

Classical vs high-dimensional theory

Statistical Learning

CLAMSES - University of Milano-Bicocca

Aldo Solari

References

- Sur, Candès, E. J. (2019). A modern maximum-likelihood theory for high-dimensional logistic regression. Proceedings of the National Academy of Sciences, 116, 14516–14525

Classical theory

- It concerns the behaviour when the *sample size* $n \rightarrow \infty$
- Suppose $y_1, \dots, y_n \stackrel{i.i.d.}{\sim} y \in \mathbb{R}^p$ with mean $\mu = \mathbb{E}(y)$ and finite variance $\Sigma = \mathbb{V}\text{ar}(y)$
- *Law of large numbers*: the sample mean $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n y_i$ converges in probability to μ
- *Central limit theorem*: the rescaled deviation $\sqrt{n}(\hat{\mu}_n - \mu)$ converges in distribution to a centered Gaussian with covariance matrix Σ
- *Consistency of maximum likelihood estimation*
- Etc.

Suppose that we are given $n = 1000$ samples from a statistical model in $p = 500$ dimensions

Will theory that requires $n \rightarrow \infty$ with the dimension p remaining fixed provide useful predictions?

High-dimensional data

- The data sets arising in many parts of modern science have a “high-dimensional flavor”, with p on the same order as, or possibly larger than n

$$p \gg n$$

- Classical “large n , fixed p ” theory fails to provide useful predictions
- Classical methods can break down dramatically in high-dimensional regimes

Neyman-Scott problem (Bartlett, 1937)

- Let (x_i, y_i) be independent $N(\mu_i, \sigma^2)$ for $i = 1, \dots, n$
- The MLE for μ_i is $\hat{\mu}_i = (x_i + y_i)/2$
- The MLE for σ^2 is

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n s_i^2$$

where

$$s_i^2 = [(x_i - \hat{\mu}_i)^2 + (y_i - \hat{\mu}_i)^2]/2 = (x_i - y_i)^2/4$$

- Then

$$\mathbb{E}(\hat{\sigma}^2) = \sigma^2/2$$

and the MLE is inconsistent as $n \rightarrow \infty$

High-dimensional logistic regression

- The logistic model assumes that the probability of a case conditional on the covariates is given by

$$P(y_i = 1|x_i) = \rho(x_i^t \beta)$$

where $\rho(z) = e^z/(1 + e^z)$ is the sigmoidal function

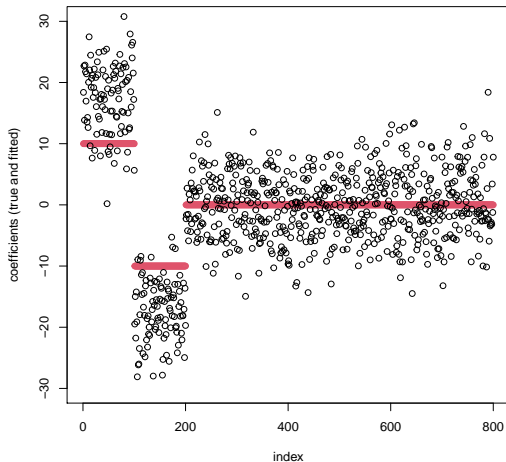
- When p is fixed and $n \rightarrow \infty$, the MLE obeys

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(0, \mathcal{I}_{\beta}^{-1})$$

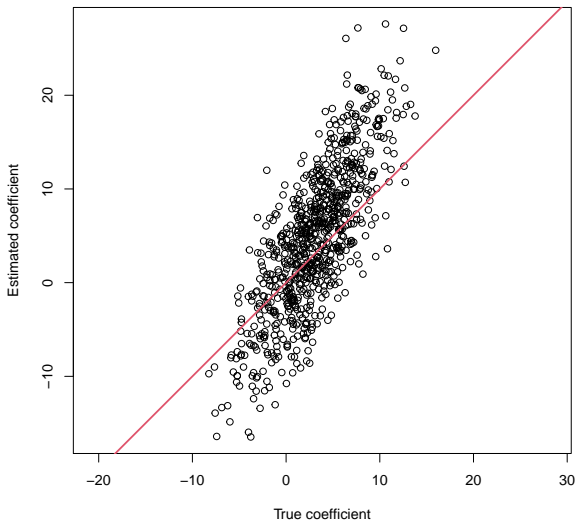
where \mathcal{I}_{β} is the $p \times p$ Fisher information matrix evaluated at the true β

- Does this approximation holds in high-dimensional settings where p is not vanishingly small compared with n ?
- We set $n = 4000$ and $p = 800$, so that $p/n = 1/5$, with $x_{ij} \stackrel{\text{iid}}{\sim} N(0, 1/n)$

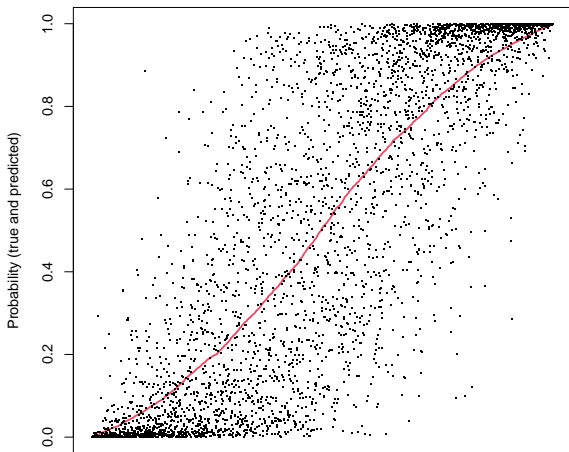
Unbiasedness?



$$\beta_1 = \dots = \beta_{100} = 10, \beta_{101} = \dots = \beta_{200} = -10, \beta_{201} = \dots = \beta_{800} = 0$$



$$\beta_i \stackrel{\text{iid}}{\sim} N(3, 16)$$



$$\beta_i \stackrel{\text{iid}}{\sim} N(3, 16)$$

Distribution of the LRT?

- Testing $H_j : \beta_j = 0$ against $\beta_j \neq 0$
- Log-likelihood ratio statistic

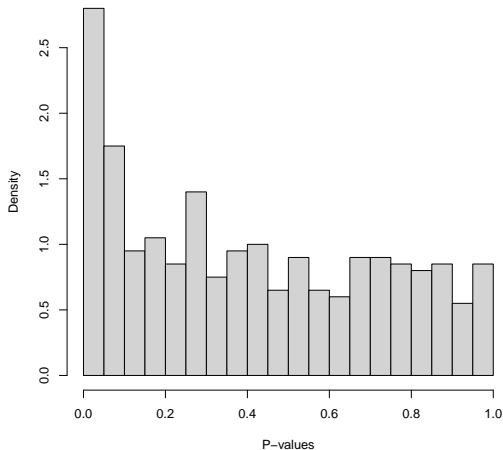
$$\Lambda_j = \ell(\hat{\beta}) - \ell(\hat{\beta}_{-j})$$

- Wilks Theorem: with p fixed and $n \rightarrow \infty$, the test has asymptotic null distribution

$$2\Lambda_j \xrightarrow{d} \chi_1^2$$

- If the χ^2 approximation were true, then we would expect to observe uniformly distributed null p -values

Distribution of the null p -values?



Half of β are i.i.d. $N(7, 1)$, the other half o

Linear discriminant analysis in high-dimensions

Classification problem

- Let's turn to the classification problem involving the allocation of the observed unit x to one of two classes A and B
- For a Bayesian analysis suppose that the prior probabilities are $\pi_A \equiv P(y = A)$ and $\pi_B \equiv P(y = B)$ with $\pi_A + \pi_B = 1$. Then the posterior probabilities satisfy

$$\frac{P(y = B|x)}{P(y = A|x)} = \frac{\pi_B f_B(x)}{\pi_A f_A(x)}$$

giving the class with the higher posterior probability

- As a special case, suppose that the two classes are distributed as multivariate Gaussians $x_A \sim N(\mu_A, I_p)$ and $x_B \sim N(\mu_B, I_p)$, with $\pi_A = \pi_B = 1/2$

Optimal decision

- The optimal decision rule is to threshold the log-likelihood ratio

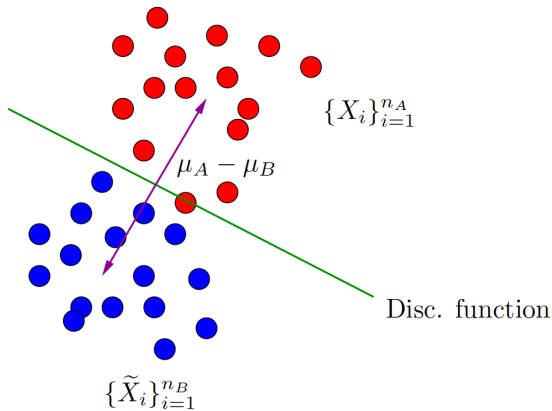
$$\Psi(x) = \langle \mu_A - \mu_B, \left(x - \frac{\mu_A + \mu_B}{2} \right) \rangle$$

where $\langle x, z \rangle = x^t z = \sum_{j=1}^p x_j z_j$ denotes the Euclidean inner product in \mathbb{R}^p

- If $\Psi(x) > 0$ then classify A , otherwise B
- Error probability of the optimal rule:

$$\text{Err}(\Psi) = \frac{1}{2} \text{P}(\Psi(x_A) < 0) + \frac{1}{2} \text{P}(\Psi(x_B) \geq 0) = \Phi\left(-\frac{\gamma}{2}\right)$$

where $\gamma = \|\mu_A - \mu_B\|$ and Φ is the cdf of a standard normal variable



$$\langle \mu_A - \mu_B, \left(x - \frac{\mu_A + \mu_B}{2} \right) \rangle = 0$$

source: Wainwright

Linear Discriminant Analysis

- Fisher's LDA: uses the plug-in principle based on n_A samples from class A and n_B samples from class B

$$\hat{\Psi}(x) = \langle \hat{\mu}_A - \hat{\mu}_B, x - \frac{\hat{\mu}_A + \hat{\mu}_B}{2} \rangle$$

- Error probability of LDA (is itself a random variable)

$$\text{Err}(\hat{\Psi}) = \frac{1}{2}P(\hat{\Psi}(x_A) < 0) + \frac{1}{2}P(\hat{\Psi}(x_B) \geq 0)$$

- Classical theory: if $(n_A, n_B) \rightarrow \infty$ and p remains fixed, then $\hat{\mu}_A \xrightarrow{prob.} \mu_A$, $\hat{\mu}_B \xrightarrow{prob.} \mu_B$ and the asymptotic error probability is $\text{Err}(\hat{\Psi}) \xrightarrow{prob.} \text{Err}(\Psi) = \Phi(-\gamma/2)$

- If $x \sim N(\mu, I_p)$, then

$$\begin{aligned}\hat{\Psi}(x) &= \langle \hat{\mu}_A - \hat{\mu}_B, x - \frac{\hat{\mu}_A + \hat{\mu}_B}{2} \rangle \\ &= \hat{d}^t(x - \hat{m}) \sim N(\hat{d}^t(\mu - \hat{m}), \hat{d}^t \hat{d})\end{aligned}$$

where $\hat{d} = \hat{\mu}_A - \hat{\mu}_B$ and $\hat{m} = \frac{\hat{\mu}_A + \hat{\mu}_B}{2}$

High-Dimensional Theory

- What happens if $(n_A, n_B, p) \rightarrow \infty$ with
 - $p/n_A \rightarrow \delta$ with $\delta \geq 0$
 - $p/n_B \rightarrow \delta$
 - $\|\mu_A - \mu_B\|_2 \rightarrow \gamma > 0$
- Kolmogorov (1960) showed that

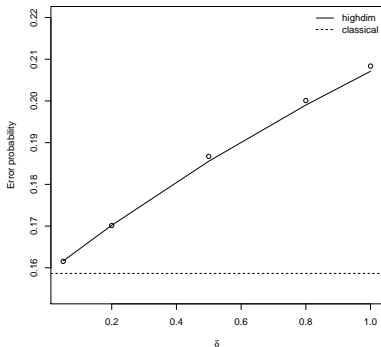
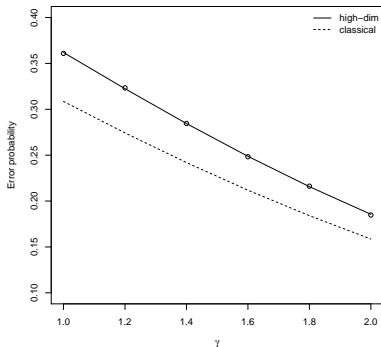
$$\text{Err}(\hat{\Psi}) \xrightarrow{\text{prob.}} \Phi\left(-\frac{\gamma^2}{2\sqrt{\gamma^2 + 2\delta}}\right)$$

- If $p/n \rightarrow 0$, then the asymptotic error probability is $\Phi(-\gamma/2)$ as is predicted by classical theory
- If $p/n \rightarrow \delta > 0$, then the asymptotic error probability is strictly larger than $\Phi(-\gamma/2)$

The error probability of $\hat{\Phi}$, for the finite triple

$$(p, n_A, n_B) = (400, 800, 800)$$

is better described by the classical $\Phi(-\gamma/2)$, or the high-dimensional analog $\Phi(-\gamma^2/(2\sqrt{\gamma^2 + 2\delta}))$?



circles correspond to the empirical error probabilities, averaged over 10 trials

What can help us in high dimensions?

- An important fact is that high-dimensional phenomena are unavoidable
- If the ratio p/n stays bounded strictly above zero, then it is not possible to achieve the optimal classification rate
- Our only hope is that the data is endowed with some form of *low-dimensional structure*

- What is the underlying cause of the inaccuracy of the prediction for the LDA in high-dimensions?
- The squared Euclidean error

$$\|\hat{\mu} - \mu\|^2 = \sum_{j=1}^p (\hat{\mu}_j - \mu_j)^2$$

concentrates sharply around p/n , i.e. for $t \in (0, 1)$

$$\mathbb{P} \left(\left| \|\hat{\mu} - \mu\|^2 - \frac{m}{n} \right| \geq \frac{p}{n} t \right) = \mathbb{P} \left(\left| \frac{1}{p} \sum_{j=1}^p Z_j^2 - 1 \right| \geq t \right) \leq 2e^{-\frac{pt^2}{8}}$$

where $Z_j = \sqrt{n}(\hat{\mu}_j - \mu_j) \sim N(0, 1)$; for the upper bound see Wainwright (2019), Example 2.11

Sparsity

- Suppose that the p -vector μ is *sparse*, with only s of its p entries being nonzero, for some sparsity parameter $s \ll p$
- If sparsity holds, we can obtain a better estimator by thresholding the sample means

$$\tilde{\mu}_j = \hat{\mu}_j \mathbb{1}\{|\hat{\mu}_j| > \lambda\}$$

where

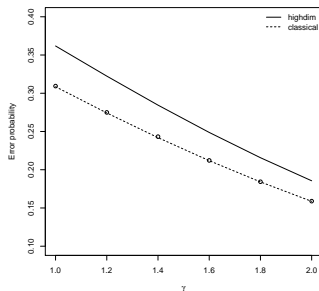
$$\lambda = \sqrt{\frac{2 \log p}{n}}$$

Thresholded mean

Suppose to replace $\hat{\mu}$ by the thresholded mean $\tilde{\mu}$, then

$$\tilde{\Psi}(x) = \langle \tilde{\mu}_A - \tilde{\mu}_B, x - \frac{\tilde{\mu}_A + \tilde{\mu}_B}{2} \rangle$$

approaches the optimal $\text{Err}(\Psi)$ if $\log \binom{p}{s} / n \rightarrow 0$. For $s = 5$:



circles correspond to the empirical error probabilities, averaged over 10 trials