# Lecture 2: Learning a High-Dimensional Gaussian

- Last time: testing high dimensional distributions
  (= large domain size) is <u>hard</u>
  (requires many samples).

- Today: a strong assumption under which efficient
  learning is possible.

## Gaussian Distributions

- Reminder: Gaussian in one dimension :



$$Pr(x) \propto e^{-(x-\mu)^2/2\sigma^2}$$

- CLT: add together many indep. r.v.'s $\rightarrow$ Gaussian

- Also holds in high dimensions!

- So if getting samples from a high-dimensional population
  where high-dimensional features act like sums of indep
  r.v's, Gaussian assumption can be reasonable.

- Multivariate Gaussian: any affine transformation of
  $$Pr(x) \propto e^{-\|x\|^2/2}$$

- Notation: transform by $x \rightarrow Ax + \mu$, distn is called

$$N(\mu, A^2) \quad (\text{traditionally}: N(\mu, \Sigma).)$$

- Fact: $\underset{x \sim N(\mu, \Sigma)}{\mathbb{E}} x = \mu, \quad \underset{x \sim N(\mu, \Sigma)}{\mathbb{E}} (x-\mu)(x-\mu)^T = \Sigma$

  "Gaussian is determined by its 1st and 2nd moments"

---

Task: given $x_1 \ldots x_n \sim N(\mu, \Sigma)$ for some unknown $\mu \in \mathbb{R}^d$, $\Sigma \in \mathbb{R}^{d \times d}$, find some distn $D$ s.t. $TV(D, N(\mu, \Sigma)) \leq \varepsilon$.

> How big does $n = n(d, \varepsilon)$ need to be?

Theorem: $n = O\left(\frac{d^2}{\varepsilon^2}\right)$ suffices.

(Compare with $n = \Omega(\sqrt{2^d})$ for uniformity testing on $\{0,1\}^d$, only easier than learning general distn on $\{0,1\}^d$.)

(And, can do in polynomial time.)

- Let $\hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} x_i, \quad \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^{n} (x_i - \hat{\mu})(x_i - \hat{\mu})^T$.

- Lemma (directly implies Theorem above):
$$\underset{x_1 \ldots x_n \sim N(\mu, \Sigma)}{\mathbb{E}} TV\left(N(\hat{\mu}, \hat{\Sigma}), N(\mu, \Sigma)\right) \leq O\left(\frac{d}{\sqrt{n}}\right).$$

Proof outline:

① $TV\left(N(\hat{\mu},\hat{\Sigma}), N(\mu,\Sigma)\right) \leq$

$\quad TV(N(\hat{\mu},\Sigma), N(\mu,\Sigma)) + TV(N(\hat{\mu},\hat{\Sigma}), N(\hat{\mu},\Sigma))$

$= \underbrace{TV\left(N(\Sigma^{-1/2}\hat{\mu}, I), N(\Sigma^{-1/2}\mu, I)\right)}_{Ⓐ} + \underbrace{TV(N(0,\hat{\Sigma}), N(0,\Sigma))}_{Ⓑ}$

$\leq O\left(\|\Sigma^{-1/2}(\hat{\mu}-\mu)\|\right) \qquad\qquad \leq O\left(\|I - \Sigma^{-1/2}\hat{\Sigma}\Sigma^{-1/2}\|_F\right)$

② $\mathbb{E}\|\Sigma^{-1/2}(\hat{\mu}-\mu)\|^2 \leq O\left(\frac{d}{n}\right)$

③ $\mathbb{E}\|I - \Sigma^{-1/2}\hat{\Sigma}\Sigma^{-1/2}\|_F \leq O\left(\frac{d}{\sqrt{n}}\right)$

---

Reminder: $\|M\|_F^2 = \text{Tr } MM^T = \text{Tr } M^T M$

$\qquad\qquad = \sum_{i,j \leq d} M_{ij}^2 = \sum \lambda_i(M)^2$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \uparrow$ Singular values of $M$

"Frobenius norm"

Reminder 2: $\|v\|^2 = \text{Tr } vv^T = \sum_i v_i^2$

Ⓐ Goal: $TV\left(N(\mu, I), N(\tau, I)\right) \le O\left(\|\mu - \tau\|\right)$

Enough: $KL\left(N(\mu, I) \| N(\tau, I)\right) \le O\left(\|\mu - \tau\|^2\right)$

$$\overset{\shortparallel}{=} \underset{x \sim N(\mu, I)}{\mathbb{E}} \log \frac{e^{-\|x - \mu\|^2/2}}{e^{-\|x - \tau\|^2/2}}$$

$$\frac{1}{2}\left[\underset{x \sim N(\mu, I)}{\mathbb{E}} -\|x - \mu\|^2 + \|x - \tau\|^2\right]$$

$$= \frac{1}{2} \mathbb{E}\left[-\|x - \mu\|^2 + \|x - \mu + \mu - \tau\|^2\right]$$

$$= \frac{1}{2} \mathbb{E}\left[-\|x - \mu\|^2 + \|x - \mu\|^2 + 2\langle x - \mu, \mu - \tau\rangle + \|\mu - \tau\|^2\right]$$

$$= \frac{1}{2}\|\mu - \tau\|^2 \qquad\qquad \square$$

(B) Enough: $KL\left(N(0,\Sigma) \| N(0,\Gamma)\right) \le O\left(\|I - \Gamma^{-1/2}\Sigma\Gamma^{-1/2}\|_F^2\right)$

Recall that $P_{N(0,\Sigma)}(x) = \left(\frac{1}{2\pi}\right)^{d/2} \frac{1}{\sqrt{\det\Sigma}} e^{-\|\Sigma^{-1/2}x\|^2/2}$

$$KL(\cdots) = \mathbb{E}_{x\sim N(0,\Sigma)} \log \frac{\frac{1}{\sqrt{\det\Sigma}} e^{-\|\Sigma^{-1/2}x\|^2/2}}{\frac{1}{\sqrt{\det\Gamma}} e^{-\|\Gamma^{-1/2}x\|^2/2}}$$

$$= \frac{1}{2}\left[\log\frac{\det\Sigma}{\det\Gamma} + \mathbb{E}_{x\sim N(0,\Sigma)} \|\Gamma^{-1/2}x\|^2 - \|\Sigma^{-1/2}x\|^2\right]$$

$$= \frac{1}{2}\left\{-\log\det\left(\Gamma^{-1/2}\Sigma\Gamma^{-1/2}\right) + \mathbb{E}\,Tr\left(\Gamma^{-1/2}xx^T\Gamma^{-1/2} - \Sigma^{-1/2}xx^T\Sigma^{-1/2}\right)\right\}$$

$$= -\frac{1}{2}\log\det\left(\Gamma^{-1/2}\Sigma\Gamma^{-1/2}\right) + \frac{1}{2}Tr\left(\Gamma^{-1/2}\Sigma\Gamma^{-1/2} - I\right)$$

let $\lambda_1 \ldots \lambda_d$ be the eigs. of $\Gamma^{-1/2}\Sigma\Gamma^{-1/2} - I$ Then,

$$= -\frac{1}{2}\sum \log(1+\lambda_i) + \frac{1}{2}\sum\lambda_i$$

$$= O\left(\sum\lambda_i^2\right) \quad \text{if} \quad |\lambda_i| < 1 \quad \forall i \; (\text{Convergence of Taylor}).$$

ok to assume, since o.w. $\|I - \Gamma^{-1/2}\Sigma\Gamma^{-1/2}\|_F \ge 1 \ge TV(\cdots)$.

□

② $E\| \Sigma^{-1/2}(\hat{\mu} - \mu)\|^2 = E\| \frac{1}{n}\sum_{i=1}^{n} \Sigma^{-1/2} x_i - \mu\|^2$

$x_i = \Sigma(z_i + \mu)$ for $z_i \sim N(0, I)$, so

$= E\underset{z_1 \dots z_n}{\|} \frac{1}{n}\sum_{i=1}^{n} z_i\|^2 = \frac{1}{n^2} E\sum_{i,j} z_i^T z_j = \frac{1}{n^2}\sum_i E\|z_i\|^2$

$= \frac{d}{n}$

$\square$

$$\text{(3)} \quad I - \Sigma^{-1/2} \hat{\Sigma} \Sigma^{-1/2} = I - \frac{1}{n} \sum_{i=1}^{n} \Sigma^{-1/2} (X_i - \hat{\mu})(X_i - \hat{\mu})^T \Sigma^{-1/2}$$

$$= I - \frac{1}{n} \sum_{i=1}^{n} \Sigma^{-1/2} (X_i - \mu + \mu - \hat{\mu})(X_i - \mu + \mu - \hat{\mu})^T \Sigma^{-1/2}$$

$$= I - \frac{1}{n} \sum_{i=1}^{n} \Sigma^{-1/2} (X_i - \mu)(X_i - \mu)^T - \Sigma^{-1/2} (\mu - \hat{\mu})(\mu - \hat{\mu})^T \Sigma^{-1/2}$$

$$\underbrace{- \frac{1}{n} \sum_{i=1}^{n} \Sigma^{-1/2} (X_i - \mu)(\mu - \hat{\mu})^T \Sigma^{-1/2} - \frac{1}{n} \sum_{i=1}^{n} \Sigma^{-1/2} (\mu - \hat{\mu})(X_i - \mu) \Sigma^{-1/2}}$$

$$= \Sigma^{-1/2} (\hat{\mu} - \mu)(\mu - \hat{\mu})^T \Sigma^{-1/2}$$

By triangle inequality,

$$\mathbb{E} \| I - \Sigma^{-1/2} \hat{\Sigma} \Sigma^{-1/2} \|_F \leq \mathbb{E} \| I - \Sigma^{-1/2} \left( \frac{1}{n} \sum_{i=1}^{n} (X_i - \mu)(X_i - \mu)^T \right) \Sigma^{-1/2} \|_F$$

$$+ \mathbb{E} \| \Sigma^{-1/2} (\hat{\mu} - \mu) \|_2^2$$

$$\leq \mathbb{E} \| I - \Sigma^{-1/2} \left( \frac{1}{n} \sum_{i=1}^{n} (X_i - \mu)(X_i - \mu)^T \right) \Sigma^{-1/2} \|_F + O\left(\frac{d}{n}\right) \text{ by } \text{(2)}.$$

$\Sigma^{-1/2}(X_i - \mu)$ is distributed as $Z \sim N(0, I)$, so get

$$\leq \mathbb{E}_{z_1 \dots z_n} \| I - \frac{1}{n} \sum z_i z_i^T \|_F + O\left(\frac{d}{n}\right).$$

$$\leq \left( \mathbb{E}_{z_1 \dots z_n} \| I - \frac{1}{n} \sum z_i z_i^T \|_F^2 \right)^{1/2} + O\left(\frac{d}{n}\right)$$

$$= O\left(\frac{d}{n}\right) + \left( \frac{1}{n^2} \sum_{i,j \leq n} \mathbb{E} \left[ \text{Tr} \left( z_i z_i^T - \mathbb{E} z_i z_i^T \right) \cdot (z_j z_j^T - \mathbb{E} z_j z_j^T) \right] \right)^{1/2}$$

$$= O\left(\frac{d}{n}\right) + \left(\frac{1}{n} \mathbb{E} \, \text{Tr} \left(zz^\top - \mathbb{E}zz^\top\right)^2\right)^{1/2}$$

$$= O\left(\frac{d}{n}\right) + \left(\frac{1}{n} \mathbb{E} \sum_{i,j \leq d} \left(z(i)z(j) - \mathbb{E}\, z(i)z(j)\right)^2\right)^{1/2}$$

$$= O\left(\frac{d}{\sqrt{n}}\right) \qquad\qquad \square$$

What if we only cared about learning the "shape", not learning in TV dist?

Task: Given $X_1 \ldots X_n \sim N(0, \Sigma)$, find $\hat{\Sigma}$ s.t.

$$\| \Sigma^{-1/2} \hat{\Sigma} \Sigma^{-1/2} - I \|_{op} \leq \varepsilon$$

Reminder: $\| M \|_{op} = \max_v \dfrac{v^T M v}{\| v \|^2}$. So,

$$\forall v, \quad v^T \hat{\Sigma} v = (1 \pm \varepsilon) v^T \Sigma v$$

- Simultaneously estimate the variance in every direction.

- good for e.g. PCA.

Theorem: $n = O\left(\dfrac{d}{\varepsilon^2}\right)$ samples suffice — compare to $\dfrac{d^2}{\varepsilon^2}$ for T.V.

# Interlude: Matrix Concentration

- $n$ copies of a $\mathbb{R}$-valued random variable $X$ — how close is $\frac{1}{n} \sum X_i$ to $\mathbb{E} X$?

$$\left[ \mathbb{E} \left( \frac{1}{n} \sum X_i - \mathbb{E} X \right)^2 \right]^{1/2} = \sqrt{\frac{\text{Var}(x)}{n}}$$

If $X$ is a little bit "nice" (bdd, subGaussian, etc.), get concentration — $\frac{1}{n} \sum X_i$ acts like Gaussian

$$\Pr \left( \left| \frac{1}{n} \sum X_i - \mathbb{E} X \right| > t \right) \lesssim e^{-t^2 n / \text{Var}(x)}$$

- $n$ copies of a $\mathbb{R}^d$-valued random vector $X$:

$$\left[\mathbb{E}\left\| \frac{1}{n}\sum X_i - \mathbb{E}X \right\|^2\right]^{1/2} = \sqrt{\frac{\mathbb{E}\|X-\mathbb{E}X\|^2}{n}} = \sqrt{\frac{\text{Tr Cov}(X)}{n}}$$

- $n$ copies of a $\mathbb{R}^{d\times d}$ random matrix $X$ (symmetric)?

$$\mathbb{E}\left\| \frac{1}{n}\sum X_i - \mathbb{E}X \right\|_{op} \leq \ {\color{red}??}$$

- What should we expect? Suppose $X = \begin{pmatrix} X^{(1)} & & 0 \\ & \ddots & \\ 0 & & X^{(d)} \end{pmatrix}$

is diagonal.

$$\frac{1}{n}\sum X_i = \begin{pmatrix} \frac{1}{n}\sum X_i^{(1)} & & & 0 \\ & \ddots & & \\ & & \ddots & \\ 0 & & & \frac{1}{n}\sum X_i^{(d)} \end{pmatrix} \quad \text{has}$$

Singular vals $\quad \frac{1}{n}\sum X_i^{(j)}$.

For each $j \leq d$, expect

$$\Pr\left(\left| \frac{1}{n}\sum X_i^{(j)} - \mathbb{E}X^{(j)} \right| > t \right) \lesssim e^{-t^2 n / \text{Var}(X^{(j)})} \lesssim \frac{1}{100d}$$

if $\quad t = \Theta\left( \sqrt{\frac{\text{Var}(X^{(j)})\cdot \log d}{n}} \right)$

By union bound, wp $\geq 0.99$,

$$\left\| \frac{1}{n}\sum X_i - \mathbb{E}X \right\| \leq O\left( \max_j \sqrt{\text{Var}[X^{(j)}]} \cdot \sqrt{\frac{\log d}{n}} \right)$$

Is something like this true in general, ie even for non-diagonal $X$ ?

Matrix Bernstein Inequality : Let $X$ be $d \times d$ random matrix, $\mathbb{E}X = 0$, and $\|X\| \leq R$ w.p. 1, $X_1 \dots X_n$ indep. copies.

$$\mathbb{E} \left\| \frac{1}{n}\sum X_i \right\|_{op} \leq O\left( \frac{\|\mathbb{E}XX^T\| + \|\mathbb{E}X^TX\|}{\sqrt{n}} \cdot \sqrt{\log d} \right.$$
$$\left. + \frac{R \log d}{n} \right)$$

Application to prove thm, up to $\log d$ factors:

WLOG, $\Sigma = I$, and goal is to show

$$\mathbb{E} \left\| \frac{1}{n} \sum x_i x_i^T - I \right\|_{op} \leq O\left(\sqrt{\frac{d \log d}{n}}\right).$$

$\mathbb{E}[x_i x_i^T - I] = 0$, but don't have $\|x_i x_i^T - I\|_{op} \leq R$ wp 1.

But let's pretend... $\|x_i x_i^T - I\|_{op} \leq O\left(\frac{d}{n}\right)$ wp 1 (almost true).

Then, just need to calculate $\left\| \mathbb{E}(xx^T - I)(xx^T - I) \right\|$

$$\leq \left\| \mathbb{E} \|x\|^2 xx^T \right\| + O(1) \leq O(d)$$

which proves the theorem $\quad\quad\quad\quad$ □