

Problem 1

calculating the power of a chi-squared

$$X_1, \dots, X_n \sim N(\mu, \sigma)$$

$$H_0: \sigma \leq 1$$

$$H_A: \sigma > 1$$

$$\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 > \frac{1}{n-1} \chi_{\alpha, n-1}^2$$

\Rightarrow estimate the power of this test
if $\sigma = 2$, $n = 10$, $\alpha = .05$

\Rightarrow generate $(m \times n)$ array of
 $X =$ normal random variable
with $\mu = 0$, $\sigma = 2$

\Rightarrow calculate critical region

$$C = \frac{1}{n-1} \text{chi2.ppf}(1-\alpha, n-1)$$

\Rightarrow for eq sum in X ,
calculate S^2 `np.var(X, axis=1, ddof=1)`
 $S^2 =$

`np.mean(S^2 > C)`

Problem 3

$$H_0: \mu \leq 2.5$$

$$H_A: \mu > 2.5$$

\hookrightarrow calculating an upper confidence interval and checking if
 $2.5 \in \text{interval}$

\hookrightarrow in interval, fail to reject

\hookrightarrow not in, reject

upper confidence:

(lower bound, ∞)

lower confidence

($-\infty$, upper bound)

$S_2 = \text{cray}$ of sample variances

	4.7	T
$\left(\frac{1}{n-1}\right) \chi^2_{.05, 9}$	3.2	T
	2.8	F
\downarrow	1.7	F
3.0	2.9	F
	3.4	T

$$S_2 > 3.0$$

confidence interval for the expected value
 $y|x$

$$y_i = \alpha + \beta x_i + \epsilon_i$$

$\hat{\alpha}, \hat{\beta} \rightarrow$ best guesses

confidence $\hat{\alpha} + \hat{\beta} x_i$

prediction interval accounts for the variation amongst individuals at x_i

bootstrapping where we resample the prior
on the predictor

$$\hat{\alpha} + \hat{\beta} x_i \rightarrow \text{confidence for}$$

$$\alpha + \beta x_i$$

prediction interval
given x_0

① generate m bs samples
 $\hat{\alpha}_j, \hat{\beta}_j$ for $j=1, \dots, m$

② generate a prediction

$$\hat{y}_j = \hat{\alpha}_j + \hat{\beta}_j \cdot x_0 + \hat{\epsilon}_{0j} \Rightarrow$$

$$\hat{\epsilon}_{ij} = y_{ij} - \hat{\alpha}_j - \hat{\beta}_j x_{ij}$$

↳ for each of these j
samples, I randomly
choose one of these
 $\hat{\epsilon}_{ij}$, and I add it

$$\text{Var}(X) = E[(X - \mu_x)^2]$$

$$\text{Cov}(X, Y) = E[(X - \mu_x) \cdot (Y - \mu_y)]$$

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$$

$$SD(X) SD(Y)$$

$$(X_1, \dots, X_k)$$

Covariance matrix Σ
 $k \times k$ matrix

$$\Sigma_{ii} = \text{Var}(X_i) = \sigma_i^2$$

$$\Sigma_{ji} = \Sigma_{ij} = \text{Cov}(X_i, X_j) = \rho_{ij} \sigma_i \sigma_j$$

\uparrow
 $\text{corr}(X_i, X_j)$
 $\cdot SD(X_i) \cdot SD(X_j)$

Multiple regression

$$\Rightarrow y = \alpha + \beta x_i \quad (x_1, y_1), \dots, (x_n, y_n)$$

what if we have more than 1 predictors?
 — d predictors

$$(x_{11}, x_{12}, \dots, x_{1d}, y_1), (x_{21}, x_{22}, \dots, x_{2d}, y_2),$$

$$\dots (x_{n1}, x_{n2}, \dots, x_{nd}, y_n)$$

$$Y_i = \alpha + \beta_1 X_{i1} + \dots + \beta_d X_{id} + \epsilon_i$$

$$X_{i1} = 1 \text{ for every } i=1, \dots, N \quad \epsilon_i \sim N(0, \sigma)$$

$$= \beta_1 X_{i1} + \dots + \beta_d X_{id} + \epsilon_i \quad \textcircled{=}$$

X is $n \times d$ matrix

Y is a $n \times 1$ column vector

β is a $d \times 1$ column vector

ϵ is a $n \times 1$ column vector

$$Y = X\beta + \epsilon \quad \triangle \text{ design matrix}$$

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}$$

$=$

$$\begin{bmatrix} X_{11} & X_{12} & \dots & X_{1d} \\ X_{21} & & & \\ \vdots & & & \\ \vdots & & & \\ \vdots & & & \\ X_{n1} & \dots & \dots & X_{nd} \end{bmatrix}$$

$$\begin{bmatrix} \beta_1 \\ \vdots \\ \vdots \\ \beta_d \end{bmatrix}$$

$$+ \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$\text{minimize}_{\beta} \sum_{i=1}^n (y_i - \beta_1 x_{i1} - \dots - \beta_d x_{id})^2$$

$$\text{minimize}_{\beta} (Y - X\beta)^T (Y - X\beta)$$

$$\frac{d}{d\beta} (Y - X\beta)^T (Y - X\beta) = 0$$

$$X^T (Y - X\beta) = 0$$

$$X^T Y = X^T X \beta$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

$X^T X \Rightarrow$ has to be invertible

columns of X need to be linearly independent

$\hat{\beta}$ is a random variable
 d -dimensional normal random variable

$$\hat{\beta} \sim N(\beta, \sigma^2 (X^T X)^{-1})$$

↳ this can be problematic if there are columns X_i and X_j

columns of X have standardised so that are highly correlated

$$X^T X = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$(X^T X)^{-1} = \frac{1}{1} \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\hat{\beta} \sim N(\beta, \sigma^2 I)$$

$$X^T X = \begin{pmatrix} 1 & .95 \\ .95 & 1 \end{pmatrix}$$

$$(X^T X)^{-1} = \frac{1}{1 - .95^2} \begin{pmatrix} 1 & -.95 \\ -.95 & 1 \end{pmatrix}$$

$$\begin{pmatrix} 1 & -.95 \end{pmatrix}$$

$$\approx 11 \quad \begin{pmatrix} \dots & \dots \\ -0.95 & 1 \end{pmatrix}$$

$$X^T X = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$y_i = 0.5 x_{i1} + 0.5 x_{i2} + \epsilon_i$$

$$\hat{\beta}_1 = 0.5 \quad \hat{\beta}_2 = 0.49$$

$$X^T X = \begin{pmatrix} 1.95 & \\ .95 & 1 \end{pmatrix}$$

$$\hat{\beta}_1 = 3.1$$

$$\hat{\beta}_2 = -1.8$$

↳ tends to really create issues with out-of-sample predictive

d is relatively big \rightarrow spurious correlation or going to exist between the columns

idea: penalize the magnitude of the regression coefficients

$$\underset{\beta}{\text{minimize}} \sum_{i=1}^n (y_i - \beta_1 x_{i1} - \dots - \beta_d x_{id})^2 + \lambda \sum_{j=1}^d \beta_j^2$$

λ is what we call a "hyperparameter"

↳ not a parameter that determines a dist'n

$$\min_{\beta} (Y - X\beta)^T (Y - X\beta) + \lambda \beta^T \beta$$

↳ for this to make sense, we need to standardize the columns of $X \Rightarrow$ all have unit $\sum_{i=1}^n x_{ij}^2 = 1$

$$\hat{\beta}_{\lambda} = (X^T X + \lambda \mathbf{I}_d)^{-1} X^T Y$$

$$X^T X = \begin{pmatrix} 1 & .95 \\ .95 & 1 \end{pmatrix}$$

$$(X^T X)^{-1} \approx \mathbf{I} = \begin{pmatrix} 1 & .95 \\ -.95 & 1 \end{pmatrix}$$

$$\lambda = 1$$

$$(X^T X + \mathbf{I}) = \begin{pmatrix} 2 & .95 \\ .95 & 2 \end{pmatrix}$$

$$(X^T X + \mathbf{I})^{-1} = \frac{1}{3.1} \begin{pmatrix} 2 & -.95 \\ -.95 & 2 \end{pmatrix}$$

reduced variance of β_i by about a factor of 15 :

What is the downside of this?

$$\hat{\beta}_\lambda = (X^T X + \lambda I)^{-1} X^T Y$$

$$(X^T X + \lambda I) \hat{\beta}_\lambda = X^T Y$$

$$(X^T X)^{-1} (X^T X + \lambda I) \hat{\beta}_\lambda = \underbrace{(X^T X)^{-1} X^T Y}_{\hat{\beta}_0} = \hat{\beta}_0$$

$$\underline{\underline{(I + (X^T X)^{-1}) \hat{\beta}_\lambda = \hat{\beta}_0}}$$

$$\hat{\beta}_\lambda = (I + (X^T X)^{-1})^{-1} \hat{\beta}_0$$

↳ this is a biased estimator of β

How do we pick λ ?

↳ trying to evaluate which λ will generalize best out-of-sample

