

# Statistical Learning

## Lecture 1: Statistical Decision and Estimation Theory

Mário A. T. Figueiredo

Instituto Superior Técnico & Instituto de Telecomunicações

University of Lisbon, **Portugal**

March 2024

# Statistical Decision Theory

- Specify the set of possible *decisions/action*  
**Examples:** *classify* email as spam or not; *predict* the next word in a text.
- Specify possible underlying truth (*state of nature*)  
**Examples:** email message *is/isn't* spam; actual word.
- Specify *observations*.  
**Example:** message and metadata (from, subject, IP, ...); context text.
- Specify how *observations* are related to *state of nature*.  
**Example:** word statistics.
- Specify payoff/consequences of the decision.  
**Example:** *loss* of sending a good email to the spam folder (vice-versa).
- **Goal:** derive optimal decision rules; characterize decision rules.

# Statistical Decision Theory

- **Unknown** underlying reality/truth, the *state of nature*:  $s \in \mathcal{S}$
- *Observations*: a random variable  $X \in \mathcal{X}$ , with **known**  $f_X(\cdot|s)$ ;  
It can be a pdf, pmf, hybrid, univariate, multivariate,...  
Often called **likelihood function**
- Set of allowed *decisions/action*  $\mathcal{A}$
- *Decision rules*  $\delta : \mathcal{X} \rightarrow \mathcal{A}$ , where  $\delta \in \mathcal{D}$ :  
$$\delta(x) = a \quad \text{means "if } x \text{ is observed, decide } a\text{"}$$
- *Loss function*, quantifying the decision consequences,  $L : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ :  
$$L(s, a) = \text{"the loss of deciding } a \text{ under state of nature (truth) } s\text{"}$$

**Statistical decision problem**: given  $(\mathcal{S}, \mathcal{A}, \mathcal{X}, f_X, \mathcal{D}, L)$ ,  
choose optimal  $\delta \in \mathcal{D}$ .

# Statistical Decision Problems: Examples

- **Medical decision problem:**  $\mathcal{S} = \{\text{cancer, no cancer}\}$ ,  
 $\mathcal{A} = \{\text{surgery, chemotherapy, do nothing}\}$ ,  
 $X = \text{lab test results, images, ...}$
- **Fingerprint id system:**  $\mathcal{S} = \{\text{user 1, user 2, ..., user N, stranger}\}$   
 $\mathcal{A} = \{\text{login user 1, ..., login user N, not allowed}\}$   
 $X = \text{data acquired at the fingerprint sensor}$
- **Parameter estimation:**  $\mathcal{S} = \mathcal{A} = \mathbb{R}$   
 $X = \text{observations; known } f_X(\cdot|s), \text{ for } s \in \mathcal{S}$
- **Binary decision:**  $\mathcal{S} = \mathcal{A} = \{0, 1\}$   
 $X = \text{observed data; } f_X(\cdot|s) \text{ observation model.}$
- **Binary decision with reject option:**  $\mathcal{S} = \{0, 1\}$ ;  $\mathcal{A} = \{0, 1, R\}$

**Critical aspect:** (unrealistically) assumes perfect knowledge of  $f_X$

# Statistical Decision Problems: Frequentist Risk

**Assumptions:**  $s \in \mathcal{S}$  is fixed but unknown

decisions should have low expected (w.r.t.  $X \sim f_X$ ) loss

The **risk**:  $R(s, \delta) = \mathbb{E}_X [L(s, \delta(X))|s]$

$$\mathbb{E}_X [L(s, \delta(X))|s] = \begin{cases} \int_{\mathcal{X}} L(s, \delta(x)) f_X(x|s) dx & \Leftarrow \mathcal{X} \text{ is continuous} \\ \sum_{x \in \mathcal{X}} L(s, \delta(x)) f_X(x|s) & \Leftarrow \mathcal{X} \text{ is discrete;} \end{cases}$$

**Question:** how to use the risk to compare/order different decisions?

The risk depends on  $s$ , and  $s$  is **unknown!**

# Domination and Admissibility

The *risk*:  $R(s, \delta) = \mathbb{E}_X [L(s, \delta(X)) | s]$

...induces a **partial order**: **domination**

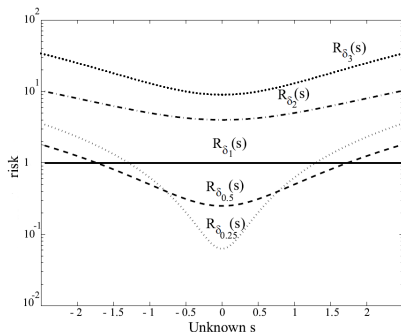
Definition: A rule  $\delta \in \mathcal{D}$  is **dominated** by another rule  $\delta' \in \mathcal{D}$  if

- (i) for any  $s \in \mathcal{S}$ ,  $R(s, \delta') \leq R(s, \delta)$ , and
- (ii) there exists  $s_0 \in \mathcal{S}$  such that  $R(s_0, \delta') < R(s_0, \delta)$ .

Definition: a rule  $\delta \in \mathcal{D}$  is **admissible** if it is not dominated by any other rule in  $\mathcal{D}$

## Admissible Rules: Parameter Estimation Example

- unknown real parameter  $s \in \mathcal{S} = \mathbb{R}$ , to be estimated ( $\mathcal{A} = \mathcal{S} = \mathbb{R}$ )
- single Gaussian observation  $X \in \mathbb{R}$ , with  $f_X(x) = \mathcal{N}(x; s, 1)$
- allowed decision rules are linear:  $\mathcal{D} = \{\delta_k(x) = kx, k \in \mathbb{R}_+\}$
- loss function:  $L(s, a) = (s - a)^2$
- the risk is  $R(s, \delta_k) = \mathbb{E}[(s - kX)^2 | s] = s^2(1 - k)^2 + k^2$



Any  $\delta_k$ , for  $k > 1$ , is not admissible. All  $\delta_k$ , for  $k \leq 1$ , are not ordered

# Minimax Rules

Definition: **minimax risk**:

$$\min_{\delta \in \mathcal{D}} \max_{s \in \mathcal{S}} R(s, \delta)$$

(with inf and sup in case  $\mathcal{D}$  and  $\mathcal{S}$  are open)

Definition:  $\delta^*$  is a **minimax rule** if

$$\max_{s \in \mathcal{S}} R(s, \delta^*) = \min_{\delta \in \mathcal{D}} \max_{s \in \mathcal{S}} R(s, \delta)$$

Back to the previous estimation example:  $R(s, \delta_k) = s^2(1 - k)^2 + k^2$

$$\sup_{s \in \mathbb{R}} R(s, \delta_k) = \sup_{s \in \mathbb{R}} \{s^2(1 - k)^2 + k^2\} = \begin{cases} 1 & \Leftarrow k = 1 \\ \infty & \Leftarrow k \neq 1 \end{cases}$$

...thus  $\delta_1$  is the (unique) minimax rule.



# Frequentist Characterization of Estimators

Estimation problems ( $\mathcal{S} = \mathcal{A} = \mathbb{R}$ ): goal is to **estimate**  $s \in \mathbb{R}$ .

Observations:  $X = (X_1, \dots, X_n) \in \mathbb{R}^n$ , with  $f_X(\cdot|s)$

- **Consistent** estimator:  $\lim_{n \rightarrow \infty} \delta(X) = s$ , in probability ( $\delta(X)$  is a r.v.),  
i.e.,

$$\forall \varepsilon > 0, \lim_{n \rightarrow \infty} \mathbb{P}(|\delta(X) - s| > \varepsilon) = 0$$

- **Bias** of an estimator:  $\mathbb{E}[\delta(X)] - s$  (if  $\mathbb{E}[\delta(X)] = s$ ,  $\delta$  is **unbiased**)
- **Variance** of an estimator:  $\mathbb{E}\left[\left(\delta(X) - \mathbb{E}[\delta(X)]\right)^2\right]$
- **MVU** (*minimum variance unbiased*) estimator,  $\delta_{\text{MVU}} \in \mathcal{D}$

$$\mathbb{E}[\delta_{\text{MVU}}(X)] = s \quad \text{and} \quad \mathbb{E}\left[(\delta_{\text{MVU}}(X) - s)^2\right] \leq \mathbb{E}\left[(\delta(X) - s)^2\right]$$

for any other unbiased  $\delta \in \mathcal{D}$ . Also known as **efficient estimator**.

# Frequentist Characterization of Estimators: Example 1

Independent and identically distributed (i.i.d.) observations:  $X = (X_1, \dots, X_n)$ ,

$$f_X(x|s) = \prod_{i=1}^n \mathcal{N}(x_i; s, \sigma^2) = \left(\frac{1}{2\sigma^2}\right)^n \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - s)^2\right)$$

Sample mean estimator:  $\delta(x) = \frac{x_1 + \dots + x_n}{n}$

- **Unbiased:**  $\mathbb{E}[\delta(X)] = \frac{1}{n} (\mathbb{E}(X_1) + \dots + \mathbb{E}(X_n)) = s$
- **Consistent:**  $\delta(X) \sim \mathcal{N}(s, \phi_n^2)$ ,

$$\phi_n^2 = \text{var}\left(\frac{X_1}{n} + \dots + \frac{X_1}{n}\right) = n \frac{\sigma^2}{n^2} = \frac{\sigma^2}{n} \xrightarrow{n \rightarrow \infty} 0$$

- **Minimum variance?** ...more later.

## Frequentist Characterization of Estimators: Example 2

Observations:  $X = (X_1, \dots, X_n) \in \mathbb{R}^n$ , where  $s = (s_1, s_2) = (\mu, \sigma^2)$  and

$$f_X(x|s) = \prod_{i=1}^n \mathcal{N}(x_i; \mu, \sigma^2) = \left(\frac{1}{2\sigma^2}\right)^n \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right)$$

Sample variance:  $\delta_2(x) = \frac{1}{n} \sum_{i=1}^n (x_i - \delta_1(x))^2$ , where  $\delta_1(x) = \frac{1}{n} \sum_{i=1}^n x_i$

We have seen that  $\delta_1$  is unbiased and consistent; what about  $\delta_2$ ?

$$\bullet \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n (X_i - \frac{1}{n} \sum_{j=1}^n X_j)^2 \right] = \frac{n-1}{n} \sigma^2 \quad (\text{using } \mathbb{E}[X_i^2] = \mu^2 + \sigma^2; \\ \mathbb{E}[X_i X_j] = \mu^2)$$

$$\bullet \text{Unbiased variance estimator: } \hat{\sigma}^2(x) = \frac{1}{n-1} \sum_{i=1}^n \left(x_i - \frac{1}{n} \sum_{j=1}^n x_j\right)^2$$

# Minimum Variance and the Cramér-Rao Bound

Let  $s \in \mathbb{R}$  be a (scalar) parameter to be estimated from  $X \sim f_X(\cdot|s)$

**Cramér-Rao bound:**

$$\underbrace{\mathbb{E}[\delta(X)] = s}_{\delta \text{ is unbiased}} \Rightarrow \text{var}[\delta(X)] \equiv \mathbb{E}[(\delta(X) - s)^2] \geq \frac{1}{I(s)},$$

**Fisher information:**  $I(s) = -\mathbb{E} \left[ \frac{\partial^2 \log f_X(X|s)}{\partial s^2} \right]$  (log = log<sub>e</sub>)

(curvature)

**Example:**  $X = (X_1, \dots, X_n) \in \mathbb{R}^n$ , with  $f_X(x|s) = \prod_{i=1}^n \mathcal{N}(x_i; s, \sigma^2)$

$$-\log f_X(x|s) = K + \frac{\sum_{i=1}^n (X_i - s)^2}{2\sigma^2} \Rightarrow I(s) = \frac{\mathbb{E} \left[ \frac{\partial^2 \sum_{i=1}^n (s - X_i)^2}{\partial s^2} \right]}{2\sigma^2} = \frac{n}{\sigma^2}$$

thus, since  $\delta(x) = \frac{x_1 + \dots + x_n}{n}$  is unbiased and has variance  $\frac{\sigma^2}{n}$ , it is MVU

## Cramér-Rao Bound: Proof (simplified)

- **Score:**  $V(X) = \frac{\partial \log f_X(X|s)}{\partial s} = \frac{1}{f_X(X|s)} \frac{\partial f_X(X|s)}{\partial s}$  (random)
- $\mathbb{E}[V(X)] = \int_{\mathcal{X}} f_X(x|s) \frac{1}{f_X(x|s)} \frac{\partial f_X(x|s)}{\partial s} dx = \frac{\partial}{\partial s} \int_{\mathcal{X}} f_X(x|s) dx = 0$
- $\text{cov}(V(X), \delta(X)) = \mathbb{E}[V(X)\delta(X)] - \underbrace{\mathbb{E}[V(X)]}_{0} \mathbb{E}[\delta(X)] = \mathbb{E}[V(X)\delta(X)]$
- $\mathbb{E}[V(X)\delta(X)] = \int_{\mathcal{X}} \delta(x) \frac{\partial f_X(x|s)}{\partial s} dx = \frac{\partial}{\partial s} \underbrace{\int_{\mathcal{X}} \delta(x) f_X(x|s) dx}_{\mathbb{E}[\delta(X)] = s} = 1$
- $\text{cov}(V(X), \delta(X)) = 1$

## Cramér-Rao Bound: Proof (continuation)

- Fisher information is the negative score variance:

$$\begin{aligned} I(s) &= -\mathbb{E} \left[ \frac{\partial^2 \log f_X(X|s)}{\partial s^2} \right] \\ &= -\mathbb{E} \left[ \frac{\partial}{\partial s} \left( \frac{1}{f_X(X|s)} \frac{\partial f_X(X|s)}{\partial s} \right) \right] \\ &= \mathbb{E} \left[ \frac{1}{(f_X(X|s))^2} \left( \frac{\partial f_X(X|s)}{\partial s} \right)^2 \right] - \mathbb{E} \left[ \frac{1}{f_X(X|s)} \frac{\partial^2 f_X(X|s)}{\partial s^2} \right] \\ &= \underbrace{\mathbb{E} \left[ \left( \frac{\partial \log f_X(X|s)}{\partial s} \right)^2 \right]}_{\text{var}(V(X))} - \underbrace{\frac{\partial^2}{\partial s^2} \int_{\mathcal{X}} f_X(x|s) dx}_0 \end{aligned}$$

- Cauchy-Schwartz:**  $|\text{cov}(V(X), \delta(X))| \leq \sqrt{\text{var}(V(X)) \text{var}(\delta(X))}$

$$1 = |\text{cov}(V(X), \delta(X))| \Rightarrow \text{var}(\delta(X)) \geq \frac{1}{I(s)} \quad \square$$

## Fisher Information and Cramér-Rao Bound: Vector Case

Let  $s \in \mathbb{R}^p$  be a (vector) parameter to be estimated from  $X \sim f_X(\cdot|s)$

- **Score:**  $V(X) = \nabla_s \log f_X(X|s) = \frac{\nabla_s f_X(X|s)}{f_X(X|s)}$  (random vector)

- As in the scalar case:  $\mathbb{E}[V(X)] = 0$

- **Fisher information matrix:**

$$\begin{aligned} I(s) &= -\mathbb{E} [\nabla_s^2 \log f_X(X|s)] \\ &= \mathbb{E} [(\nabla_s \log f_X(X|s)) (\nabla_s \log f_X(X|s))^T] = \text{cov}(V(X)) \end{aligned}$$

- **Cramér-Rao bound:** for  $\delta : \mathcal{X} \rightarrow \mathbb{R}^n$ ,

$$\underbrace{\mathbb{E}[\delta(X)] = s}_{\text{unbiased}} \Rightarrow \text{cov}(\delta(X)) \succeq I(s)^{-1},$$

*i.e.*  $A - B$  is positive semi-definite:  $u^T A u \geq u^T B u$ , for any  $u \in \mathbb{R}^n$

# Fisher Information and Cramér-Rao Bound: Vector Case

Let  $s \in \mathbb{R}^p$  be a (vector) parameter to be estimated from  $X \sim f_X(\cdot|s)$

- **Cramér-Rao bound:** for  $\delta : \mathcal{X} \rightarrow \mathbb{R}^n$ ,

$$\underbrace{\mathbb{E}[\delta(X)] = s}_{\text{unbiased}} \Rightarrow \text{cov}(\delta(X)) \succeq I(s)^{-1},$$

- If  $r = a^T s$  is a linear function of  $s$ ,  $\gamma(X) = a^T \delta(X)$  is its estimate.

$$\text{var}(\gamma(X)) = a^T \text{cov}(\delta(X)) a \geq a^T I(s)^{-1} a$$

- For  $n$  i.i.d. observations,  $X_i \sim f_{X_i}(\cdot|s) = f_X(\cdot|s)$ , for  $i = 1, \dots, n$ ,

$$\log f_{X_1, \dots, X_n}(x_1, \dots, x_n|s) = \log \prod_{i=1}^n f_X(x_i|s) = \sum_{i=1}^n \log f_X(x_i|s),$$

$$I_{(n)}(s) = \mathbb{E} \left[ \nabla_s^2 \sum_{i=1}^n \log f_X(X_i|s) \right] = \sum_{i=1}^n \mathbb{E} \left[ \nabla_s^2 \log f_X(X|s) \right] = n I_{(1)}(s)$$



# Fisher Information: Exponential Families

Canonical exponential family with parameter  $\eta \in \mathbb{R}^p$ ,

$$f_X(x|\eta) = \frac{1}{Z(\eta)} h(x) \exp(\eta^T \phi(x))$$

- **Fisher information** matrix (notice that  $\nabla_z^2 z^T a = 0$ ):

$$I(\eta) = -\mathbb{E} [\nabla_\eta^2 \log f_X(X|\eta)] = \nabla_\eta^2 \log Z(\eta)$$

- ...but, it is easy to show (try it), that

$$\frac{\partial^2}{\partial \eta_i \partial \eta_j} \log Z(\eta) = \text{cov}(\phi_i(X) \phi_j(X)),$$

thus

$$I(\eta) = \nabla_\eta^2 \log Z(\eta) = \text{cov}(\phi(X))$$

# Maximum Likelihood (ML) Estimation

Canonical exponential family with parameter  $\eta \in \mathbb{R}^n$ ,

$$f_X(x|\eta) = \frac{1}{Z(\eta)} h(x) \exp(\eta^T \phi(x))$$

- **Maximum likelihood** estimate:

$$\hat{\eta}(x) = \delta(x) = \arg \max_{\eta'} f_X(x|\eta') = \arg \max_{\eta'} [-\log Z(\eta') + (\eta')^T \phi(x)]$$

- ...thus  $\hat{\eta}(x) = \text{solution}_{\eta'} [\nabla \log Z(\eta') = \phi(x)]$

$$= \text{solution}_{\eta'} [\mathbb{E}[\phi(X)|\eta'] = \phi(x)]$$

- **Example:**  $f_X(x) = \mathcal{N}(x; \eta, 1) \propto \exp(\eta x)$

- ▶  $\phi(x) = x$
- ▶  $\mathbb{E}[\phi(X)|\eta'] = \eta'$
- ▶  $\hat{\eta}(x) = \text{solution}_{\eta'} (\eta' = x) = x$

# Maximum Likelihood (ML) Estimation: Consistency

Canonical exponential family with  $n$  i.i.d. observations  $X = (X_1, \dots, X_n)$

$$\log f_X(x|\eta) = -n \log Z(\eta) + \eta^T \sum_{i=1}^n \phi(x_i)$$

- **Maximum likelihood** estimate:

$$\hat{\eta}(x) = \text{solution}_{\eta'} \left( \mathbb{E}[\phi(X)|\eta'] = \frac{1}{n} \sum_{i=1}^n \phi(x_i) \right) \quad (\text{moment matching})$$

- Weak law of large numbers:

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \phi(x_i) = \mathbb{E}[\phi(X)|\eta]$$

- **Consistency:**

$$\lim_{n \rightarrow \infty} \hat{\eta}(x) = \text{solution}_{\eta'} \left( \mathbb{E}[\phi(X)|\eta'] = \mathbb{E}[\phi(X)|\eta] \right) = \eta$$

...also true with more generality, under some additional conditions.

# ML Estimation: Asymptotic Normality

Canonical exponential family with  $n$  i.i.d. observations  $X = (X_1, \dots, X_n)$

$$\log f_X(x|\eta) = -n \log Z(\eta) + \eta^T \sum_{i=1}^n \phi(x_i)$$

- **Asymptotic normality:** let  $\hat{\eta}(x)$  be the ML estimation function

$$\lim_{n \rightarrow \infty} \sqrt{n}(\hat{\eta}(X) - \eta) \stackrel{d}{=} Z \sim \mathcal{N}(0, I(\eta)^{-1})$$

- ...thus (true with more generality, under some conditions)

$$\lim_{n \rightarrow \infty} \hat{\eta}(X) = Z \sim \mathcal{N}\left(\eta, \frac{I(\eta)^{-1}}{n}\right)$$

- **Convergence in distribution:**  $Z_1, Z_2, \dots, Z_n, \dots$  have distribution functions  $F_1, F_2, \dots, F_n, \dots$ ;  $Z$  has distribution function  $F$

$$\lim_{n \rightarrow \infty} Z_n \stackrel{d}{=} Z \quad \text{if} \quad \lim_{n \rightarrow \infty} F_n(x) = F(x),$$

at every point  $x$  where  $F$  is continuous.

# The Bias-Variance Trade-off

Let  $s \in \mathbb{R}$  be a scalar to be estimated,  $X \sim f_X(\cdot|s)$ , and  $\delta : \mathcal{X} \rightarrow \mathbb{R}$ .

- Expected estimate:  $\bar{s} = \mathbb{E}[\delta(X)]$
- Bias-variance decomposition (of the squared error loss)

$$\begin{aligned}\mathbb{E} [(\delta(X) - s)^2] &= \mathbb{E} [(\delta(X) - \bar{s} + \bar{s} - s)^2] \\ &= \mathbb{E} [(\delta(X) - \bar{s})^2] + \mathbb{E} [(s - \bar{s})^2] + 2\mathbb{E} [(\delta(X) - \bar{s})(\bar{s} - s)] \\ &= \mathbb{E} [(\delta(X) - \bar{s})^2] + (s - \bar{s})^2 + 2(\bar{s} - s) \underbrace{(\mathbb{E}[\delta(X)] - \bar{s})}_0 \\ &= \underbrace{\mathbb{E} [(\delta(X) - \bar{s})^2]}_{\text{variance}} + \underbrace{(s - \bar{s})^2}_{\text{(squared) bias}}\end{aligned}$$

- The minimum expected loss may be achieved by a biased estimator.
- Generalizes trivially to the vector case.

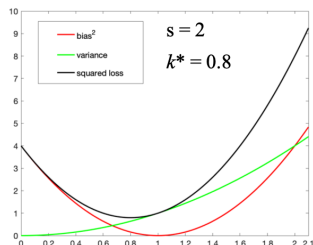
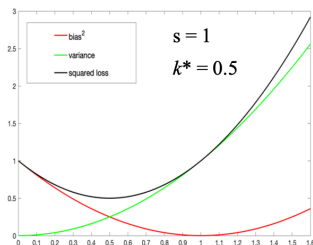
# The Bias-Variance Trade-off: Example

Let  $s \in \mathbb{R}$ ; observations  $X \sim \mathcal{N}(s, 1)$ ; linear estimator:  $\delta_k(x) = kx$ .

- Bias-variance decomposition:

$$\mathbb{E} [(\delta_k(X) - s)^2] = \underbrace{k^2}_{\text{variance}} + \underbrace{s^2(k-1)^2}_{\text{bias}^2} \quad (\text{see slide 7})$$

- Optimal  $k$  (minimum squared error loss):  $k^* = s^2 / (s^2 + 1)$



- Of course, in practice,  $s$  is unknown.

## Stein's "Paradox" and James-Stein Estimation

Estimate  $\theta \in \mathbb{R}^n$  from  $f_X(x|\theta) = \prod_{i=1}^n \mathcal{N}(x_i; \theta_i, 1)$ , with  $n \geq 3$ .

- ML estimate (unbiased):  $\hat{\theta}_{\text{ML}}(x) = \delta_{\text{ML}}(x) = x$ .
- **James-Stein estimate**:  $\hat{s}_{\text{JS}}(x) = \delta_{\text{JS}}(x) = x \underbrace{\left(1 - \frac{n-2}{\|x\|_2^2}\right)}_{<1}+$  (shrinkage)
- Risk of  $\delta_{\text{ML}}$ :  $R(\theta, \delta_{\text{ML}}) = \mathbb{E} [\|\delta_{\text{ML}}(X) - \theta\|_2^2] = \mathbb{E} [\|X - \theta\|_2^2] = n$
- Risk of  $\delta_{\text{JS}}$ :  $R(\theta, \delta_{\text{JS}}) = n - \underbrace{(n-2)^2 \mathbb{E} \left[ \frac{1}{\|X\|_2^2} \right]}_{>0} < R(\theta, \delta_{\text{ML}})$
- ...the ML estimator, in this problem, is **inadmissible**, which seems paradoxical, given conditionally independent observations.

# James-Stein Estimation: A Bias-Variance View

Estimate  $\theta \in \mathbb{R}^n$  from  $f_X(x|\theta) = \prod_{i=1}^n \mathcal{N}(x_i; \theta_i, 1)$ , with  $n \geq 3$ .

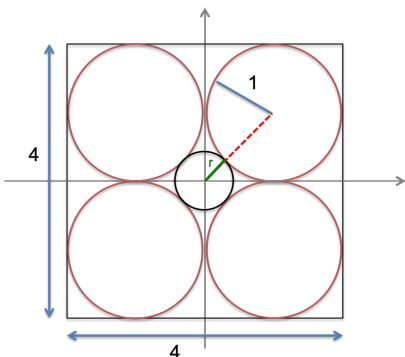
- Estimators of the form  $\delta_\lambda(x) = \lambda x$ , for  $\lambda \in \mathbb{R}_+$
- **Bias-variance** decomposition:

$$\begin{aligned} \mathbb{E} [\|\delta_\lambda(X) - \theta\|^2] &= \underbrace{\mathbb{E} [\|\delta_\lambda(X) - \underbrace{\lambda \mathbb{E}[X]}_{\mathbb{E}[\delta_\lambda(X)]}\|^2]}_{\text{variance}} + \underbrace{\|\theta - \lambda \mathbb{E}[X]\|^2}_{\text{(squared) bias}} \\ &= \lambda^2 n + (\lambda - 1)^2 \|\theta\|^2 \end{aligned}$$

- Optimal choice (not usable,  $\theta$  is unknown):  $\lambda^* = \frac{\|\theta\|^2}{\|\theta\|^2 + n} \in ]0, 1[$
- The **James-Stein** factor  $(1 - (n - 2)/\|x\|^2)_+$  is (approximately) an unbiased estimate of  $\lambda^*$ .
- Only for  $\|\theta\| \rightarrow \infty$ , does  $\lambda \rightarrow 1$ . For finite  $\|\theta\|$ : shrinking is good.



# Dessert: High-Dimensional Intuition Failure



- J. Michael Steele, *The Cauchy-Schwarz Master Class*, Cambridge, 2004.

- In  $\mathbb{R}^2$ :
  - ▶ 4 unit-radius circles inside  $[-2, 2]^2$
  - ▶ circle of radius  $r$  centred at the origin, confined by the unit-radius circle
  - ▶ distance from origin to centers:  $\sqrt{2}$
  - ▶  $r = \sqrt{2} - 1$
- In  $\mathbb{R}^d$ :
  - ▶  $2^d$  unit-radius spheres inside  $[-2, 2]^d$
  - ▶ sphere of radius  $r$  centred at the origin, confined by the unit-radius spheres
  - ▶  $r = \sqrt{d} - 1$
  - ▶ for  $d \geq 10$ , we have  $r > 2$

# Recommended Reading

- K. Murphy, “Machine Learning: A Probabilistic Perspective”, MIT Press, 2012 (chapter 6).
- L. Wasserman, “All of Statistics: A Concise Course in Statistical Inference”, Springer, 2004.
- M. Figueiredo, “Lectures Notes on Bayesian Estimation and Classification”, unpublished manuscript, available at <https://fenix.tecnico.ulisboa.pt/disciplinas/AEsta-4/2014-2015/2-semester/lecture-notes>