

High-dimensional Gaussians

Daniel Hsu

COMS 4772

1

Gaussian distributions

2

Gaussian (normal) distributions

- ▶ $Z \sim N(0, 1)$ means Z follows a *standard Gaussian distribution*, i.e., has probability density

$$z \mapsto \frac{1}{\sqrt{2\pi}} e^{-z^2/2}.$$

- ▶ If Z_1, Z_2, \dots, Z_d are iid $N(0, 1)$ random variables, then say $\mathbf{Z} := (Z_1, Z_2, \dots, Z_d)$ follows a *standard multivariate Gaussian distribution* on \mathbb{R}^d , i.e., $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I})$.
- ▶ Other Gaussian distributions on \mathbb{R}^d arise by applying (invertible) linear maps and translations to \mathbf{Z} :

$$\mathbf{z} \mapsto \overbrace{\mathbf{A}\mathbf{z}}^{\text{linear map}} \mapsto \underbrace{\mathbf{A}\mathbf{z} + \boldsymbol{\mu}}_{\text{translation}}.$$

- ▶ $\mathbf{X} := \mathbf{A}\mathbf{Z} + \boldsymbol{\mu} \sim N(\boldsymbol{\mu}, \mathbf{A}\mathbf{A}^\top)$ has

$$\mathbb{E}(\mathbf{X}) = \boldsymbol{\mu} \quad \text{and} \quad \text{cov}(\mathbf{X}) = \mathbf{A}\mathbf{A}^\top.$$

3

Shape of Gaussian distributions

- ▶ Let $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, $\boldsymbol{\mu} \in \mathbb{R}^d$, and $\boldsymbol{\Sigma} \succ \mathbf{0}$.
- ▶ Contours of equal density are ellipsoids around $\boldsymbol{\mu}$:

$$\{\mathbf{x} \in \mathbb{R}^d : (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) = r^2\}.$$

- ▶ Let eigenvalues of $\boldsymbol{\Sigma}$ be $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d > 0$, corresponding (orthonormal) eigenvectors be $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d$.
 - ▶ $\text{var}(\langle \mathbf{v}_i, \mathbf{X} \rangle) = \lambda_i$. (This is true even if \mathbf{X} is not Gaussian.)
 - ▶ If $Y_i := \langle \mathbf{v}_i, \mathbf{X} - \boldsymbol{\mu} \rangle$, then $Y_i \sim N(0, \lambda_i)$.
 - ▶ Y_1, Y_2, \dots, Y_d are independent;
 $\mathbf{Y} := (Y_1, Y_2, \dots, Y_d) \sim N(\mathbf{0}, \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_d))$.
- ▶ What about concentration properties?

4

Concentration of spherical Gaussians

- ▶ Spherical Gaussian: $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$.
- ▶ Pick any $\delta \in (0, 1)$. Then

$$\begin{aligned} \text{for any } \mathbf{u} \in S^{d-1}, \quad \mathbb{P}\left(\langle \mathbf{u}, \mathbf{X} - \boldsymbol{\mu} \rangle \leq \sigma \sqrt{2 \ln(1/\delta)}\right) &\geq 1 - \delta, \\ \mathbb{P}\left(\|\mathbf{X} - \boldsymbol{\mu}\|_2^2 \leq \sigma^2 d \left(1 + 2\sqrt{\frac{\ln(1/\delta)}{d}} + \frac{2 \ln(1/\delta)}{d}\right)\right) &\geq 1 - \delta, \\ \mathbb{P}\left(\|\mathbf{X} - \boldsymbol{\mu}\|_2^2 \geq \sigma^2 d \left(1 - 2\sqrt{\frac{\ln(1/\delta)}{d}}\right)\right) &\geq 1 - \delta. \end{aligned}$$

(Standard tail bounds for $\mathcal{N}(0, 1)$ and $\chi^2(d)$ distributions.)

- ▶ Behaves like spherical shell around $\boldsymbol{\mu}$ of radius $\sigma\sqrt{d}$ and thickness $O(\sigma d^{1/4})$.

5

Concentration of general Gaussians

- ▶ General Gaussian: $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.
- ▶ Concentration of $\langle \mathbf{u}, \mathbf{X} - \boldsymbol{\mu} \rangle$ for $\mathbf{u} \in S^{d-1}$ depends on \mathbf{u} :

$$\langle \mathbf{u}, \mathbf{X} - \boldsymbol{\mu} \rangle \sim \mathcal{N}(0, \mathbf{u}^\top \boldsymbol{\Sigma} \mathbf{u}).$$

- ▶ Concentration of $\|\mathbf{X} - \boldsymbol{\mu}\|_2^2$: *a mismatch of norms*.
 - ▶ $\|\boldsymbol{\Sigma}^{-1/2}(\mathbf{X} - \boldsymbol{\mu})\|_2^2 \sim \chi^2(d)$.
 - ▶ $\|\mathbf{X} - \boldsymbol{\mu}\|_2^2$ distributed as linear combination of independent $\chi^2(1)$ random variables:

$$\sum_{i=1}^d \lambda_i Z_i^2$$

where Z_1, Z_2, \dots, Z_d are iid $\mathcal{N}(0, 1)$.

- ▶ $\mathbb{E} \|\mathbf{X} - \boldsymbol{\mu}\|_2^2 = \sum_{i=1}^d \lambda_i$.
- ▶ $\|\mathbf{X} - \boldsymbol{\mu}\|_2^2$ is $(4 \sum_{i=1}^d \lambda_i^2, 4\lambda_1)$ -subexponential.

6

Eccentricity of general Gaussians

- ▶ For $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, with probability $1 - \delta$,

$$\|\mathbf{X} - \boldsymbol{\mu}\|_2^2 \in \bar{\lambda}d \left(1 \pm O\left(\sqrt{\frac{\kappa \log(1/\delta)}{d}} + \frac{\kappa \log(1/\delta)}{d} \right) \right),$$

where $\bar{\lambda} := \frac{1}{d} \sum_{i=1}^d \lambda_i$ and $\kappa := \lambda_1/\bar{\lambda}$.

- ▶ κ measure *eccentricity* of equal density ellipsoids: $1 \leq \kappa \leq d$.

7

Using multivariate Gaussians as a statistical model

- ▶ $\mathcal{P} := \{\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) : \boldsymbol{\mu} \in \mathbb{R}^d, \boldsymbol{\Sigma} \succ \mathbf{0}\}$
- ▶ Parameter estimation given data $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$
- ▶ Maximum likelihood estimators:

$$\hat{\boldsymbol{\mu}} := \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i, \quad \hat{\boldsymbol{\Sigma}} := \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^\top.$$

- ▶ Accuracy when data is iid sample from $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$:

$$\|\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}\|_2 \leq ?, \quad \|\hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}\|_? \leq ?$$

- ▶ $\hat{\boldsymbol{\mu}} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}/n)$.
- ▶ $\|\hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}\|_2 = \max_{\mathbf{u} \in S^{d-1}} \left| \mathbf{u}^\top (\hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}) \mathbf{u} \right|$.
- ▶ Note that $\mathbb{E}(\hat{\boldsymbol{\Sigma}}) \neq \boldsymbol{\Sigma}$, but almost.

8

Multiple Gaussian populations

9

Multiple populations

- ▶ Often data do not come from just a single population, but rather **several different populations**.
- ▶ **If data are “labeled” by population**, then can partition data, and (say) fit a Gaussian distribution to each part (or whatever).
- ▶ **What if data are not labeled?**

10

Simple case: multiple Gaussian populations

- ▶ Suppose data come from k populations P_1, P_2, \dots, P_k .
- ▶ Further, for extreme simplicity, suppose $P_i = N(\boldsymbol{\mu}_i, \mathbf{I})$.
- ▶ When can we separate data from P_i and P_j ($i \neq j$)?
 - ▶ Easier when means $\boldsymbol{\mu}_i$ and $\boldsymbol{\mu}_j$ are farther apart.
- ▶ **Strict separation condition:**
 - ▶ Whenever \mathbf{a} and \mathbf{b} come from same P_i , and \mathbf{c} comes from different P_j ,
$$\|\mathbf{a} - \mathbf{b}\|_2 < \|\mathbf{a} - \mathbf{c}\|_2.$$
- ▶ Under strict separation, Kruskal's minimum spanning tree (where edge weight = Euclidean distance) *connects data from same population, before connecting across populations.*
- ▶ How far apart should $\boldsymbol{\mu}_i$ and $\boldsymbol{\mu}_j$ be to have strict separation?

11

Disjoint spherical shells

- ▶ Recall: $N(\boldsymbol{\mu}_i, \mathbf{I}) \approx$ thin spherical shell around $\boldsymbol{\mu}_i$ of radius \sqrt{d} .
- ▶ If $\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|_2 \gg \sqrt{d}$, then " $N(\boldsymbol{\mu}_i, \mathbf{I}) \cap N(\boldsymbol{\mu}_j, \mathbf{I}) \approx 0$ ".
 - ▶ (This can be easily formalized.)
- ▶ But this reasoning ignores approximate orthogonality!

12

Approximate orthogonality

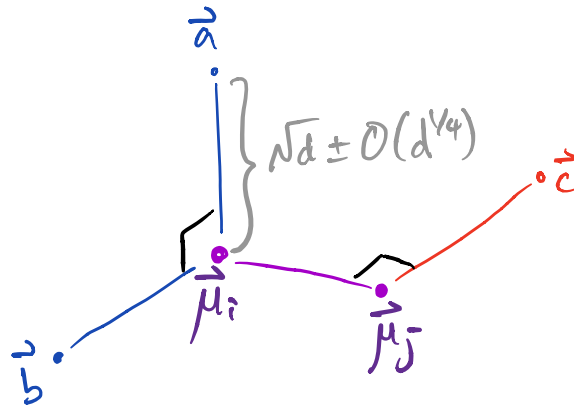


Figure 1: Distances between points from spherical Gaussian populations

13

Probabilistic analysis

- ▶ Let $\mathbf{A}, \mathbf{B} \sim N(\mu_i, \mathbf{I})$ and $\mathbf{C} \sim N(\mu_j, \mathbf{I})$ (all independent).

- ▶ Write

$$\mathbf{A} = \mu_i + \mathbf{Z}_A, \quad \mathbf{C} = \mu_j + \mathbf{Z}_C,$$

$$\mathbf{B} = \mu_i + \mathbf{Z}_B,$$

where $\mathbf{Z}_A, \mathbf{Z}_B, \mathbf{Z}_C$ are iid $N(\mathbf{0}, \mathbf{I})$.

- ▶ Then

$$\begin{aligned} & \|\mathbf{A} - \mathbf{C}\|_2^2 - \|\mathbf{A} - \mathbf{B}\|_2^2 \\ &= \|\mu_i - \mu_j\|_2^2 + 2\langle \mu_i - \mu_j, \mathbf{Z}_A - \mathbf{Z}_C \rangle + \|\mathbf{Z}_A - \mathbf{Z}_C\|_2^2 \\ & \quad - \|\mathbf{Z}_A - \mathbf{Z}_B\|_2^2. \end{aligned}$$

- ▶ With high probability, this is at least

$$\|\mu_i - \mu_j\|_2^2 - O(\|\mu_i - \mu_j\|_2) - O(\sqrt{d}),$$

which is positive when $\|\mu_i - \mu_j\|_2 \gg d^{1/4}$.

14

Probabilistic analysis (continued)

- ▶ Need previous concentration to hold for all triples in n data: union bound over $O(n^3)$ events means we need $\log(n)$ factors in separation, specifically

$$\|\mu_i - \mu_j\|_2 \geq C \left((d \log(n))^{1/4} + \log(n) \right) \quad \text{for all } i \neq j,$$

where $C > 0$ is a sufficiently large absolute constant.

15

Mixture models

- ▶ Can think of overall population as a *mixture distribution*

$$\pi_1 N(\mu_1, I) + \pi_2 N(\mu_2, I) + \cdots + \pi_k N(\mu_k, I),$$

where π_i is expected proportion from $N(\mu_i, I)$.

- ▶ Usually MLE for mixture distribution parameters $\{(\pi_i, \mu_i)\}_{i=1}^k$ is computationally intractable in general.
- ▶ But with **strict separation**:
 - ▶ First separate data by *mixture component* source.
 - ▶ Then estimate π_i and μ_i using separated data.

16

Another approach

- ▶ Project data to line spanned by some $\mathbf{u} \in S^{d-1}$.
- ▶ With “good” \mathbf{u} , projected means remain separated.
 - ▶ Use classical statistical methods to estimate projected means.
- ▶ Do this for d nearby but linearly independent \mathbf{u} ; can then back-out estimates of original means.

17

Projection pursuit

18

Exploratory data analysis (Tukey)

- ▶ Many classical data analysis methods based on finding “interesting” features of data set.
- ▶ E.g., visually inspect many one-dimensional projections of data.
- ▶ Called **projection pursuit**.
- ▶ Folklore: most projections are not interesting.

19

Two examples

- ▶ $\mathbf{X} = (X_1, X_2, \dots, X_d)$ is Rademacher (i.e., uniform on $\{\pm 1\}^d$).
- ▶ $\mathbf{u}_1 := (1/\sqrt{d}, 1/\sqrt{d}, \dots, 1/\sqrt{d})$:

$$\langle \mathbf{u}_1, \mathbf{X} \rangle = \frac{1}{\sqrt{d}} \sum_{i=1}^d X_i.$$

- ▶ By central limit theorem, this is approximately $N(0, 1)$.
- ▶ $\mathbf{u}_2 := (1, 0, \dots, 0)$:

$$\langle \mathbf{u}_2, \mathbf{X} \rangle = X_1.$$

- ▶ Very different from $N(0, 1)$.
- ▶ **“Theorem”**: Most projections are more like \mathbf{u}_1 rather than \mathbf{u}_2 .

20

Projection pursuit asymptotics (Diaconis-Freedman, 1984)

- ▶ Suppose X_1, X_2, \dots, X_d are independent random variables.
 - ▶ Assume $\mathbb{E}(X_i) = 0$, $\mathbb{E}(X_i^2) = 1$, $\mathbb{E}|X_i|^3 \leq \rho < \infty$.
- ▶ For nearly all $\mathbf{u} \in S^{d-1}$,

$$\sup_{t \in \mathbb{R}} \left| \mathbb{P}(\langle \mathbf{u}, \mathbf{X} \rangle \leq t) - \Phi(t) \right| \leq \tilde{O}\left(\frac{\rho}{\sqrt{d}}\right),$$

where Φ is $N(0, 1)$ CDF.

21

Application to mixture models

- ▶ Suppose $\mathbf{X} \sim \pi_1 P_1 + \pi_2 P_2 + \dots + \pi_k P_k$, where each P_i is a *product distribution*.
- ▶ \mathbf{X} generally does not have independent coordinates.
- ▶ But for most $\mathbf{u} \in S^{d-1}$, $\langle \mathbf{u}, \mathbf{X} \rangle$ is close to

$$\pi_1 N_1 + \pi_2 N_2 + \dots + \pi_k N_k$$

for some univariate normal distributions N_1, N_2, \dots, N_k .

22