

한국통계학회 공개강연

# 평균 낼 수 없는 것들의 통계

정성규

서울대학교 통계학과

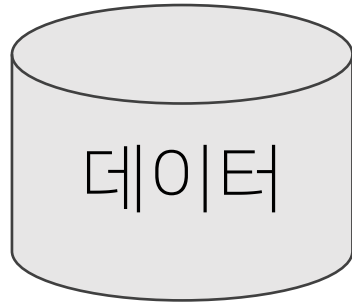
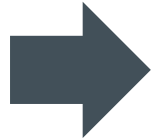
Part 1,

통계와 평균

## 통계학

### 통계

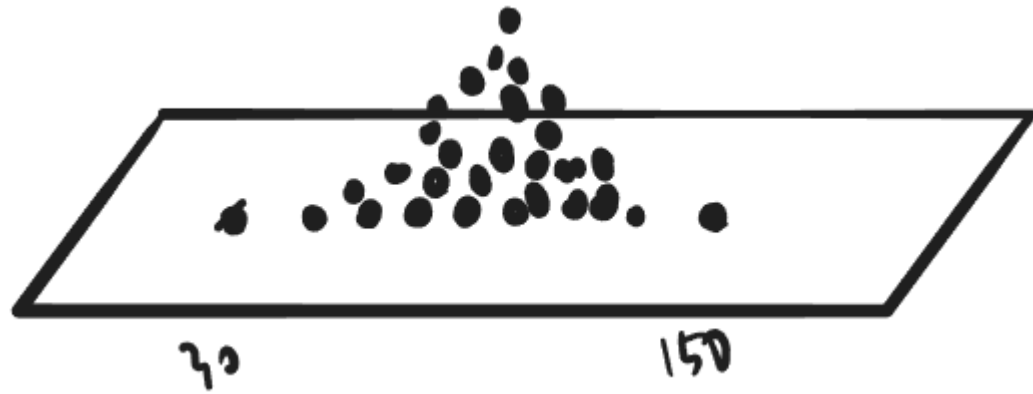
사회 현상  
자연 현상



요약된 정보

요약된 "정보"가 얼마나 정확한가?

# 90



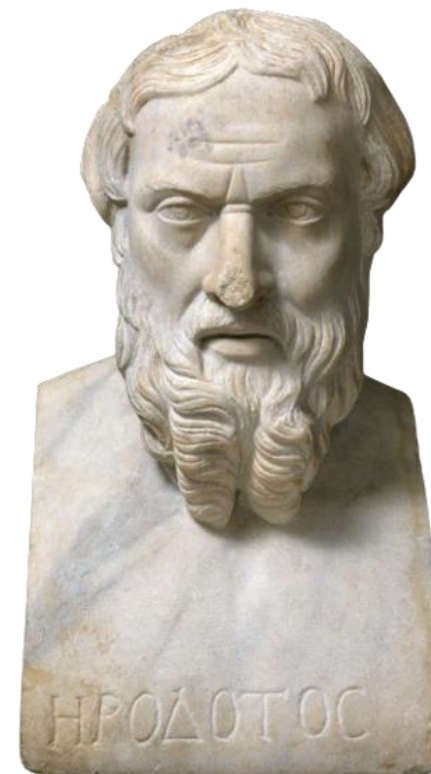
평균 70

≈ 기준, 대푯값

페르미 추정 문제:

**Q:** 이집트 왕조는 **341**세대가 있다. 이 왕조는 몇 년이나 지속되었을까?

**A:** 세 세대당 “평균” **100**년 지속된다면, **11340**년이다.\*



헤로도토스 (BC 484 – BC 425)

\* Rubin, E. (1968). STATISTICAL WORLD OF HERODOTUS. *American Statistician*, 22(1), 31.

# 평균과 미지수의 추정

미지의 참값 추정 문제:

빛의 속도 측정 실험: 29.6만km/s, 30만km/s, 30.1만km/s

빛의 속도 추정값은? 중앙값 30만km/s?

평균 29.9만km/s?

가우스 “관측오차의 제곱합이 가장 작은 추정값”

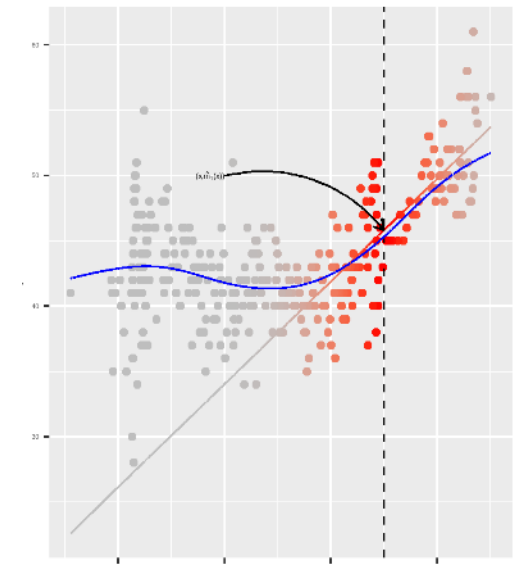
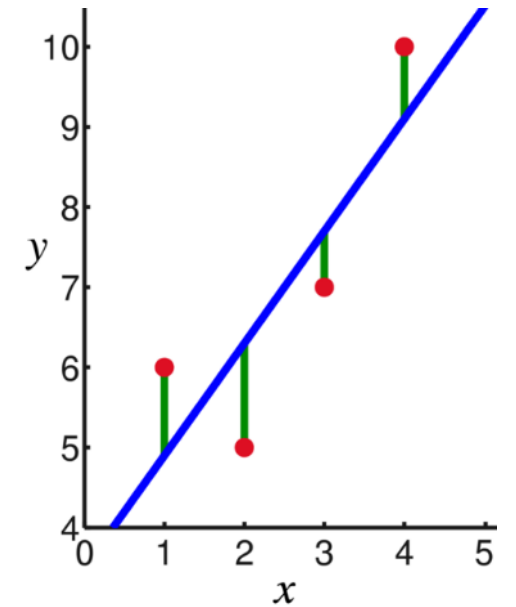
$$\text{평균 } 29.9 = \frac{29.6+30+30.1}{3}$$

$$= \operatorname{argmin}_{\mu} (29.6 - \mu)^2 + (30 - \mu)^2 + (30.1 - \mu)^2$$



카를 프리드리히 가우스 (1777-1855)

- 가우스의 최소제곱추정법: 평균 = 최소제곱법의 해.
- 선형회귀분석: 최소제곱추정된 계수 =  $(x^T x)^{-1} x^T y$   
= (가중)평균
- 비모수회귀분석 (국소다항회귀, 스플라인 평활):  
최소제곱추정된 회귀곡선  $\hat{m}(x) = \sum_{i=1}^n w_i(x) y_i$   
= (가중)평균



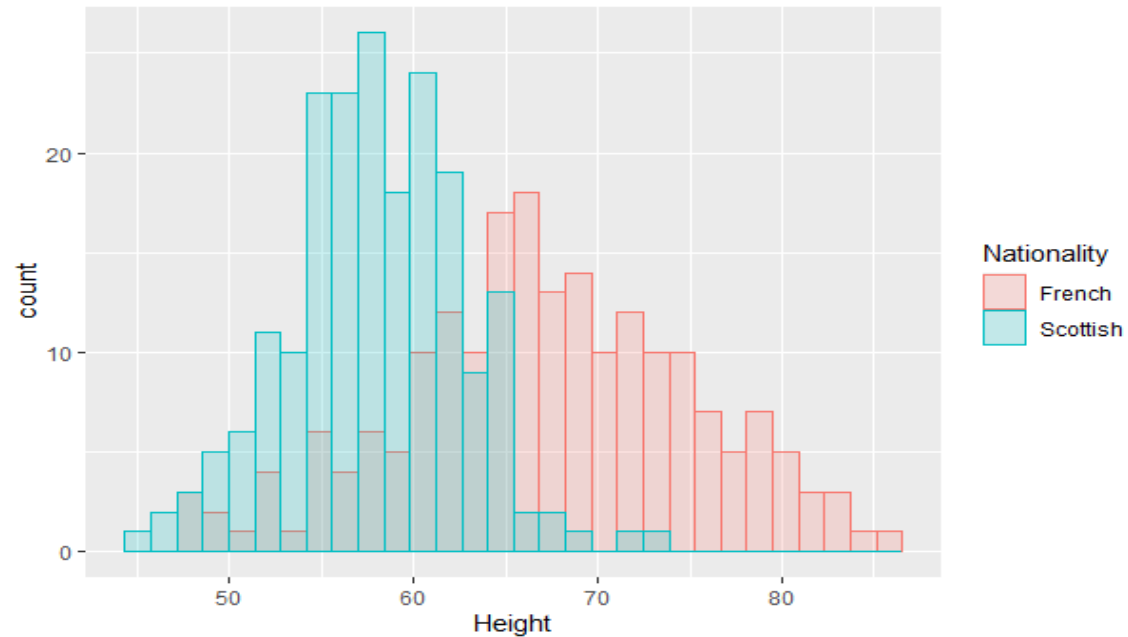
# 평균 인간

“프랑스군 신병”

Vs

“스코틀랜드군 신병”

인간과 그 재능 발달에 대한 연구 (1842)



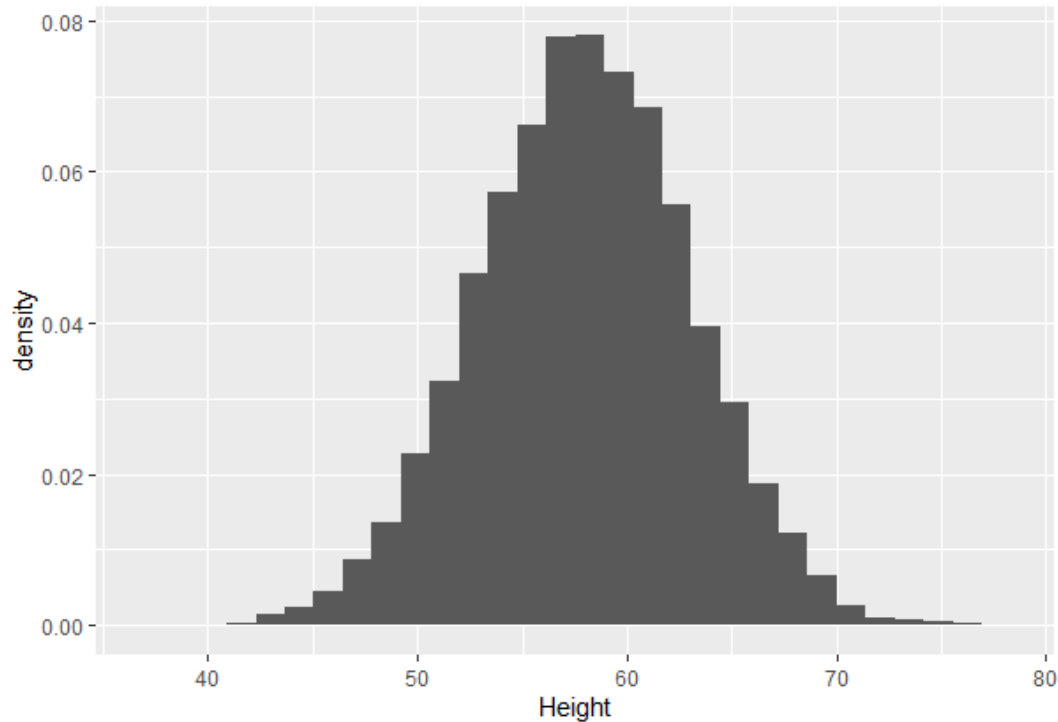
아돌프 케틀레  
1796-1874



- **“평균 인간”** : 육체적·정신적으로 한 집단의 무게중심에 해당하는 사람,  
즉 집단의 대표 Stereotype이자 신이 의도한 이상

"평균적인 프랑스군 신병" vs "평균적인 스코틀랜드군 신병"

# 평균 인간은 존재하는가?



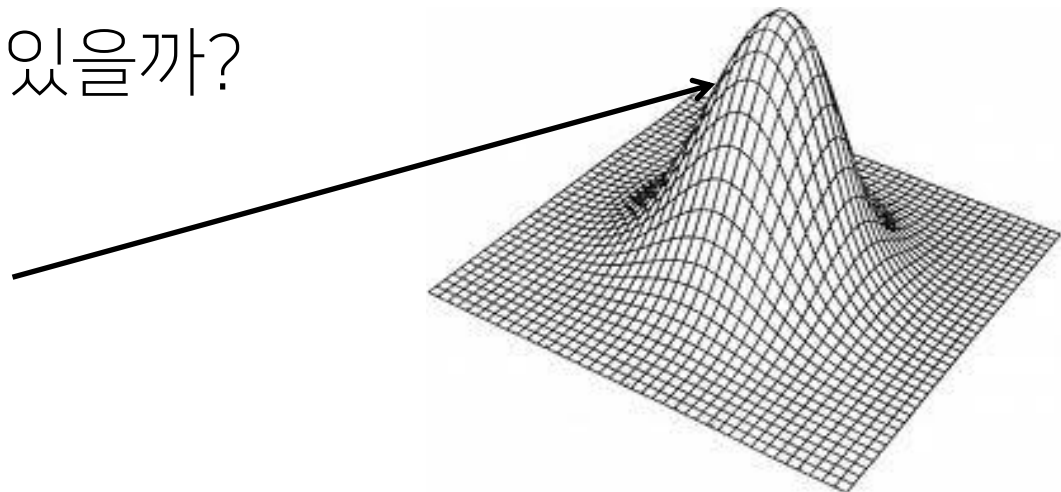
- 키가 평균(58-60인치)인 사람의 비율 =  $1/5$
- 몸무게가 평균인 사람의 비율 =  $1/5$
- 특성 두 개가 모두 평균인 사람의 비율  
=  $1/5 \times 1/5 = 1/25$
- 특성 4개가 평균인 사람의 비율?

## 고차원 데이터

$$\underline{Z} = \begin{pmatrix} Z_1 \\ \vdots \\ Z_d \end{pmatrix} \sim N_d(\mu, I_d)$$

차원  $d$ 가 매우 클 때, 데이터가 어디 있을까?

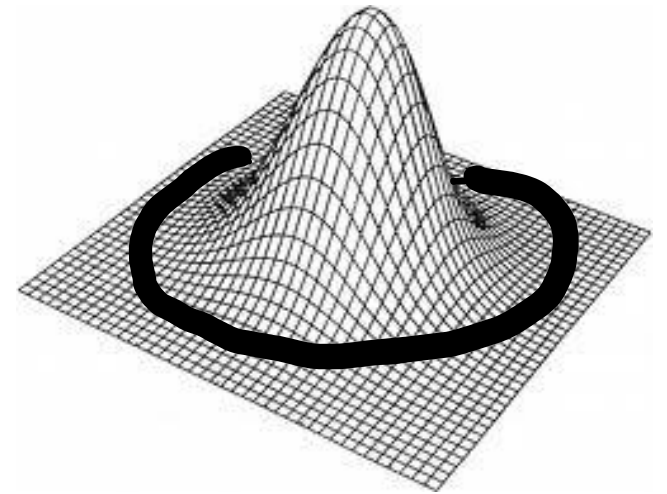
밀도함수의 꼭대기?



차원  $d$ 가 늘어날 때,

$$\|\underline{z} - \underline{\mu}\| = \sqrt{d} + O_p(1)$$

- 데이터는 (반지름이  $\sqrt{d}$ 인) 구체의 표면 위에 존재
- 구체의 중심 (“평균”) 은 여전히 가장 높은 밀도!
- 역설의 해결: 르벡 측도와 비교한 밀도함수
- “평균인간은 없다”의 수학적 증명?

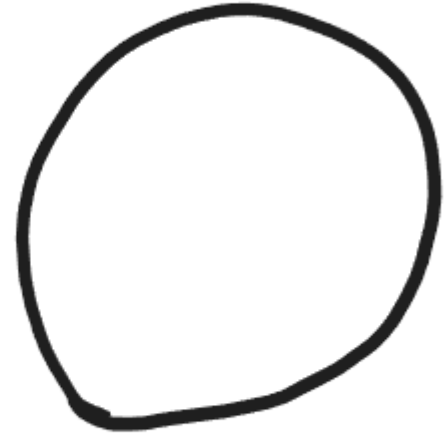
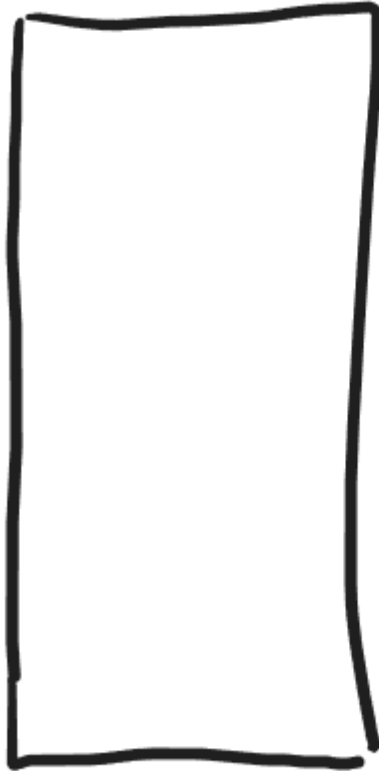


Part 2,

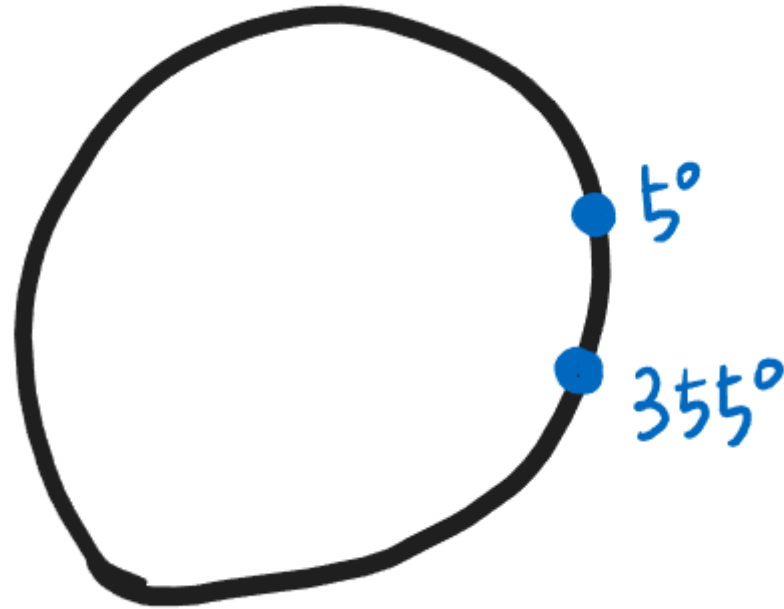
평균 낼 수 없는 것들



J. S. Marron

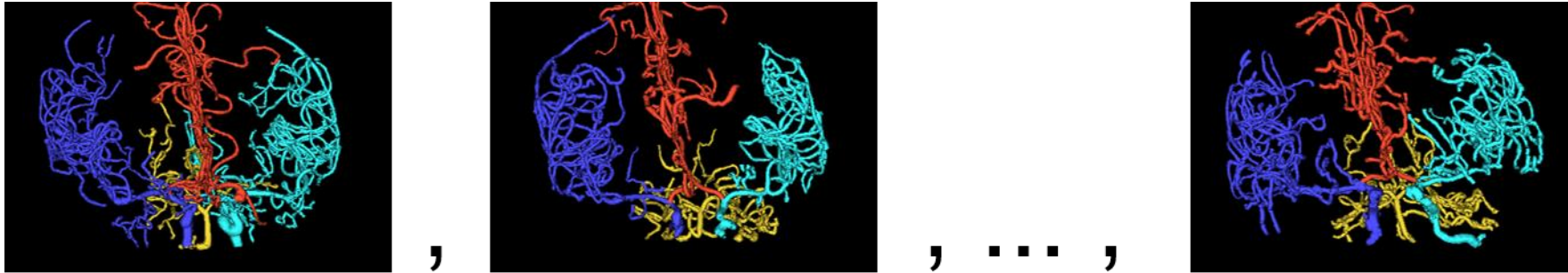


# 방향 데이터와 평균

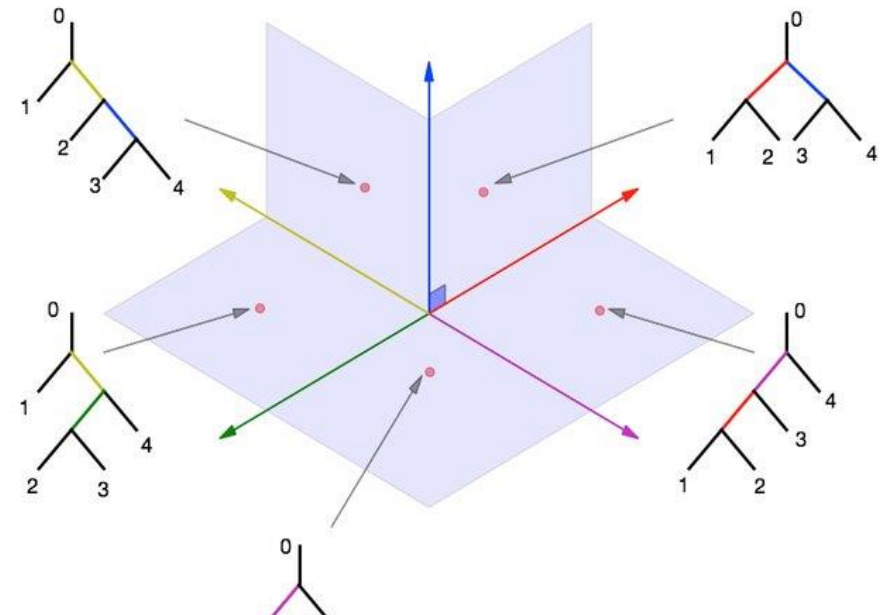


산술평균  $\neq$  대표값

- Brain artery trees (Bendich et al. 2016)



- Phylogenetic trees (Dinh et al. 2018)

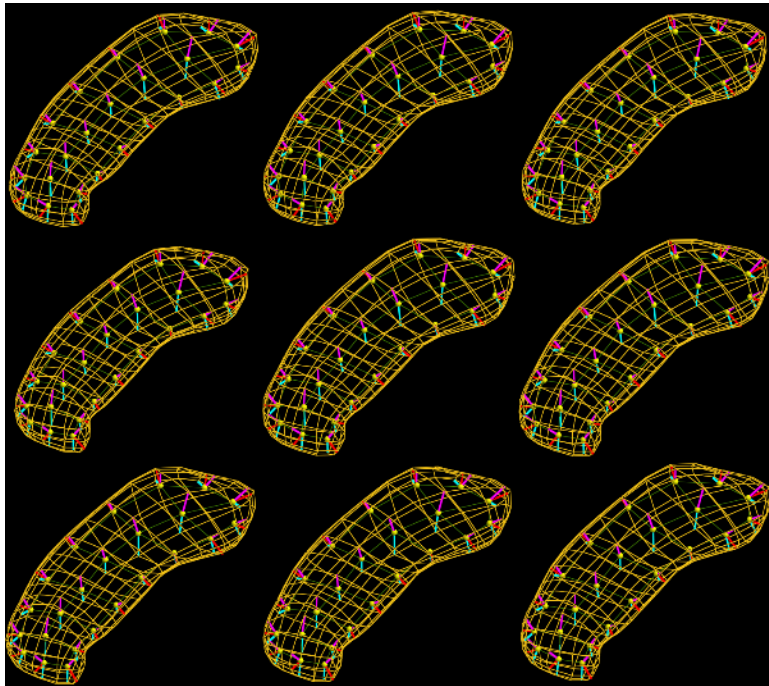


Bendich, P., Marron, J. S., Miller, E., Pieloch, A., and Skwerer, S. (2016), "Persistent homology analysis of brain artery trees," *Ann. Appl. Stat.*, 10, 198–218.

Dinh, V., Ho, L. S. T., Suchard, M. A., & Matsen IV, F. A. (2018). Consistency and convergence rate of phylogenetic inference via regularization. *Annals of statistics*, 46(4), 1481.



- Shapes (Pizer et al. 2013)



Stephen M. Pizer, Sungkyu Jung, Dibyendusekhar Goswami, Xiaojie Zhao, Ritwik Chaudhuri, James N. Damon, Stephan Huckemann, J. S. Marron (2013). "Nested Sphere Statistics of Skeletal Models", in *Innovations for Shape Analysis: Models and Algorithms*. M. Breus, Bruckstein and Maragos (Eds), Springer, Berlin, 93-115.

- Diffusions (Groisser et al. 2021)

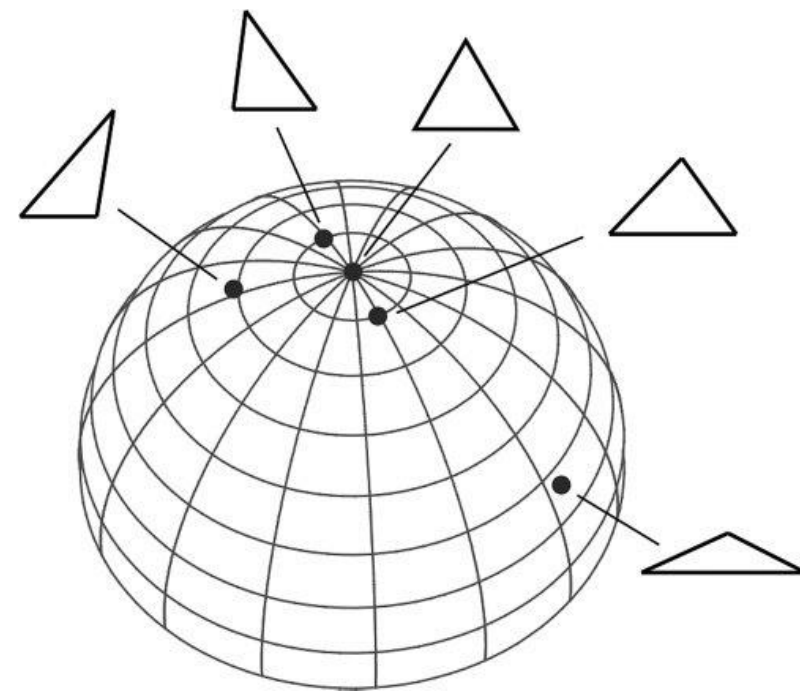
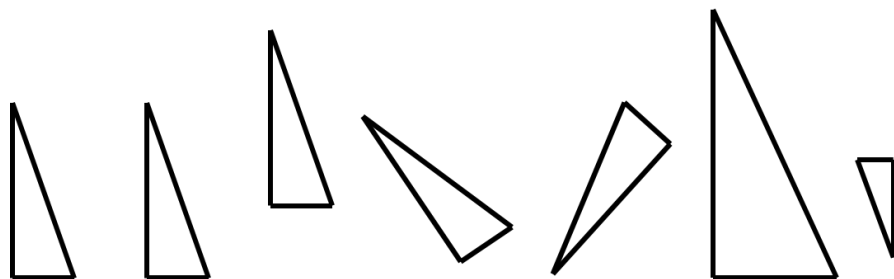


David Groisser, Sungkyu Jung and Armin Schwartzman (2021). "Uniqueness questions in a scaling-rotation geometry on the space of symmetric positive-definite matrices." *Differential Geometry and its Applications*. Volume 79, December, 101798

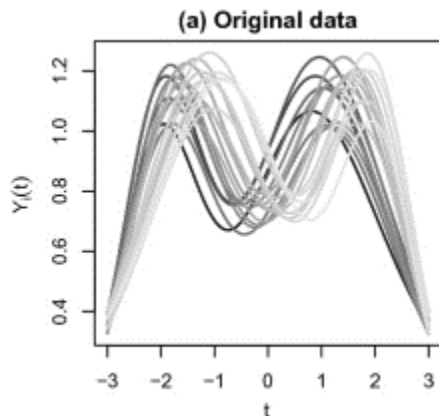
# 매니폴드 위의 데이터와 평균

- 표본 공간 = 매니폴드 (다양체)  
(분석 대상의 특성을 반영하는 기하적인 제약이 있는 공간)

Shapes

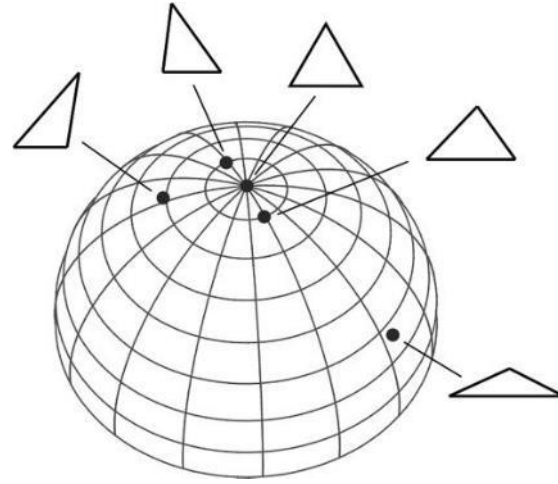


Distributions

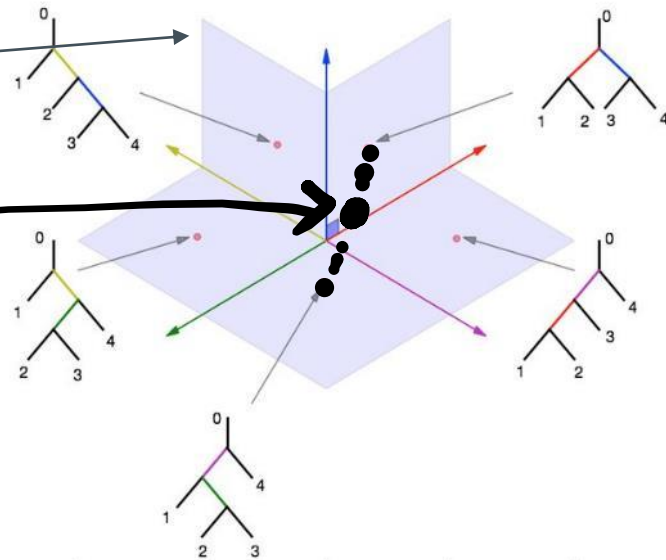


## 산술평균 $\notin$ 표본공간

- 삼각형들 = 길이가 1인 벡터들  $\in S^2$   
평균의 길이  $\neq 1$
- 트리들  $\in$  BHV tree space\*



$$\frac{\begin{array}{c} 0 \\ \diagup \quad \diagdown \\ 1 \quad 4 \\ \diagdown \quad \diagup \\ 2 \quad 3 \end{array} + \begin{array}{c} 0 \\ \diagdown \quad \diagup \\ 1 \quad 4 \\ \diagup \quad \diagdown \\ 2 \quad 3 \end{array}}{2} = ??$$

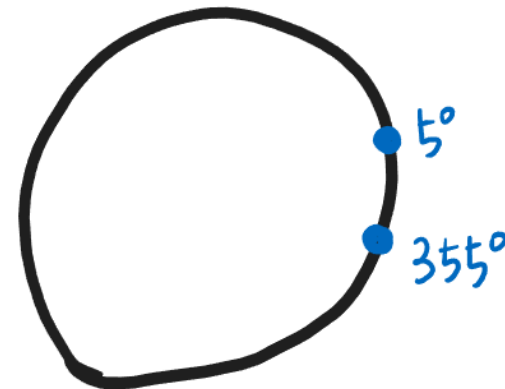


\* Billera L, Holmes S, Vogtman K. Geometry of the space of phylogenetic trees. *Adv Appl Math.* 2001;27:733–767.

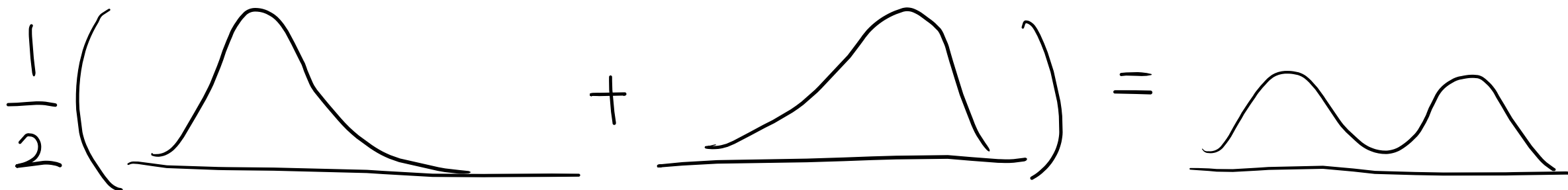
산술평균  $\neq$  대표값

- 각도들  $\in [0, 360)$

$$\text{평균} = \frac{5^\circ + 355^\circ}{2} = 180^\circ \neq \text{대표값 } (0^\circ)$$



- 실수 값을 가지는 분포들 = 밀도함수들  $\in F$



# 프레셰 평균 (Fréchet mean)

- 표본공간  $M$ 
  - $M = \mathbb{R}^p, S^p, \text{BHV tree space, 등등} \dots$
- 표본공간 위의 거리  $d(\cdot, \cdot): M \times M \rightarrow [0, \infty)$

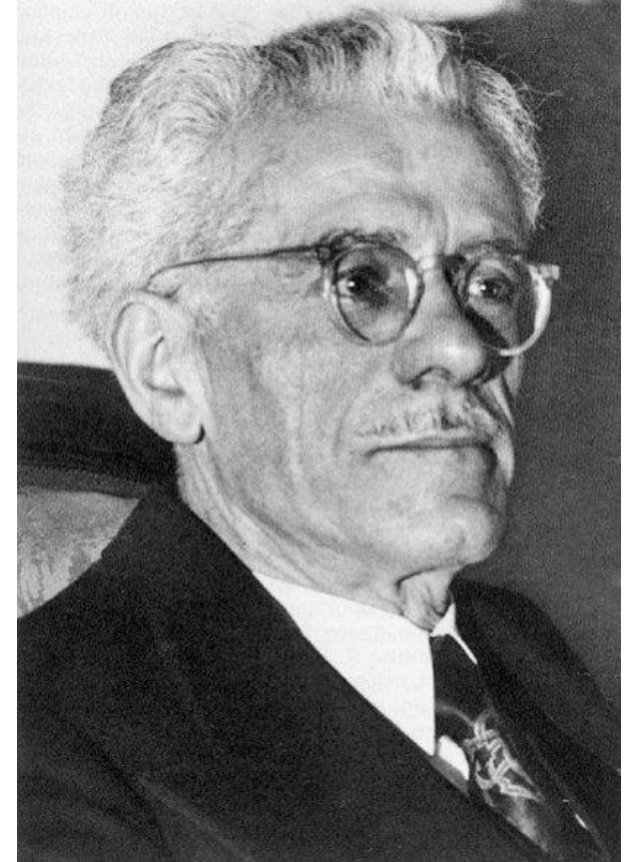
## “가우스의 최소제곱법”

만약  $M$ 이 벡터공간 ( $\mathbb{R}$ ) 이라면,  $x_1, \dots, x_n \in \mathbb{R}$  의

$$\text{평균: } \bar{x} = \arg \min_{\mu} \sum_{i=1}^n (x_i - \mu)^2$$

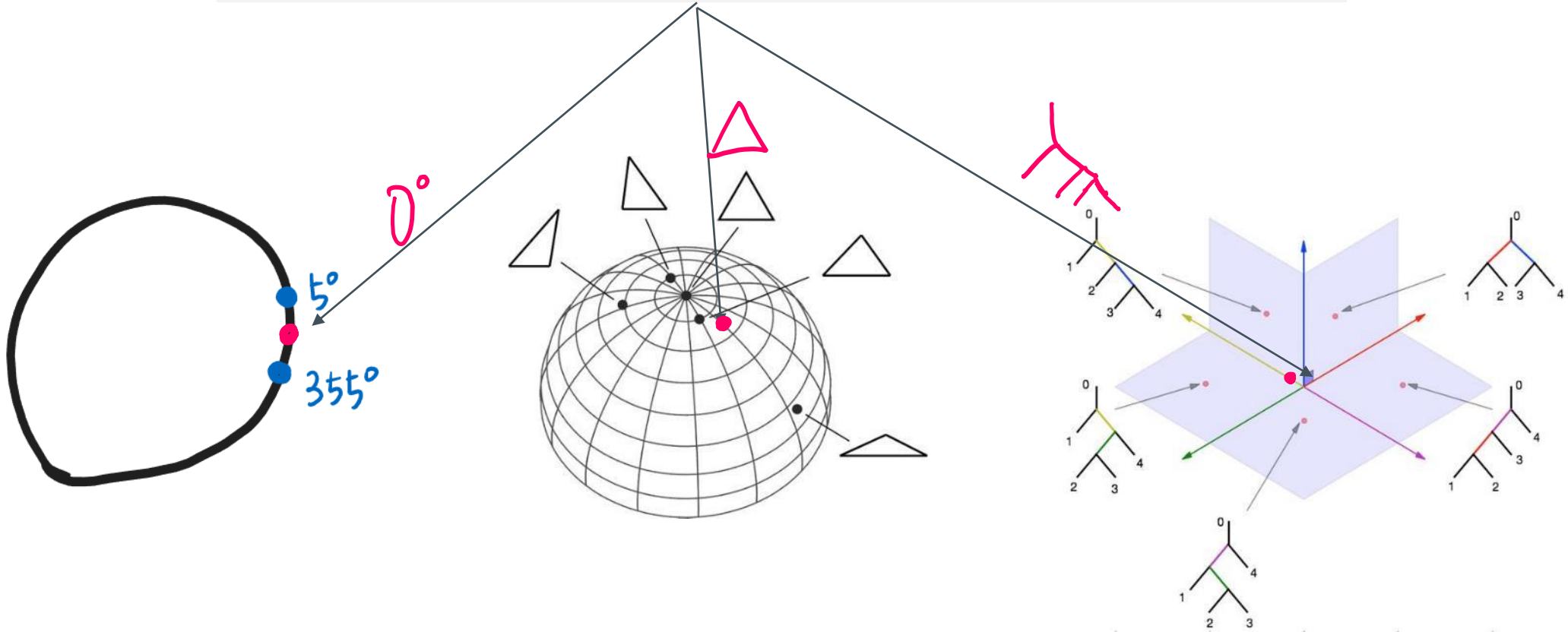
- $x_1, \dots, x_n \in (M, d)$  의 프레셰 평균:

$$\bar{x} \equiv \arg \min_{\mu \in M} \sum_{i=1}^n d(x_i, \mu)^2$$



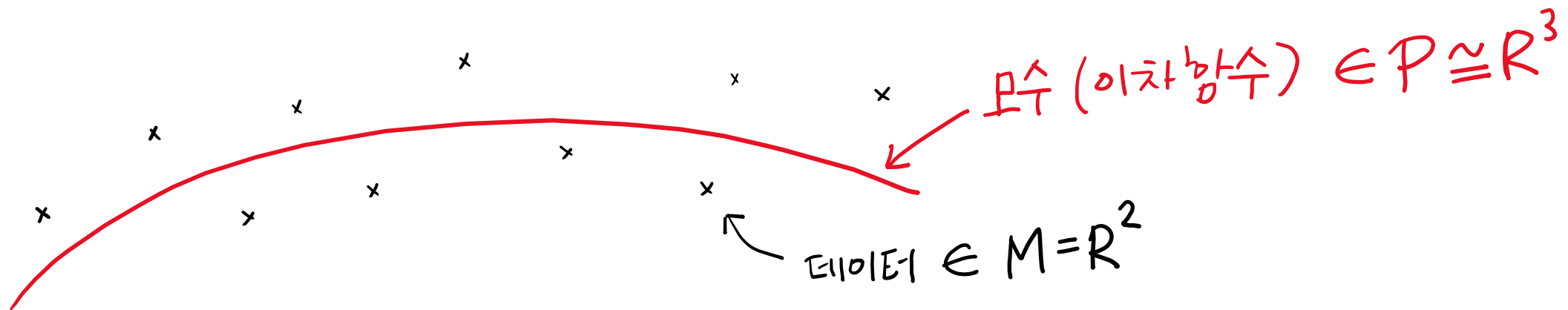
Maurice René Fréchet (1878-1973)

프레쉬 평균 = 오차제곱합을 최소화 하는 점



# 일반화된 프레쉬 평균\*

- 표본공간  $M$ 
  - $M = \mathbb{R}^p, S^p, \text{BHV tree space, 등등...}$
- 모수공간  $P$  "모수 = 분포를 기술하는 값"
  - 모수가 프레쉬 평균:  $P = M$
  - 모수가 이차회귀곡선: 표본공간  $M = \mathbb{R}^2 = \{(x,y)\}$   
 모수공간  $P = \{\text{이차함수들}\} = \{(a,b,c) \in \mathbb{R}^3\}$



# 일반화된 프레쉬 평균\*

- 표본공간  $M$ 
  - $M = \mathbb{R}^p, S^p, \text{BHV tree space, 등등...}$
- 모수공간  $P$  “모수 = 분포를 기술하는 값”
  - 모수가 프레쉬 평균:  $P = M$
  - 모수가 이차회귀곡선: 표본공간  $M = \mathbb{R}^2 = \{(x,y)\}$   
 모수공간  $P = \{\text{이차함수들}\} = \{(a,b,c) \in \mathbb{R}^3\}$
- 모수가 분포를 얼마나 잘 기술하는지에 대한 측도

$$\rho(\cdot, \cdot): P \times M \rightarrow [0, \infty)$$

예: 이차회귀곡선의 최소제곱추정:  $\rho((a,b,c), (x,y)) = |y - (a+bx+cx^2)|$

$$\arg \min_{a,b,c} \sum_{i=1}^n \rho((a,b,c), (x_i, y_i))^2$$

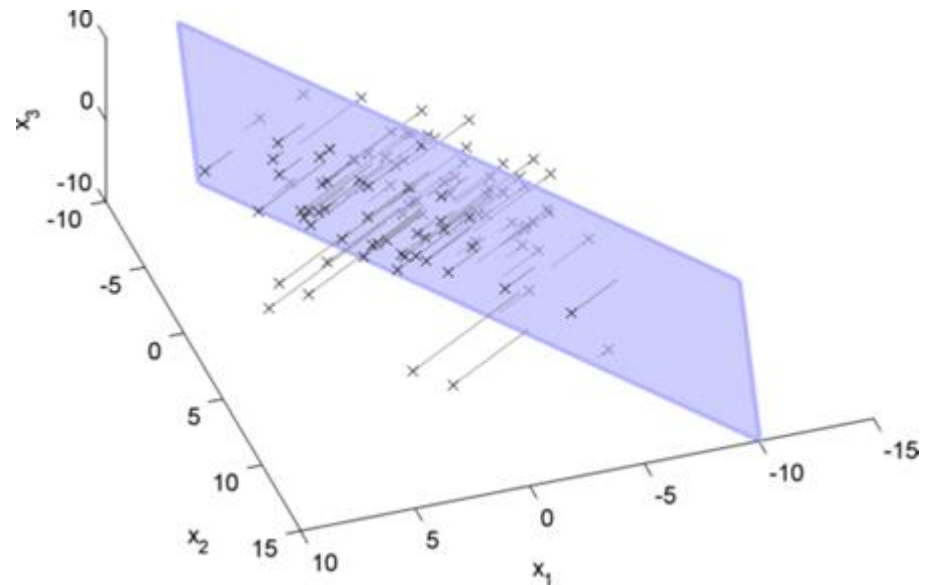


# 일반화된 프레쉬 평균

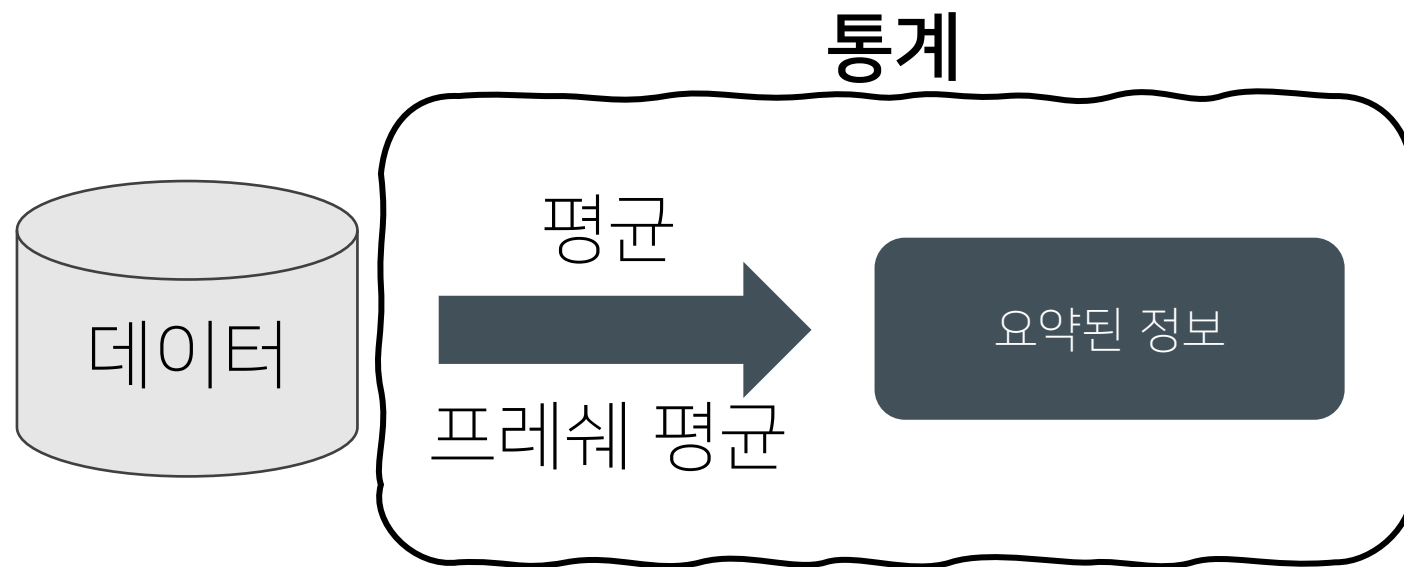
- 표본공간  $M$ 
  - $M = \mathbb{R}^p, S^p$  등의 매니폴드...
- 모수공간  $P$  “모수 = 분포를 기술하는 값”
  - 모수가 (차원축소된) “submanifold”: 모수공간  $P = \{M_q \subset M: \text{차원이 } q\}$
- $\rho$ : 작은 차원의  $M_q$  가 주어진 분포를 얼마나 잘 기술하는지에 대한 척도

예: 주성분분석

매니폴드 데이터의 차원축소



# 평균이 얼마나 정확한가?



요약된 "정보"가 얼마나 정확한가?

# 평균의 불확실성과 중심극한정리

- $x_1, x_2, \dots, x_n \sim iid (\mu, \sigma^2)$

평균  $\bar{x}_n$ 이 얼마나 달라질 수 있을까?

## “중심극한정리”

- 점근분포  $\bar{x}_n \sim N\left(\mu, \frac{\sigma^2}{n}\right)$

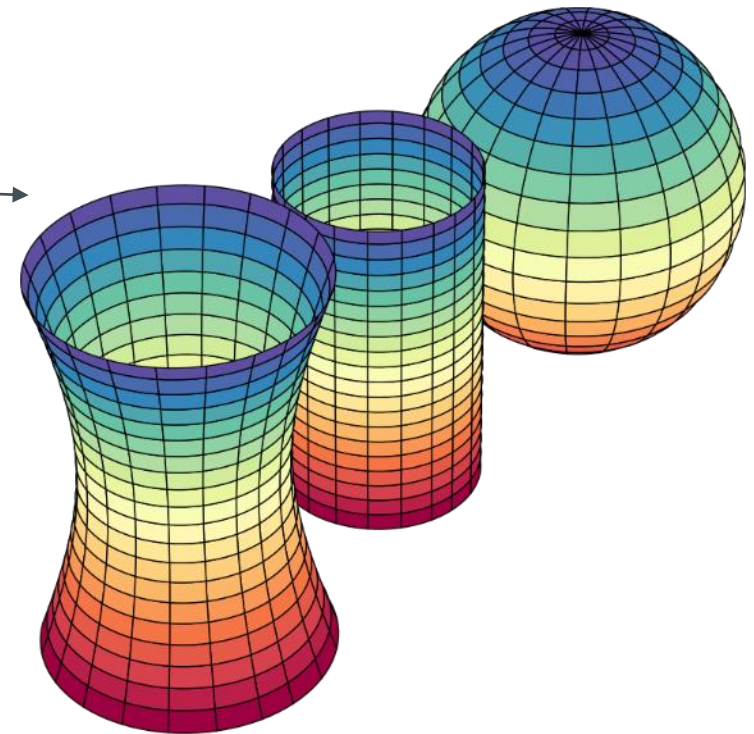
“관측값이 많을 때, 평균의 불확실성 패턴이 정규분포로 기술됨”

“평균의 표준편차가  $1/\sqrt{n}$ 의 크기 (제곱근의 법칙)”

- 모수  $\mu (= \arg \min_m E d(x, m)^2)$ 의 추정량  $\bar{x}_n$ 의 타당성

## 중심극한정리 on Riemannian manifold (리만 다양체)

- 표본공간  $M$ ,  $x_1, x_2, \dots, x_n \sim iid F$
  - 모수공간  $P$  (리만 다양체)
  - 일반화 프레쉬 표본평균  $m_n \in P$
  - 일반화 프레쉬 모평균  $m \in P$
- 평균  $m_n$ 이 얼마나 달라질 수 있을까?  
 $m_n$ 의 (점근)분포가 정규분포?  
 벡터 공간이 아닌  $P$ 에서의 정규분포?



# 리만 다양체에서의 정규분포

- 모수공간  $P$  (리만 다양체)에서의 정규분포?

- 정규분포?

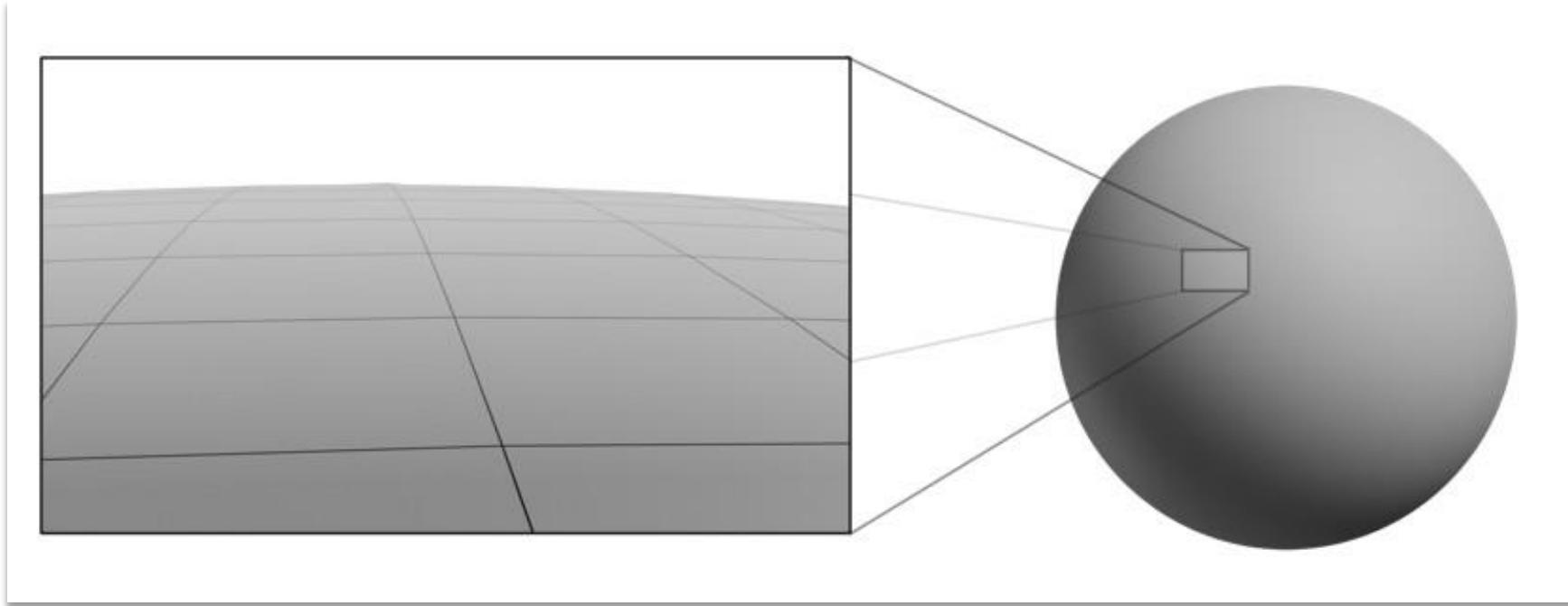
1. 밀도함수  $\propto \exp[-c(x-\mu)^2]$

2. 변동이 같은 모든 분포 중 엔트로피가 가장 큼

3. 열확산 방정식에 따라 확산하는 한 입자의 일정 시간 뒤의 위치



# 리만 다양체에서의 정규분포



3. 열확산 방정식에 따라 확산하는 한 입자의 일정 시간 뒤의 위치
  - 국소적으로 (매우 작은 평면에서) 정의되는 정규분포

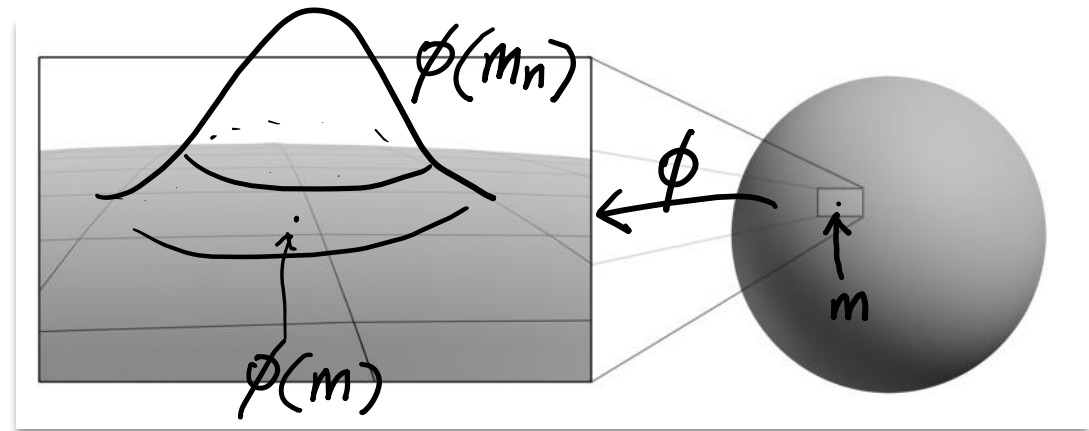
## 중심극한정리 on Riemannian manifold (리만 다양체)

- 리만 다양체: 국소적으로 평평함.
- 큰 수의 법칙\*: 프레쉬 평균  $m_n \rightarrow m$
- $m$  주변의 한없이 작고 평평한 "평면"에서의 정규분포 필요
- 중심극한정리:

$m_n$ 의 분포가  $m$ 에 한없이 가까워지므로

$m$  주변의 매우 작은 평면에서의 분포가

정규분포:  $\Phi(m_n) - \Phi(m) \sim N(0, n^{-1}\Sigma)$



# 프레쉬 평균의 불확실성?

## 중심극한정리 on Riemannian manifold (리만 다양체)

- 프레쉬 평균  $m_n \rightarrow m$
- 중심극한정리:  $\Phi(m_n) - \Phi(m) \sim N(0, n^{-1}\Sigma)$

“관측값이 많을 때, 프레쉬 평균의 불확실성 패턴이 정규분포로 기술됨”

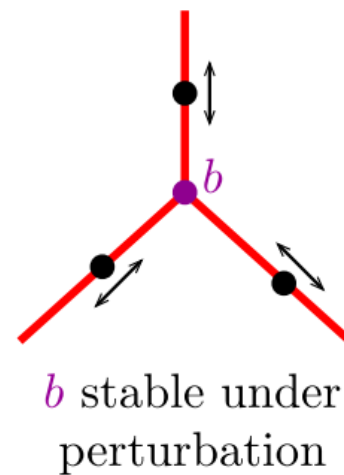
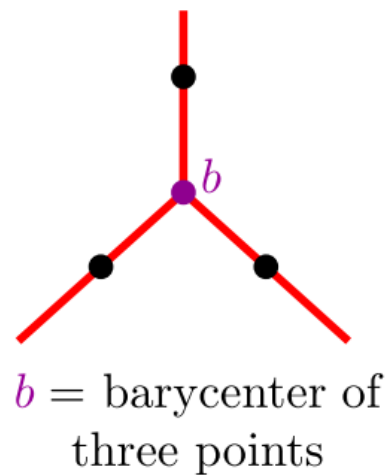
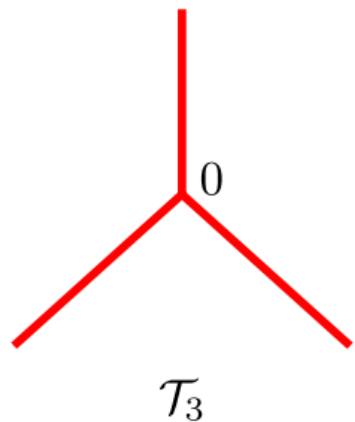
“평균의 표준편차가  $1/\sqrt{n}$  의 크기 (제공근의 법칙)”



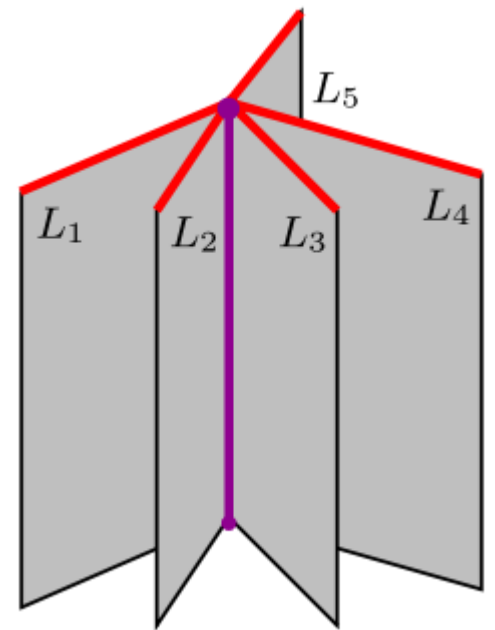
# 끈적한 평균(Sticky mean) 현상

- 층화된 모수공간  $P$ 에서  $m_n \rightarrow m$  (더 빠르게,  $1/\sqrt{n}$  보다 작은 변동)
- $P =$  거미/열린 책 (BHV tree space, Shape spaces, Space of covariance matrices)

거미 (1층: 점, 2층: 세 다각의 다리)

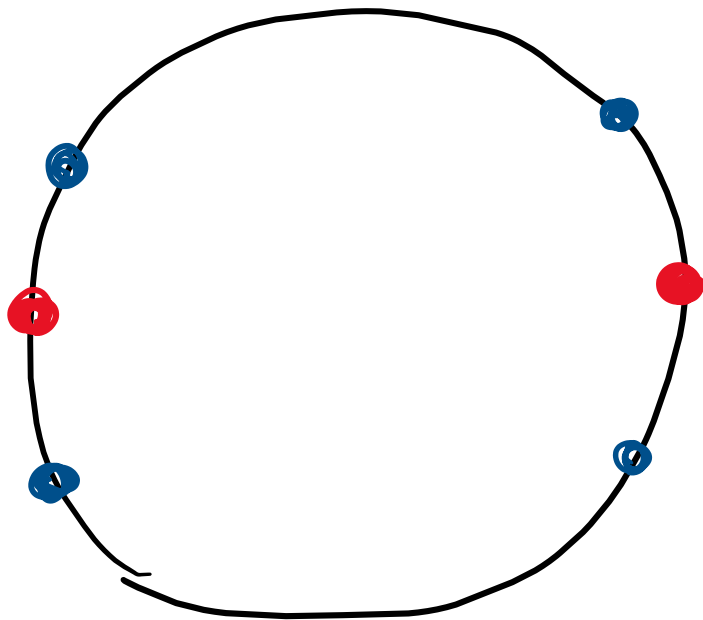


열린 책 (1층: 책등, 2층: 책장)



# 얼룩진 평균(Smeary mean) 현상

- 프레쉬 평균의 비유일성:  
분포의 범위(support)가 넓고 매니폴드의 curvature  $> 0$  일 때 종종 발생



$\{20^\circ, 160^\circ, -20^\circ, -160^\circ\}$ 의 프레쉬 평균?

$m_n = 0^\circ, 180^\circ$  두 개의 프레쉬 평균!

# 얼룩진 평균(Smeary mean) 현상

- 프레셰 평균의 비유일성:  
분포의 범위(support)가 넓고 매니폴드의 curvature  $> 0$  일 때 종종 발생
- 프레셰 평균이 유일하지만 여전히 분포의 범위가 넓을 때 “얼룩진 평균”:

$$m_n \rightarrow m \text{ (더 느리게, 더 큰 변동)}$$

표준편차가  $1/\sqrt{n}$  보다 더 큼.

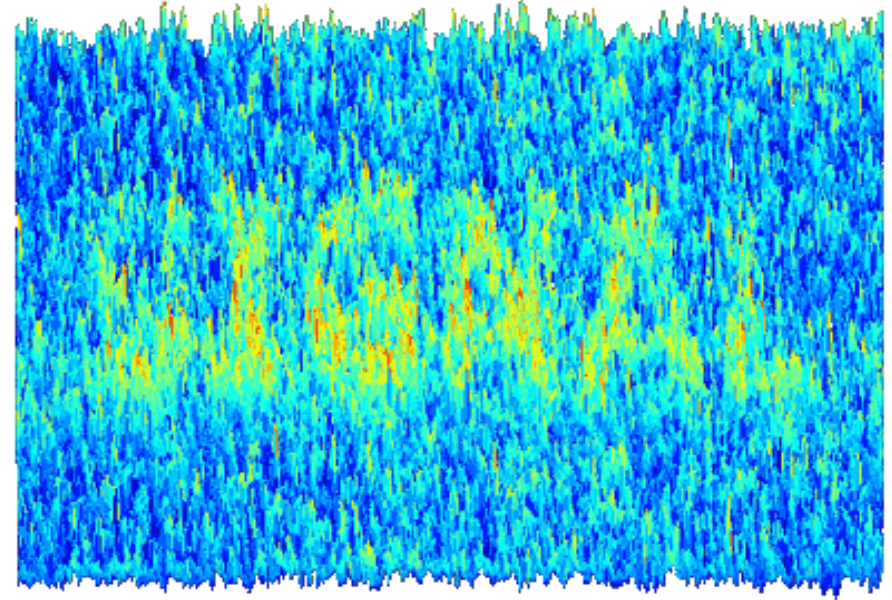
# 평균 낼 수 없는 것들: 열린 문제들

- 프레쉬 평균의 불확실성: 표본 공간  $M$ 과 모수 공간  $P$ 의 특성에 따라 다름
- Active한 연구 분야!
- “매니폴드 가설”과 빅데이터의 잠재 공간\*
- 잠재 공간이 매니폴드일 때, Manifold-specific distance/distribution?

정리

# 현대 통계학의 도전과 성취

- 데이터 = 신호와 소음
- 통계: 계산과 요약으로 "신호" 파악
- 통계학: 이 "신호"가 진짜 신호인가?
- 대표적인 요약 "평균", 산술평균 낼 수 없는 것들의 "프레쉬 평균"
- 통계학의 새로운 도전: Big Data? Complex Data?



**End of Slide**  
감사합니다