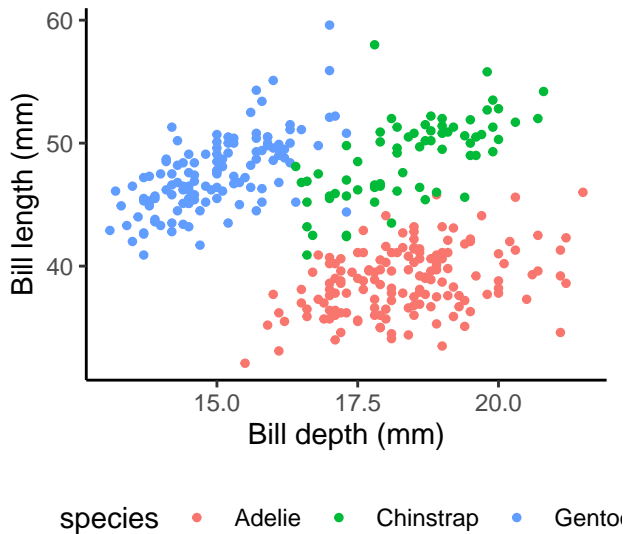
Multivariate data

Penguin data:

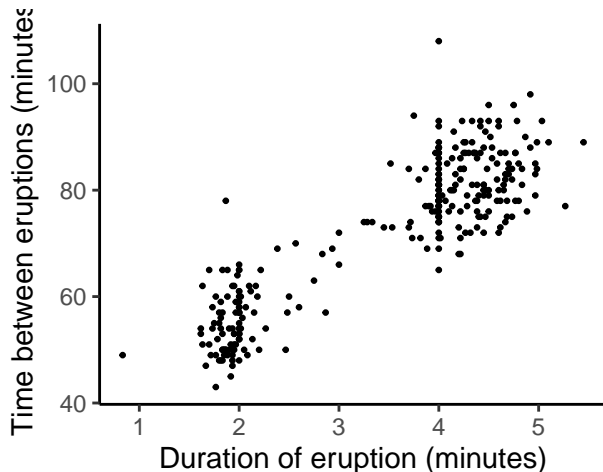


Question: What do you notice about this scatterplot?

Multivariate data



Multivariate data



Question: How should we generalize our Gaussian mixture model to multivariate data?

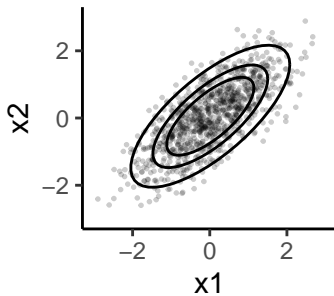
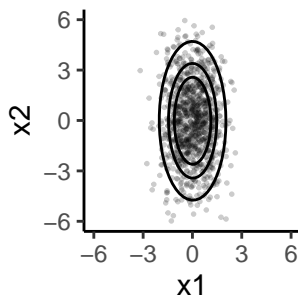
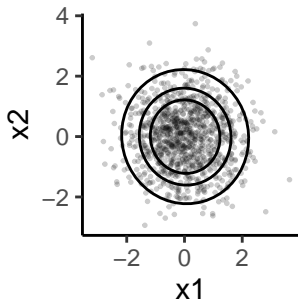
Multivariate normal distribution

Definition: Let $X = (X_1, \dots, X_k)^T$. We say that $X \sim N(\mu, \Sigma)$ if for any $\mathbf{a} \in \mathbb{R}^k$, $\mathbf{a}^T X$ follows a (univariate) normal distribution.

► $\mu =$

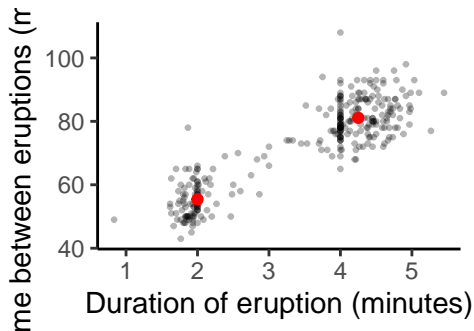
► $\Sigma =$

Multivariate normal distribution



Multivariate Gaussian mixture model

```
em_res <- mvnormalmixEM(old_faithful,  
                          lambda = c(0.5, 0.5), k=2)
```



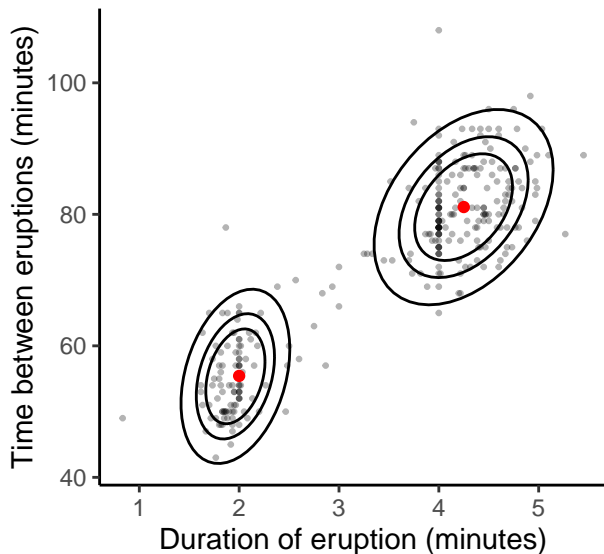
```
em_res$mu[[1]]
```

```
## [1] 1.966059 55.430118
```

```
em_res$mu[[2]]
```

```
## [1] 4.250835 81.114544
```

Multivariate Gaussian mixture model



Your turn

Implement the algorithm to fit a Gaussian mixture model:

https://sta379-s25.github.io/practice_questions/pq_31.html

- ▶ Start in class
- ▶ Welcome to work with a neighbor
- ▶ Solutions are posted on the course website