# CSCI B609:
# "Foundations of Data Science"

# Lecture 5: Dimension Reduction, Separating and Fitting Gaussians

Slides at http://grigory.us/data-science-class.html

# Grigory Yaroslavtsev

## http://grigory.us

# Gaussian Annulus Theorem

- Gaussian in $\boldsymbol{d}$ dimensions ($N_{\boldsymbol{d}}(0^{\boldsymbol{d}}, 1)$):

$$\Pr[\boldsymbol{x} = (z_1, \ldots, z_{\boldsymbol{d}})] = (2\pi)^{-\frac{\boldsymbol{d}}{2}} e^{-\frac{z_1^2 + z_2^2 + \cdots + z_{\boldsymbol{d}}^2}{2}}$$

Nearly all mass in annulus of radius $\sqrt{\boldsymbol{d}}$ and width $O(1)$:

- **Thm.** For any $\boldsymbol{\beta} \leq \sqrt{\boldsymbol{d}}$ all but $3e^{-\boldsymbol{c}\boldsymbol{\beta}^2}$ probability mass satisfies $\sqrt{\boldsymbol{d}} - \boldsymbol{\beta} \leq \lVert\boldsymbol{x}\rVert_2 \leq \sqrt{\boldsymbol{d}} + \boldsymbol{\beta}$ for constant $\boldsymbol{c}$

# Nearest Neighbors and Random Projections

- Given a database $A$ of $n$ points in $\mathbb{R}^d$
  - Preprocess $A$ into a small data structure $D$
  - Should answer following queries fast:

Given $q \in \mathbb{R}^d$ find closest $x \in A$: $argmin_{x \in A} \left|\left| q - x \right|\right|_2$

- Project each $x \in A$ onto $f(x)$, where $f: \mathbb{R}^d \to \mathbb{R}^k$

- Pick $k$ vectors $u_1, \ldots, u_k$ i.i.d: $u_i \sim N_d\left(0^d, 1\right)$
$$f(v) = (\langle u_1, v \rangle, \ldots, \langle u_k, v \rangle)$$

- Will show that w.h.p. $\left|\left| f(v) \right|\right|_2 \approx \sqrt{k} \left|\left| v \right|\right|_2$

**Return:** $argmin_{x \in A} \left|\left| f(q) - f(x) \right|\right|_2 = argmin_{x \in A} \left|\left| f(q - x) \right|\right|_2 \approx$
$\sqrt{k} \ argmin_{x \in A} \left|\left| q - x \right|\right|_2$

# Random Projection Theorem

- Pick $\boldsymbol{k}$ vectors $\boldsymbol{u}_1, \ldots, \boldsymbol{u}_k$ i.i.d: $\boldsymbol{u}_i \sim N_{\boldsymbol{d}}(0^{\boldsymbol{d}}, 1)$
$$f(\boldsymbol{v}) = (\langle \boldsymbol{u}_1, \boldsymbol{v} \rangle, \ldots, \langle \boldsymbol{u}_k, \boldsymbol{v} \rangle)$$

- Will show that w.h.p. $\left\lVert f(\boldsymbol{v}) \right\rVert_2 \approx \sqrt{\boldsymbol{k}} \left\lVert \boldsymbol{v} \right\rVert_2$

> **Thm.** Fix $\boldsymbol{v} \in \mathbb{R}^{\boldsymbol{d}}$ then $\exists c > 0$: for $\boldsymbol{\epsilon} \in (0,1)$:
> $$\Pr_{\boldsymbol{u}_i \sim N_{\boldsymbol{d}}(0^{\boldsymbol{d}}, 1)} \left[ \left\lvert \left\lVert f(\boldsymbol{v}) \right\rVert_2 - \sqrt{\boldsymbol{k}} \left\lVert \boldsymbol{v} \right\rVert_2 \right\rvert \geq \boldsymbol{\epsilon} \sqrt{\boldsymbol{k}} \left\lVert \boldsymbol{v} \right\rVert_2 \right] \leq 3\, e^{-c\boldsymbol{k}\boldsymbol{\epsilon}^2}$$

- Scaling: $\left\lVert \boldsymbol{v} \right\rVert_2 = 1$

- **Key fact**: $\langle \boldsymbol{u}_i, \boldsymbol{v} \rangle = \sum_{j=1}^{\boldsymbol{d}} \boldsymbol{u}_{ij} \boldsymbol{v}_j \sim N(0, \left\lVert \boldsymbol{v} \right\rVert_2^2) = N(0,1)$

- Apply "Gaussian Annulus Theorem" with $\boldsymbol{k} = \boldsymbol{d}$

# Nearest Neighbors and Random Projections

**Thm.** Fix $\boldsymbol{v} \in \mathbb{R}^{\boldsymbol{d}}$ then $\exists c > 0$: for $\boldsymbol{\epsilon} \in (0,1)$:

$$\Pr_{\boldsymbol{u}_i \sim N_{\boldsymbol{d}}(0^{\boldsymbol{d}},1)}\left[\left|\left|\left|f(\boldsymbol{v})\right|\right|_2 - \sqrt{\boldsymbol{k}}||\boldsymbol{v}||_2\right| \geq \boldsymbol{\epsilon}\sqrt{\boldsymbol{k}}||\boldsymbol{v}||_2\right] \leq 3\,e^{-c\boldsymbol{k}\boldsymbol{\epsilon}^2}$$

**Return:** $argmin_{x \in A}\left|\left|f(\boldsymbol{q}) - f(\boldsymbol{x})\right|\right|_2 \approx \sqrt{\boldsymbol{k}}\, argmin_{x \in A}||\boldsymbol{q} - \boldsymbol{x}||_2$

- Fix and let $\boldsymbol{v} = \boldsymbol{q} - \boldsymbol{x}_i$ for $\boldsymbol{x}_i \in A$ and let $\boldsymbol{k} = O\left(\frac{\gamma \log n}{\boldsymbol{\epsilon}^2}\right)$

$(1 \pm \boldsymbol{\epsilon})\sqrt{\boldsymbol{k}}||\boldsymbol{q} - \boldsymbol{x}_i||_2 \approx \left|\left|f(\boldsymbol{q}) - f(\boldsymbol{x})\right|\right|_2$ (prob. $1 - n^{-\gamma}$)

- Union bound:

For fixed $\boldsymbol{q}$ distances to $\boldsymbol{A}$ preserved with prob. $1 - n^{-\gamma+1}$

# Separating Gaussians

- One-dimensional mixture of Gaussians:
$$p(x) = w_1 p_1(x) + w_2 p_2(x)$$
- E.g. modeling heights of men/women
- **Parameter estimation problem**:
  - Given samples from a mixture of Gaussians
  - **Q:** Estimate means and (co)-variances
- **Sample origin problem**:
  - Given samples from **well-separated** Gaussians
  - **Q**: Did they come from the same Gaussian?

# Separating Gaussians

- Gaussian in $d$ dimensions ($N_d(0^d, 1)$):

$$\Pr[x = (z_1, \ldots, z_d)] = (2\pi)^{-\frac{d}{2}} e^{-\frac{z_1^2 + z_2^2 + \cdots + z_d^2}{2}}$$

Nearly all mass in annulus of radius $\sqrt{d}$ and width $O(1)$:

- Almost all mass in a slab $\{x| -c \leq x_1 \leq c\}$ for $c = O(1)$
- Pick $x \sim$ Gaussian and rotate coordinates to make it $x_1$
- Pick $y \sim$ Gaussian, w.h.p. projection of $y$ on $x$ is $\in [-c, c]$

$$||x - y||_2 \approx \sqrt{||x||_2^2 + ||y||_2^2}$$

# Separating Gaussians

In coordinates:

- $x = (\sqrt{d} \pm O(1), 0, 0, \ldots, 0)$

- $y = (\pm O(1), \sqrt{d} \pm O(1), 0, \ldots, 0)$

- W.h.p: $\left\|x - y\right\|_2^2 = 2d \pm O(\sqrt{d})$

# Separating Gaussians

- Two spherical unit variance Gaussians centered at $p, q$
- $\delta = \left\| p - q \right\|_2$
- $(x \sim N(p, 1), y \sim N(q, 1))$
- $x = \left( \sqrt{d} \pm O(1), 0, 0, 0 \ldots, 0 \right)$
- $y = (\pm O(1), \delta \pm O(1), \sqrt{d} \pm O(1), 0, \ldots, 0)$
- $\left\| x - y \right\|_2^2 = \delta^2 + 2d \pm O(\sqrt{d})$

# Separating Gaussians

- Same Gaussian:

$$\left|\left|x-y\right|\right|_2^2 = 2d \pm O(\sqrt{d})$$

- Different Gaussians:

$$\left|\left|x-y\right|\right|_2^2 = \delta^2 + 2d \pm O(\sqrt{d})$$

- Separation requires:

$$2d \pm O(\sqrt{d}) < \delta^2 + 2d \pm O(\sqrt{d})$$

$$O(\sqrt{d}) < \delta^2$$

$$\omega(d^{1/4}) < \delta$$

# Fitting Spherical Gaussian to Data

- Given samples $x_1, x_2, \ldots, x_n$
- **Q:** What are parameters of best fit $N(\boldsymbol{\mu}, \sigma)$?

$$\Pr[x_i = (z_1, \ldots, z_d)]$$

$$= (2\pi)^{-\frac{d}{2}} e^{-\frac{(\mu_1 - z_1)^2 + (\mu_2 - z_2)^2 + \cdots + (\mu_d - z_d)^2}{2}}$$

$$= (2\pi)^{-\frac{d}{2}} e^{-\frac{||\boldsymbol{\mu} - z||_2^2}{2}}$$

$$\Pr[x_1 = z_1, x_2 = z_2, \ldots, x_n = z_n]$$

$$= (2\pi)^{-\frac{dn}{2}} e^{-\frac{||\boldsymbol{\mu} - z_1||_2^2 + ||\boldsymbol{\mu} - z_2||_2^2 + \cdots + ||\boldsymbol{\mu} - z_d||_2^2}{2\sigma^2}}$$

# Maximum Likelihood Estimator

- PDF: $(2\pi)^{-\frac{dn}{2}} e^{-\frac{\|\mu-z_1\|_2^2 + \|\mu-z_2\|_2^2 + \cdots + \|\mu-z_d\|_2^2}{2\sigma^2}}$

- MLE for $\mu$ is $\mu = \frac{1}{n}(x_1 + x_2 + \cdots + x_n)$

- Take gradient w.r.t $\mu$ and make it $= 0$

- $\nabla_\mu \|\mu - x\|_2^2 = 2(\mu - x)$

$2(\mu - x_1) + 2(\mu - x_2) + \cdots + 2(\mu - x_n) = 0$

# MLE for Variance

- $a = ||\boldsymbol{\mu} - \boldsymbol{x_1}||_2^2 + ||\boldsymbol{\mu} - \boldsymbol{x_2}||_2^2 + \cdots + ||\boldsymbol{\mu} - \boldsymbol{x_d}||_2^2$
- $\nu = 1/2\sigma^2$

- PDF: $\dfrac{e^{-a\nu}}{\left[\int_{\boldsymbol{x}\in R^d} e^{-\nu||\boldsymbol{x}||_2^2} d\boldsymbol{x}\right]^n}$

- Log(PDF): $-a\nu - n\ln\left[\int_{\boldsymbol{x}\in R^d} e^{-\nu||\boldsymbol{x}||_2^2} d\boldsymbol{x}\right]$

- Differentiate w.r.t. $\nu$ and set derivative $= 0$

# MLE for Variance

- Log(PDF): $-a\nu - n\ln\left[\int_{x \in R^d} e^{-\nu\|x\|_2^2}\, dx\right]$

- $\frac{d}{d\nu}$ Log(PDF):

$$-a + n\frac{\int_{x \in R^d}\|x\|_2^2\, e^{-\nu\|x\|_2^2}\, dx}{\int_{x \in R^d} e^{-\nu\|x\|_2^2}\, dx}$$

- $y = \|\nu x\|_2^2$:

$$-a + \frac{n}{\nu}\frac{\int_{x \in R^d} y^2\, e^{-y^2}\, dx}{\int_{x \in R^d} e^{-y^2}\, dx} = -a + \frac{n}{\nu} \times \frac{d}{2} = 0$$

- $\nu = \frac{1}{2\sigma^2} \Rightarrow \text{MLE}(\sigma) = \sqrt{\frac{a}{nd}}$ = sample standard deviation