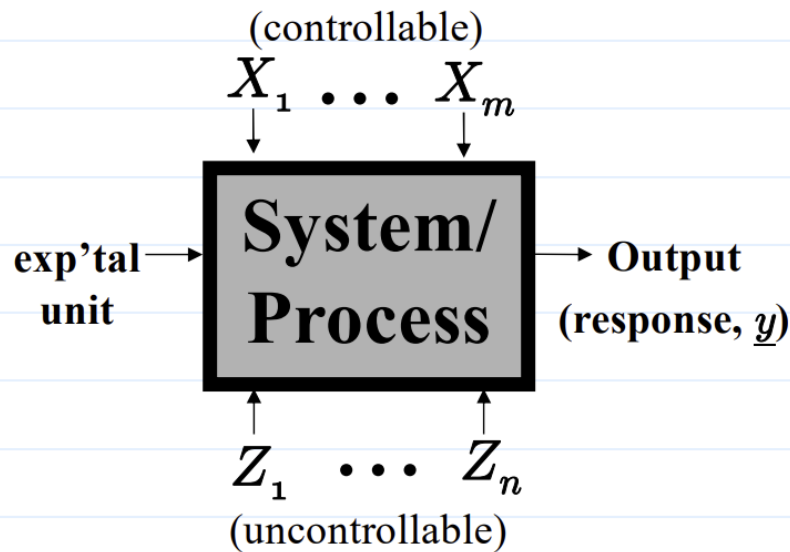


(U4284) Python程式設計 Linear Model



Speaker: 吳淳硯



Some Terminologies

- A response/output/dependent variable Y is modeled or explained by p effects/functions of m predictor/input/independent variables X_1, \dots, X_m

$$Y = \underbrace{\sum_{k=0}^p \beta_k g_k(X_1, \dots, X_m)}_{\text{deterministic part}} + \underbrace{\varepsilon}_{\text{random part}}$$

$$E(Y|X_1, \dots, X_m) = \sum_{k=0}^p \beta_k g_k, \quad \text{Var}(Y|X_1, \dots, X_m) = \sigma^2$$

- When $p = 1$, simple linear regression; $p > 1$ multiple regression
- X_1, \dots, X_m : continuous, discrete (quantitative), categorical (qualitative)
 - All Quantitative \Rightarrow Multiple regression
 - Quantitative + Qualitative \Rightarrow ANCOVA
 - All Qualitative \Rightarrow ANOVA
- More than one Y , multivariate regression.

Rationale

- A general model for the relationship between Y and X_1, \dots, X_p is

$$Y = f(X_1, \dots, X_p) + \varepsilon$$

where f is unknown and arbitrary.

- Local approximation of function f may be achievable by linear model

$$\begin{aligned} f(X_1, \dots, X_p) &\approx f(a_1, \dots, a_p) + \sum_{k=1}^p \left[\frac{\partial f}{\partial X_k}(a_1, \dots, a_p) \right] (X_k - a_k) \\ &= \beta_0 + \sum_{k=1}^p \beta_k X_k \end{aligned}$$

- Predictors can be transformed and combined in any way, linear models are actually very flexible.

Ex. $x_1 x_2, e^{x_1 + x_2}, \dots$

- Matrix Form,

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{1,1} & \cdots & x_{p,1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1,n} & \cdots & x_{p,n} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

MME (Method of Moment Estimation)

- Law of Large number

$$\frac{g(X_1) + \cdots + g(X_n)}{n} \rightarrow E(g(X))$$

- X is **exogenous** for β if $E(X\varepsilon) = 0$. While a regressor X is exogenous if it is determined separately from Y .
- Since $E(\varepsilon) = 0$ and $E(X\varepsilon) = 0$

$$E(Y - \beta_0 - \beta_1 X) = 0 \Rightarrow \frac{\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)}{n} = 0$$

$$E(X[Y - \beta_0 - \beta_1 X]) = 0 \Rightarrow \frac{\sum_{i=1}^n x_i (y_i - \beta_0 - \beta_1 x_i)}{n} = 0$$

Obtain the normal equation

$$\begin{cases} \sum_{i=1}^n y_i = n\beta_0 + \beta_1 \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i y_i = \beta_0 \sum_{i=1}^n x_i + \beta_1 \sum_{i=1}^n x_i^2 \end{cases}$$

OLS (Ordinary Least Squared)

- Consider a objective function is given by

$$S(\boldsymbol{\beta}) = \boldsymbol{\varepsilon}^\top \boldsymbol{\varepsilon} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

find the coefficient $\boldsymbol{\beta}$ which fit the equations "best", in the sense of solving the quadratic minimization problem

$$\arg \min_{\boldsymbol{\beta}} S(\boldsymbol{\beta}) \Leftrightarrow \frac{\partial S}{\partial \boldsymbol{\beta}} = 0$$

- By calculus, $\hat{\boldsymbol{\beta}}$ is the solution of

$$\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^\top \mathbf{y}$$

assume $\mathbf{X}^\top \mathbf{X}$ is nonsingular then

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \boldsymbol{\beta} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\varepsilon}$$

$$\Rightarrow \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{H} \mathbf{y}$$

- Does $\hat{\boldsymbol{\beta}}$ possess any optimal properties?

$$E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}, \quad \text{Var}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}$$

and Gauss-Markov Theorem, for any $\check{\boldsymbol{\beta}} = \mathbf{A} \mathbf{y}$ and $\|\mathbf{a}\| = 1$

$$\mathbf{a}^\top (\text{Var}(\check{\boldsymbol{\beta}}) - \text{Var}(\hat{\boldsymbol{\beta}})) \mathbf{a} \geq 0$$

What is likelihood?



- The parallel universe analogy for likelihood:

$$\mathcal{L}(\boldsymbol{\theta}|X_1, \dots, X_n) = \prod_{k=1}^n f(X_k|\boldsymbol{\theta})$$

where $\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)$

- The likelihood \mathcal{L} tells us how probable it is that this particular universe generated the observed data.
- The higher the likelihood, the more plausible the universe is for explaining your observations.

MLE (Maximum Likelihood Estimation)

- Goal of MLE

find the best parameter that make the observed data most probable.

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta} \in \Theta} \mathcal{L}(\boldsymbol{\theta} | \mathbf{x})$$

- Consistency

$$\hat{\boldsymbol{\theta}} \rightarrow \boldsymbol{\theta} \text{ in probability}$$

- Invariance

If $\hat{\boldsymbol{\theta}}$ is MLE of $\boldsymbol{\theta}$, then $\tau(\hat{\boldsymbol{\theta}})$ is the MLE of $\tau(\boldsymbol{\theta})$, where $\tau(\cdot)$ is a well-defined function.

- Asymptotic Efficiency

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \sim N(0, \mathcal{I}^{-1})$$

where \mathcal{I} is fisher information matrix

$$\begin{aligned} \mathcal{I} &= E \left(- \frac{\partial^2 \ln f(X_k | \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \right) \\ &= E \left[\left(\frac{\partial \ln f(X_k | \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right) \left(\frac{\partial \ln f(X_k | \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)^\top \right] \end{aligned}$$

What is a Gradient?

- The directional derivative measures the rate at which a function changes in a particular direction at a given point, the directional derivative of function f along a unit vector $\vec{u} = (u_1, \dots, u_n)$.

$$\mathcal{D}_{\vec{u}}f(\mathbf{x}) = \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\vec{u}) - f(\mathbf{x})}{h} = \nabla f(\mathbf{x}) \cdot \vec{u}$$

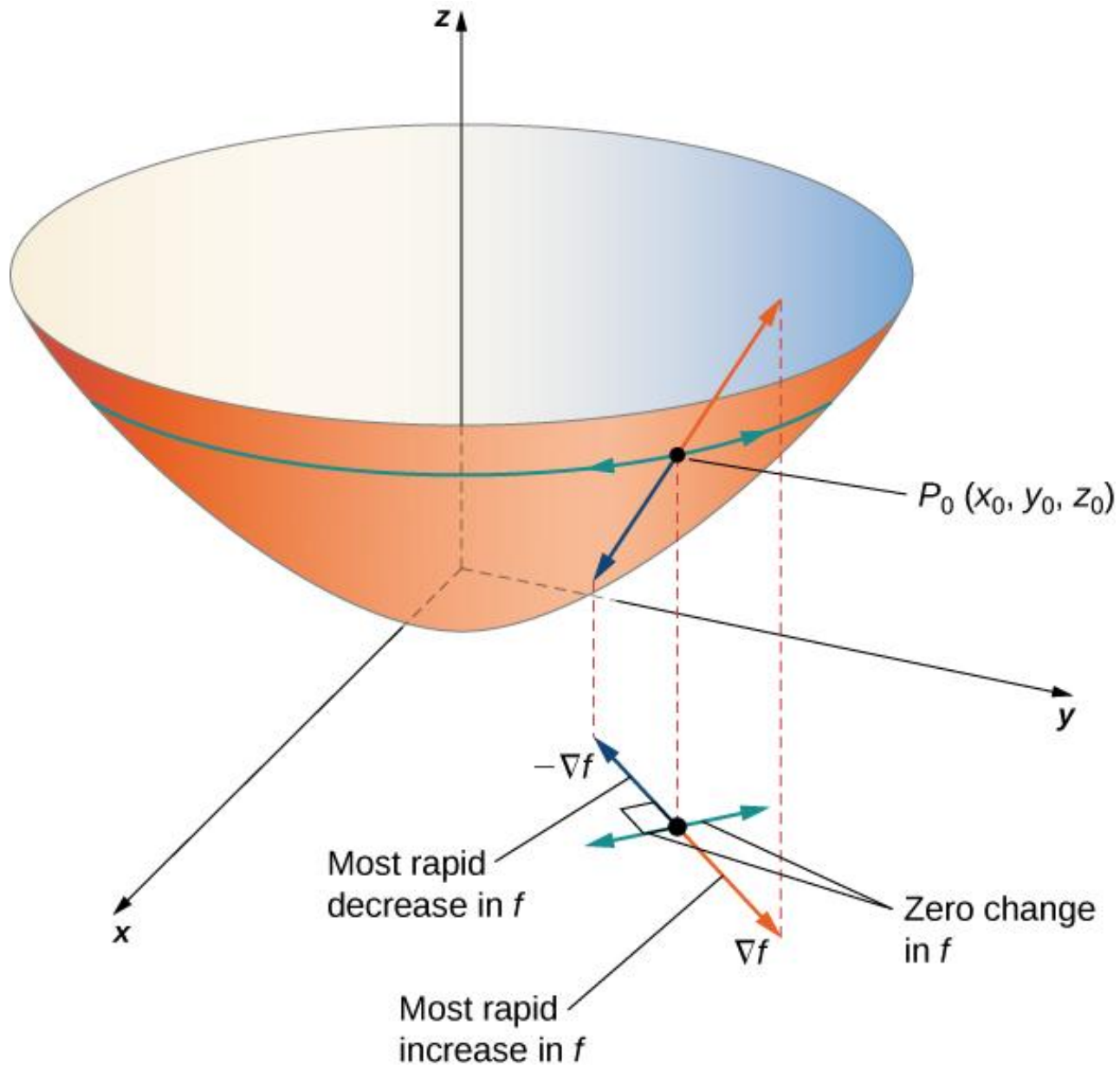
where $\nabla f(\mathbf{x})$ is gradient of f

$$\nabla f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_1}(\mathbf{x}) \quad \dots \quad \frac{\partial f}{\partial x_n}(\mathbf{x}) \right]^T$$

- Properties of the Gradient

- $\begin{vmatrix} \nabla f(\mathbf{x}) \\ -\nabla f(\mathbf{x}) \end{vmatrix}$ points in the direction of fastest $\begin{vmatrix} \text{increase} \\ \text{decrease} \end{vmatrix}$ at \mathbf{x}
 - $\begin{vmatrix} |\nabla f(\mathbf{x})| \\ -|\nabla f(\mathbf{x})| \end{vmatrix}$ is equal to $\begin{vmatrix} \text{maximum} \\ \text{minimum} \end{vmatrix} \mathcal{D}_{\vec{u}}f(\mathbf{x})$

Visual of Gradient



Gradient Descent

- Given α is the **learning rate**, consider the objective function

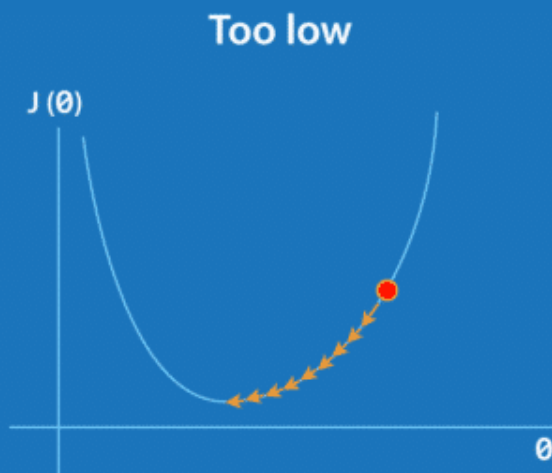
$$\mathcal{S}(\boldsymbol{\beta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

the gradient of $\mathcal{S}(\boldsymbol{\beta})$ is

$$\nabla \mathcal{S}(\boldsymbol{\beta}) = 2\mathbf{X}^\top (\mathbf{X}\boldsymbol{\beta} - \mathbf{y})$$

we have the iterative form

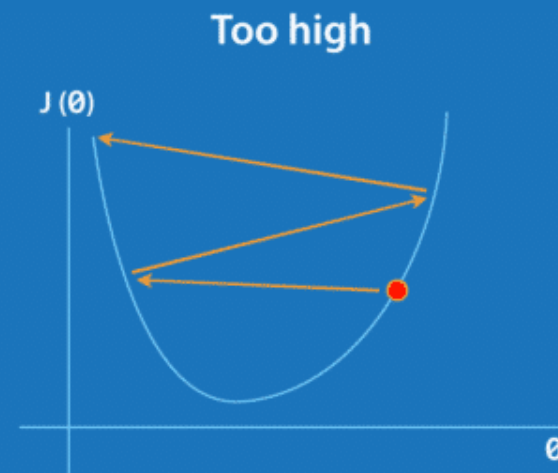
$$\boldsymbol{\beta}^{(i)} = \boldsymbol{\beta}^{(i-1)} - \alpha \nabla \mathcal{S}(\boldsymbol{\beta})$$



A small learning rate requires many updates before reaching the minimum point



The optimal learning rate swiftly reaches the minimum point



Too large of a learning rate causes drastic updates which lead to divergent behaviors

GLM

- Three component of GLM (Generalized Linear Model)
 - Random component $f(\mathbf{y})$
 - Linear predictor $\mathbf{X}\boldsymbol{\beta}$
 - Link function g such that $g(E(\mathbf{y})) = \mathbf{X}\boldsymbol{\beta}$

- Exponential Dispersion Family

$$f(y|\theta, \phi) = \exp\left(\frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi)\right)$$

where

$$\mu = E(Y) = b'(\theta), \quad \text{Var}(Y) = b''(\theta)a(\phi)$$

the parameter θ is called the natural parameter, and ϕ is called the dispersion parameter.

Random Component	Link Function	Model
Normal	Identity	Regression Analysis of variance
Exponential family	Any	Generalized linear model
Binomial	Logit	Logistic regression
Multinomial	Generalized logits	Multinomial response
Poisson	Log	Loglinear

Constraint Least Square

- An intuitively appealing method to estimate a constrained linear model is to minimize the least squared criterion subject to the constraint $\mathcal{A}\beta = c$, by Lagrange Multiplier

$$\mathcal{S}(\beta, \lambda) = \frac{1}{2}(\mathbf{y} - \mathbf{X}\beta)^\top(\mathbf{y} - \mathbf{X}\beta) + \lambda^\top(\mathcal{A}\beta - c)$$

1st order condition

$$\frac{\partial \mathcal{S}}{\partial \beta} = -\mathbf{X}^\top \mathbf{y} + \mathbf{X}^\top \mathbf{X}\beta + \mathcal{A}^\top \lambda = \mathbf{0}$$

$$\frac{\partial \mathcal{S}}{\partial \lambda} = \mathcal{A}\beta - c = \mathbf{0}$$

multiplying by $\mathcal{A}(\mathbf{X}^\top \mathbf{X})^{-1}$

$$-\mathcal{A}\hat{\beta}_{\text{ols}} + \mathcal{A}\beta + \mathcal{A}(\mathbf{X}^\top \mathbf{X})^{-1}\mathcal{A}^\top \lambda = \mathbf{0}$$

$$\Rightarrow \hat{\lambda} = [\mathcal{A}(\mathbf{X}^\top \mathbf{X})^{-1}\mathcal{A}^\top]^{-1}(\mathcal{A}\hat{\beta}_{\text{ols}} - c)$$

then

$$\hat{\beta}_{\text{cls}} = \hat{\beta}_{\text{ols}} - (\mathbf{X}^\top \mathbf{X})^{-1}\mathcal{A}^\top[\mathcal{A}(\mathbf{X}^\top \mathbf{X})^{-1}\mathcal{A}^\top]^{-1}(\mathcal{A}\hat{\beta}_{\text{ols}} - c)$$

Orthogonal Projection

- Let Ω be subspace of \mathbb{R}^n and let \mathbf{y} be a vector in \mathbb{R}^n . The expression

$$\mathbf{y} = \mathbf{y}_{\Omega} + \mathbf{y}_{\Omega^{\perp}}$$

where

$$\mathbf{y}_{\Omega} \in \Omega, \quad \mathbf{y}_{\Omega^{\perp}} \in \Omega^{\perp}$$

is called orthogonal decomposition of \mathbf{y} w.r.t Ω .

- The distance from \mathbf{y} to Ω

$$\|\mathbf{y}_{\Omega^{\perp}}\| = \|\mathbf{y} - \mathbf{y}_{\Omega}\|$$

- Let A be $m \times n$ matrix with linearly independent columns and let $\Omega = C(\mathbf{X})$. Then the $n \times n$ matrix $\mathbf{X}^{\top} \mathbf{X}$ is invertible, and for all vectors $\mathbf{x} \in \mathbb{R}^m$, we have

$$\mathbf{y}_{\Omega} = \mathbf{X}(\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{y} = \mathbf{P} \mathbf{y}$$

and \mathbf{P} is the orthogonal projection matrix with

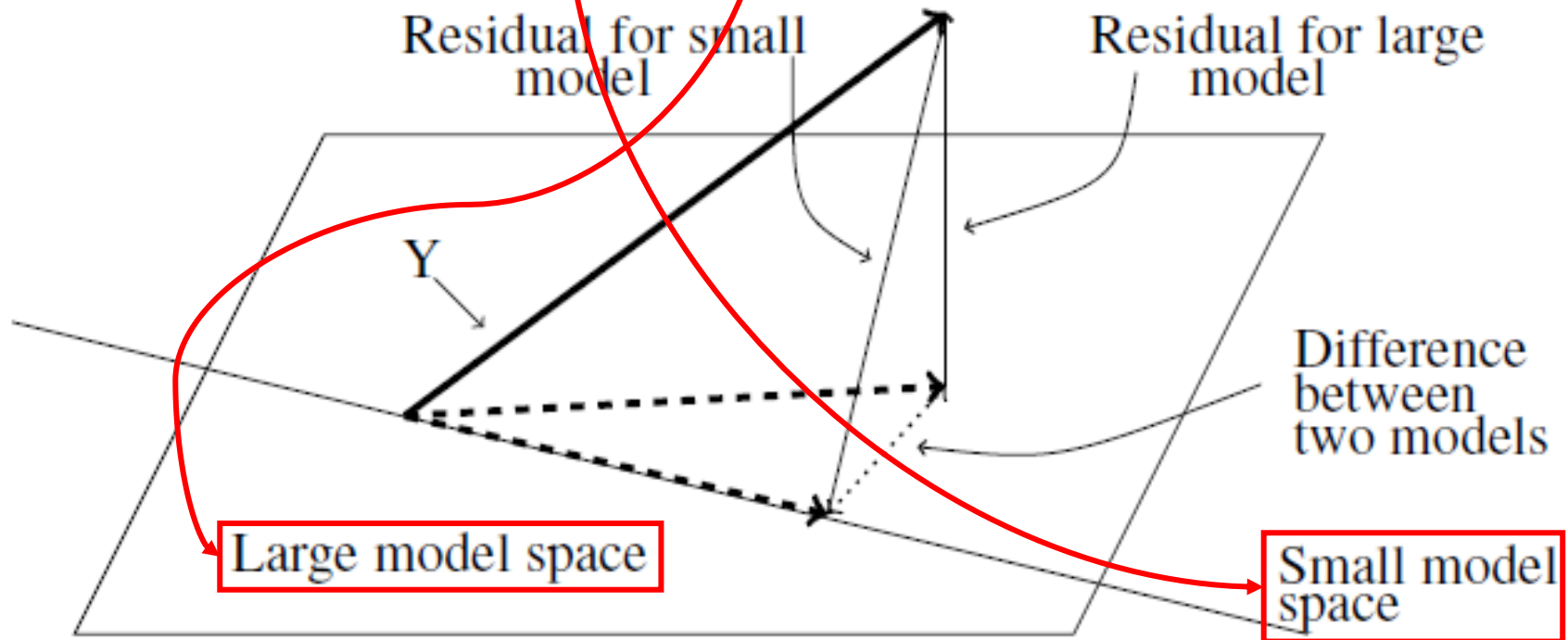
$$\mathbf{P}^2 = \mathbf{P} \text{ (idempotent)}$$

and

$$\mathbf{P}^{\top} = \mathbf{P} \text{ (symmetry)}$$

Hypothesis Testing - 1

- Geometric view of $H_0: \omega$ vs $H_1: \Omega \setminus \omega$



residual for large model $\triangleq \hat{\varepsilon}_{\Omega}$

residual for small model $\triangleq \hat{\varepsilon}_{\omega}$

difference between two models $\triangleq \hat{\varepsilon}_{\omega} - \hat{\varepsilon}_{\Omega}$

Hypothesis Testing - 2

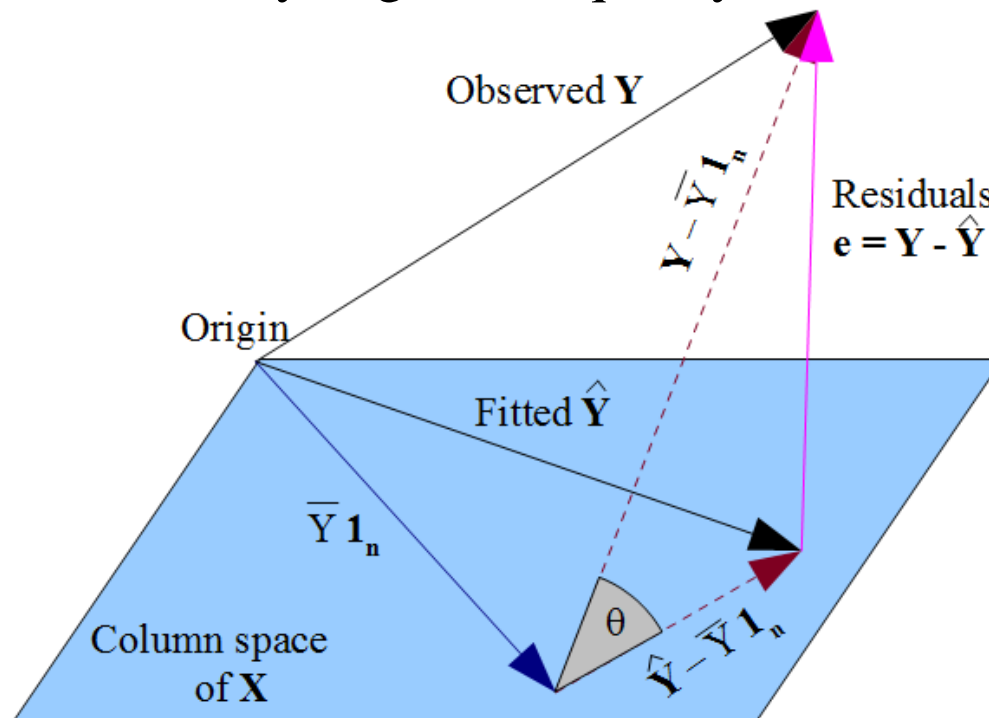
- Analysis of Variance (ANOVA)

$$\begin{aligned} SST &= \mathbf{y}^\top (\mathbf{I} - \mathbf{J}) \mathbf{y} = \mathbf{y}^\top (\mathbf{I} - \mathbf{H}) \mathbf{y} + \mathbf{y}^\top (\mathbf{H} - \mathbf{J}) \mathbf{y} \\ &= SSE + SSR \end{aligned}$$

where

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top, \quad \mathbf{J} = \frac{\mathbf{1}\mathbf{1}^\top}{n}$$

- Actually, ANOVA is a Pythagorean equality, as illustrate below



Hypothesis Testing - 3

- The length of the difference between two ω and Ω is

$$\|\hat{\varepsilon}_\omega - \hat{\varepsilon}_\Omega\|^2 = \|\hat{\varepsilon}_\omega\|^2 - \|\hat{\varepsilon}_\Omega\|^2 = \text{RSS}_\omega - \text{RSS}_\Omega$$

suggest divide with length of residual of Ω for scaling

$$\tan^2(\theta) = \frac{\|\hat{\varepsilon}_\omega - \hat{\varepsilon}_\Omega\|^2}{\|\hat{\varepsilon}_\Omega\|^2} = \frac{\text{RSS}_\omega - \text{RSS}_\Omega}{\text{RSS}_\Omega}$$

- Why not divided by RSS_ω ?

Orthogonality and central χ^2 under ω .

- Under H_0 , $\hat{\varepsilon}_\omega - \hat{\varepsilon}_\Omega = (\mathbf{H}_\Omega - \mathbf{H}_\omega)\mathbf{y} \sim N(0, (\mathbf{H}_\Omega - \mathbf{H}_\omega)\sigma^2)$

$$\text{RSS}_\omega - \text{RSS}_\Omega \perp \text{RSS}_\Omega$$

$$(\text{RSS}_\omega - \text{RSS}_\Omega)/\sigma^2 \sim \chi_{df_\omega - df_\Omega}^2, \quad \text{RSS}_\Omega/\sigma^2 \sim \chi_{df_\Omega}^2$$

where $df_\omega = \dim(\omega^\perp)$ and $df_\Omega = \dim(\Omega^\perp)$.

- General form

$$F = \frac{(\text{RSS}_\omega - \text{RSS}_\Omega)/(df_\omega - df_\Omega)}{\text{RSS}_\Omega/df_\Omega} \sim F_{df_\omega - df_\Omega, df_\Omega}$$

Hypothesis Testing - 4

- Test of all predictors

$$H_0: \beta_1 = \cdots = \beta_p$$

$$\Omega: y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \varepsilon$$

$$\omega: y = \beta_0 + \varepsilon$$

- Test one predictor

$$H_0: \beta_k = 0$$

$$\Omega: y = \beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k + \cdots + \beta_p X_p + \varepsilon$$

$$\omega: y = \beta_0 + \beta_1 X_1 + \cdots + \cancel{\beta_k X_k} + \cdots + \beta_p X_p + \varepsilon$$

- Test pair of predictors

$$H_0: \beta_i = \beta_j = 0$$

$$\Omega: y = \beta_0 + \beta_1 X_1 + \cdots + \beta_i X_i + \cdots + \beta_j X_j + \cdots + \beta_p X_p + \varepsilon$$

$$\omega: y = \beta_0 + \beta_1 X_1 + \cdots + \cancel{\beta_i X_i} + \cdots + \cancel{\beta_j X_j} + \cdots + \beta_p X_p + \varepsilon$$

- Test a subspace/subset

$$H_0: \beta_i = \beta_j$$

$$\Omega: y = \beta_0 + \beta_1 X_1 + \cdots + \beta_i X_i + \cdots + \beta_j X_j + \cdots + \beta_p X_p + \varepsilon$$

$$\omega: y = \beta_0 + \beta_1 X_1 + \cdots + \beta_i (X_i + X_j) + \cdots + \cancel{\beta_j X_j} + \cdots + \beta_p X_p + \varepsilon$$

Hypothesis Testing - 5

- Linear Hypothesis Test

$$H_0: \mathbf{R}\boldsymbol{\beta} = \mathbf{c} \text{ vs } H_1: \mathbf{R}\boldsymbol{\beta} \neq \mathbf{c}$$

$$H_0: 2\beta_1 - \beta_2 = \beta_2 - 2\beta_3 + 3\beta_4 = \beta_1 - \beta_4 = 0$$

$$H_0: \begin{bmatrix} 0 & 2 & -1 & 0 & 0 \\ 0 & 0 & 1 & -2 & 3 \\ 0 & 1 & 0 & 0 & -1 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

- Let $\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$

$$F = \frac{(\mathbf{R}\boldsymbol{\beta} - \mathbf{c})^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\boldsymbol{\beta} - \mathbf{c}) / q}{\text{RSS} / (n - k - 1)}$$

- If $H_0: \mathbf{R}\boldsymbol{\beta} = \mathbf{c}$ is false then

$$F \sim F_{q, n-k-1, \lambda}$$

$$\text{where } \lambda = (\mathbf{R}\boldsymbol{\beta} - \mathbf{c})^\top (\mathbf{R}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{R}^\top)^{-1} (\mathbf{R}\boldsymbol{\beta} - \mathbf{c}) / 2\sigma^2$$

- If $H_0: \mathbf{R}\boldsymbol{\beta} = \mathbf{c}$ is true then

$$F \sim F_{q, n-k-1}$$

Goodness-of-fit

- R^2 is a measure of the goodness of fit of a model. In regression, R^2 coefficient of determination is a statistical measure of how well the regression predictions approximate the real data points.

$$\Omega: y = \beta_0 + \beta_1 g_1(\mathbf{x}) + \cdots + \beta_p g_p(\mathbf{x}) + \varepsilon$$

$$\omega: y = \beta_0 + \varepsilon$$

$$R^2 = 1 - \frac{\text{RSS}_\Omega}{\text{RSS}_\omega} = \left(\frac{\text{RSS}_\omega - \text{RSS}_\Omega}{\text{RSS}_\omega} \right) = [\text{Cor}(y, \hat{y})]^2 \in [0, 1]$$

- Interpretation of R^2 : Proportion of total variation in y that can be explained by the effects $g_1, \cdots g_p$
 - Increase in the No. of \mathbf{X} increase the value of R^2 , R^2 alone can't be used as a meaningful comparison of models with very different No. of \mathbf{X} .
 - R^2 does not make any sense if an intercept is not in.
 - R^2 does not indicate whether the model $E(y) = \mathbf{X}\boldsymbol{\beta}$ is correct.

Confidence Interval and region - 1

- An interval/region estimation provides
 - Plausible values for parameter.
 - Uncertainty in parameter estimator.
 - Information about its length and the values it covers may be helpful.
 - Information related to testing.

- Confidence region for $R\beta$ where R is full rank matrix

$$\frac{(R\hat{\beta} - R\beta)^T (R(X^T X)^{-1} R)^{-1} (R\hat{\beta} - R\beta)}{d\hat{\sigma}^2} \leq F_{d,n-p}(\alpha)$$

- Confidence region for β

$$R = I_{p \times p}$$

Confidence region of β_i, β_j

$$R = \begin{bmatrix} 0 & \dots & 0 & \overset{i^{\text{th}}}{\underbrace{1}} & 0 & \dots & \dots & \dots & 0 \\ 0 & \dots & \dots & \dots & 0 & \underbrace{1}_{j^{\text{th}}} & 0 & \dots & 0 \end{bmatrix}$$

Confidence Interval and region - 2

- Confidence interval for prediction of mean response at \mathbf{x}_0

$$\mathbf{x}_0^\top \hat{\boldsymbol{\beta}} \pm t_{n-p} \left(\frac{\alpha}{2} \right) \left\{ \hat{\sigma} \sqrt{\mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0} \right\}$$

- What is the danger of extrapolation?

The fitted model may approximate the true model very badly or wrongly and you don't have any information about it.

- Interpolation: \mathbf{x}_0 lie within the range of \mathbf{X}
- Extrapolation: \mathbf{x}_0 lie outside the range of \mathbf{X}

- C.I. for prediction of future observation at \mathbf{x}_0

$$\mathbf{x}_0^\top \hat{\boldsymbol{\beta}} \pm t_{n-p} \left(\frac{\alpha}{2} \right) \left\{ \hat{\sigma} \sqrt{1 + \mathbf{x}_0^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_0} \right\}$$

Orthogonality

- Consider the two models

- Model 1: $Y = \beta_0 + \beta_1 X_1 + \varepsilon$

- Model 2: $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$

In general, $\hat{\beta}_1$ in the two models are not identical. If fitted model is Model 1 and true model is Model 2 then

$$E(\hat{\beta}_1) = \beta_1 + (X_1^\top X_1)^{-1} X_1^\top X_2 \beta_2$$

- An exception:

$$Y = X\beta + \varepsilon = X_1\beta_1 + X_2\beta_2 + \varepsilon$$

$$X^\top X = \begin{bmatrix} X_1^\top X_1 & X_1^\top X_2 \\ X_2^\top X_1 & X_2^\top X_2 \end{bmatrix}$$

when X_1 and X_2 are orthogonal ($X_1 \perp X_2 \Rightarrow X_1^\top X_2 = 0$).

$$(X^\top X)^{-1} = \begin{bmatrix} (X_1^\top X_1)^{-1} & \mathbf{0} \\ \mathbf{0} & (X_2^\top X_2)^{-1} \end{bmatrix}$$

then the estimation $\hat{\beta}_1$ and $\hat{\beta}_2$ independent

Identifiability

- β is called unidentifiable when $X^T X$ is singular, which means the normal equation $X^T X \beta = X^T y$ has infinite solution.
- Unidentifiable means
 - Insufficient data to estimate the parameters of interest.
 - More parameters than are necessary to model the data.
- When does unidentifiability happen?

Observational data

- Same predictor measured in different scales and both in model.
- $X_1 + X_2 = X_3$ or $X_1 + X_2 + X_3 = c$ and all three in model.
- X is supersaturated. ($p > n$)

Experimental data

$$y_{ij} = \mu_i + \varepsilon_{ij} = \mu + \alpha_i + \varepsilon_{ij}$$

- Over-parameterized: some constraint must be imposed on.

$$\sum_{i=1}^r \alpha_i = 0$$

Interpreting

- What does $\hat{\beta}$ mean?

Some β_i have physical interpretation, especially those from conceptual model. Usually β_i do have such physical interpretation. The sign of $\hat{\beta}$ indicate direction of the relationship between the term and the response.

- A naïve interpretation (causality statement)
 - **A unit increase in X_i will cause an average change of $\hat{\beta}_i$ in Y .**
- An alternative interpretation
 - **A unit increase in X_i with all the other(specified) terms held constant will be associated with an average change $\hat{\beta}_i$ in Y .**
- An interpretation from prediction viewpoint
 - Regrading the parameters and their estimates as fictional quantities, and concentrating on prediction enable a rather cautious interpretation of $\hat{\beta}$

given $(g_{1,0}, \dots, g_{i,0}, \dots, g_{p,0}) \rightarrow \hat{y}_0$
 observe $(g_{1,0}, \dots, g_{i,0} + 1, \dots, g_{p,0}) \rightarrow \hat{y}_0 + \hat{\beta}_i$

Causality \neq Association

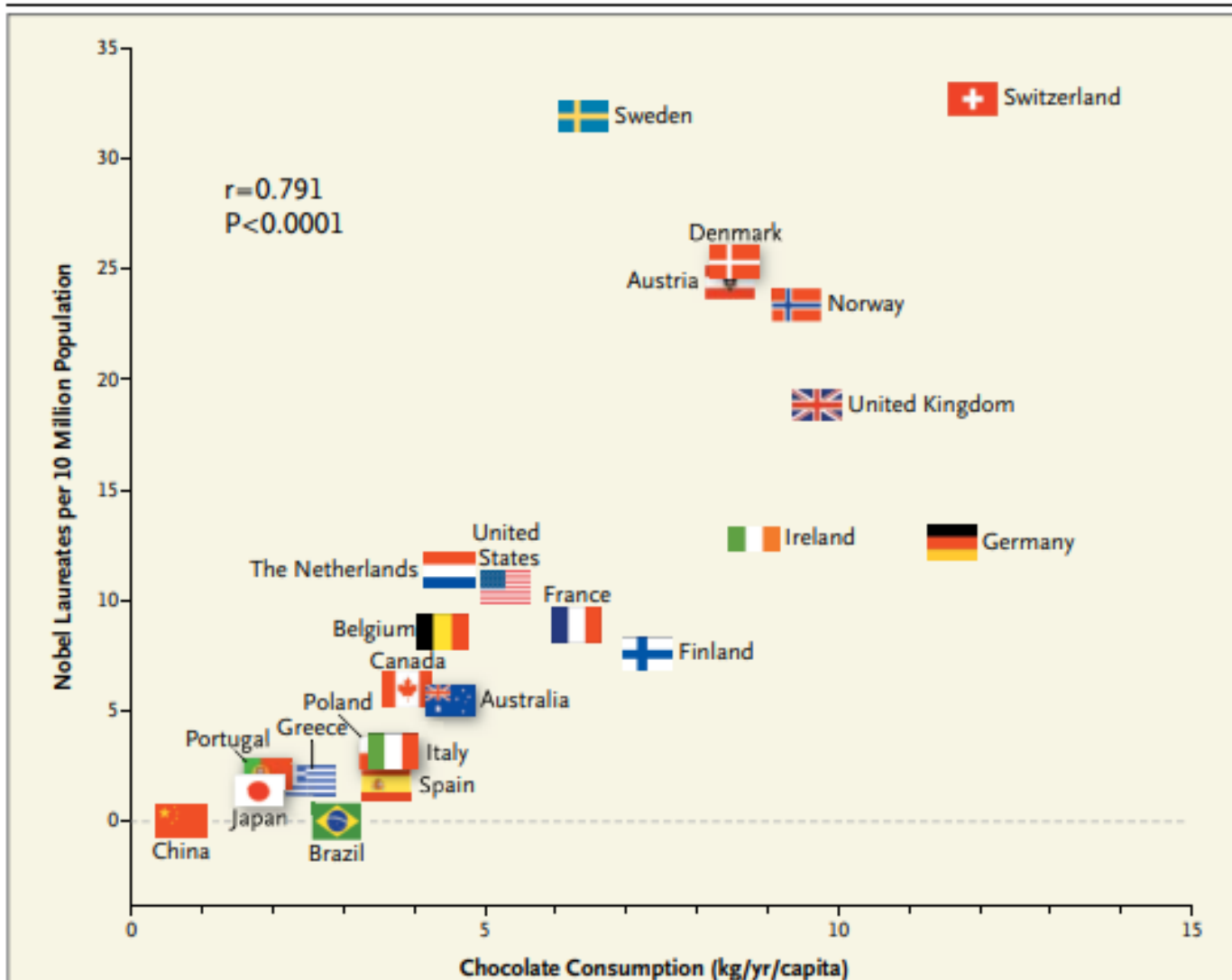


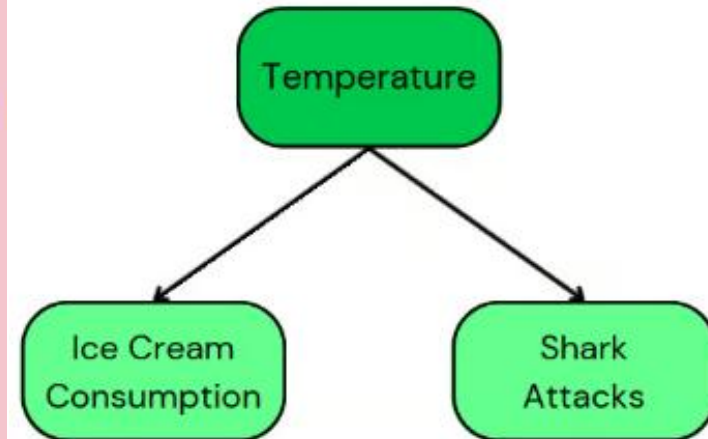
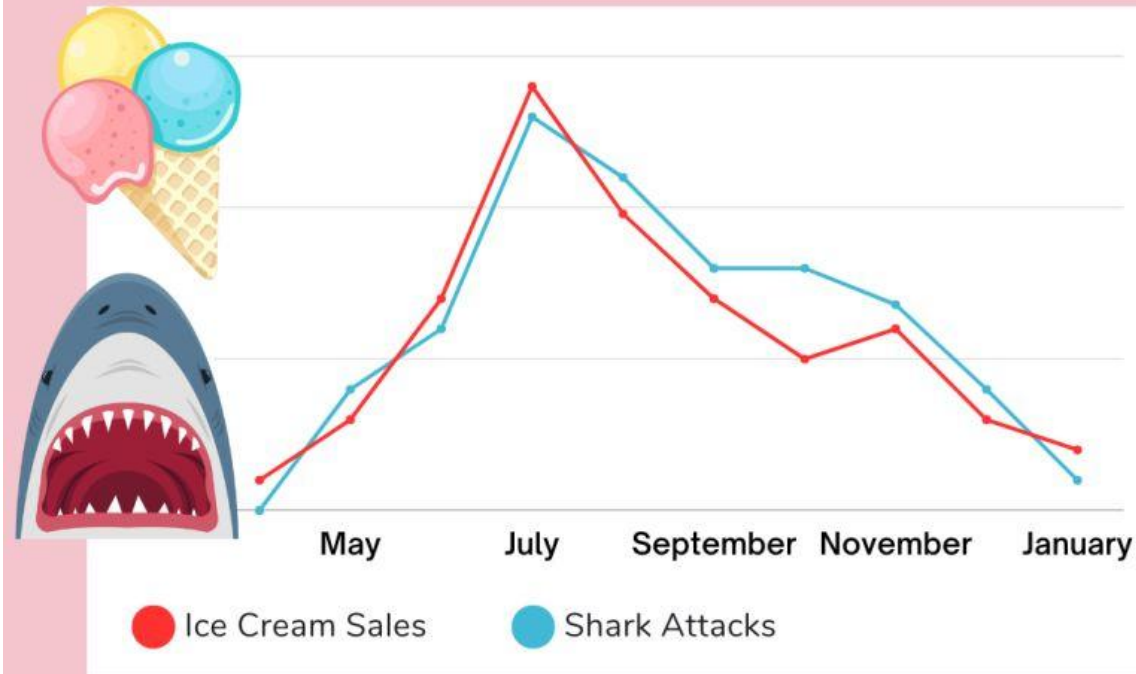
Figure 1. Correlation between Countries' Annual Per Capita Chocolate Consumption and the Number of Nobel Laureates per 10 Million Population.

Common Cause

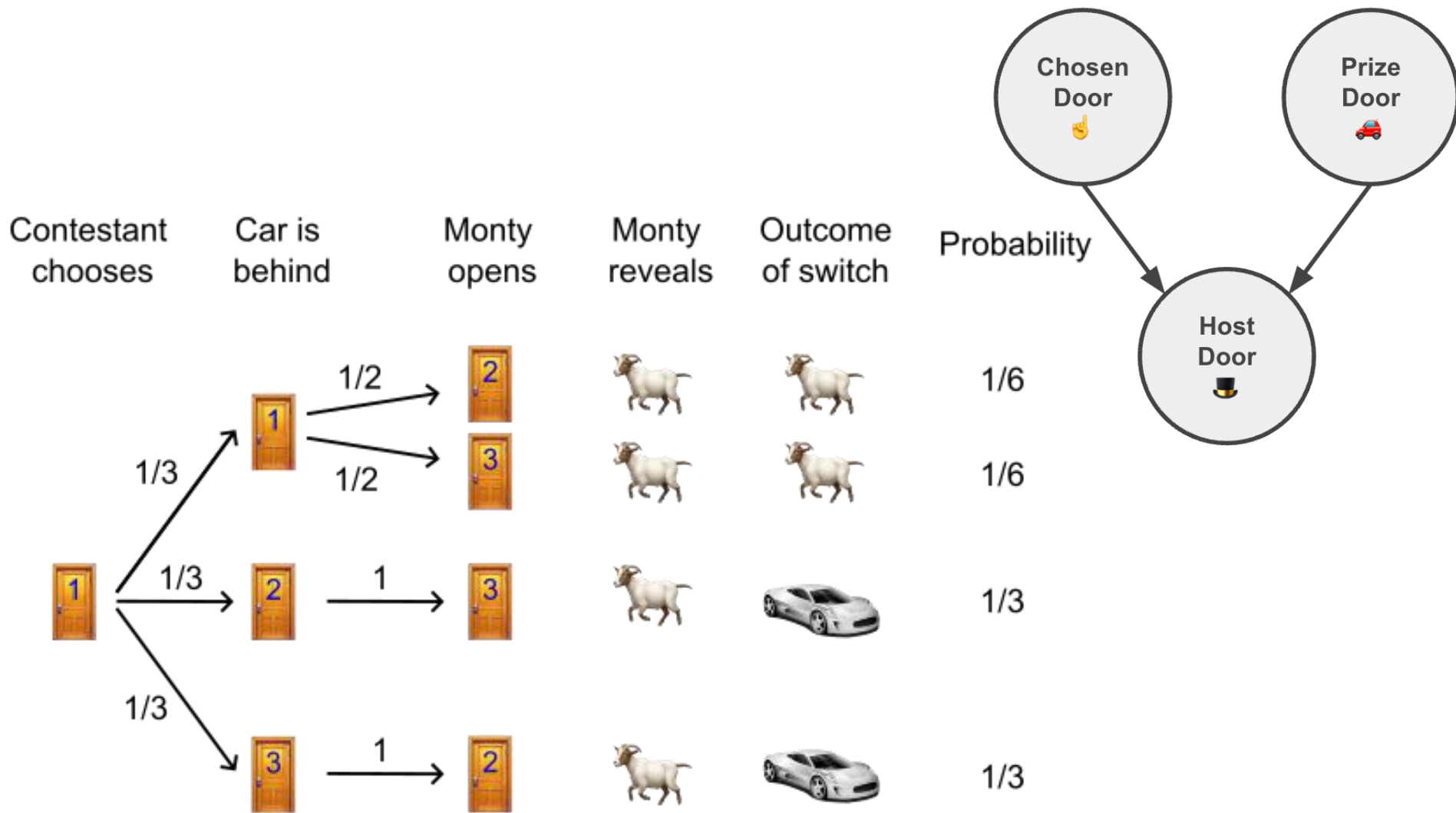
ATTENTION ATTENTION ATTENTION

ICE CREAM CAUSES SHARK ATTACKS?

Look at the data:



Collider – Monty Hall Problem



What can go wrong?

- Source and quality of the data
 - Not random sample. Ex. Biased sample, sample of convenience.
 - Important predictors may not have been observed.
 - Observational data often make causal conclusion problematic.
 - The range of \mathbf{X} and qualitative nature of some predictors may limit effective predictions, it's unsafe to extrapolate.
- Error component ε
 - May have unequal variance.
 - May be correlated.
 - May not be normally distributed.
- Structural component $\mathbf{X}\boldsymbol{\beta}$
 - Errors in \mathbf{X} . (Measurement error)
 - Series collinearity in \mathbf{X} .
- Many theory rests on the assumption that the model is correct.
- Experimenter Bias.

GLS (Generalized Least Square)

- What if $\text{Var}(\boldsymbol{\varepsilon}) \neq \sigma^2 \mathbf{I}$?
 - Time series correlation. $\varepsilon_t \sim \text{ARMA}(p, g)$
 - Repeated measurement model.
 - Spatial correlation.
 - Nested errors.
- Consider $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \boldsymbol{\Sigma}$ where $\boldsymbol{\Sigma}$ is known but σ^2 is unknown. Since $\boldsymbol{\Sigma}$ is symmetric and positive definite ($\mathbf{x}^\top \boldsymbol{\Sigma} \mathbf{x} > 0$), by Cholesky decomposition $\boldsymbol{\Sigma} = \mathbf{S}\mathbf{S}^\top$.
- Multiply both side of $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$
$$\mathbf{y}^* = \mathbf{S}^{-1}\mathbf{y} = \mathbf{S}^{-1}\mathbf{X}\boldsymbol{\beta} + \mathbf{S}^{-1}\boldsymbol{\varepsilon} = \mathbf{X}^*\boldsymbol{\beta} + \boldsymbol{\varepsilon}^*$$
then
$$\text{Var}(\mathbf{S}^{-1}\boldsymbol{\varepsilon}) = \mathbf{S}^{-1}\sigma^2\boldsymbol{\Sigma}(\mathbf{S}^{-1})^\top = \sigma^2\mathbf{S}^{-1}(\mathbf{S}\mathbf{S}^\top)(\mathbf{S}^{-1})^\top = \sigma^2\mathbf{I}$$
- Find $\boldsymbol{\beta}$ that minimize
$$\mathcal{S}(\boldsymbol{\beta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$
$$\Rightarrow \hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{y}$$

Lack of fit

- What fit is appropriate?

X too simple

- Not enough to explain the mean structure in data.
- Lack of fit. (under-fitting)
- $\hat{\sigma}^2$ over-estimated.

X too complex

- Will explain the variation caused by errors, in addition to the mean structure.
- Overfit the data.
- $\hat{\sigma}^2$ under-estimated.

- What statistic carry the information about lack of fit or overfit?
 - $\hat{\sigma}^2$

- $H_0: E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$ is correct vs $H_1: E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$ is too simple
 - Compare $\hat{\sigma}^2$ to σ^2
 σ^2 is known vs σ^2 is unknown

Diagnostics

- Check model assumptions to suggest further improvement after fitting. The building of an empirical model is an iterative process. During process, it required to check whether the current fitted model is consistent with data.
- What assumptions need to be checked?
model: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$
 - Error structure: error independent, homoscedasticity, normally distributed.
 - Mean structure: whether $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$ is a correct structure.
 - Unusual observations: where some observations do not fit the model.
- Two types of diagnostic techniques
 - Numerical vs Graphical

Leverage

- Recall residual

$$\hat{\varepsilon} = \mathbf{y} - \hat{\mathbf{y}} = \mathbf{y} - \mathbf{H}\mathbf{y} = (\mathbf{I} - \mathbf{H})\mathbf{y} = (\mathbf{I} - \mathbf{H})\boldsymbol{\varepsilon}$$

- Leverage, consider the diagonal of \mathbf{H} , let

$$h_i = \mathbf{H}_{i,i}$$

- \mathbf{x}_i whose h_i is $\begin{vmatrix} \text{large} \\ \text{small} \end{vmatrix}$ then $\text{Var}(\hat{\varepsilon}_i)$ $\begin{vmatrix} \text{small} \\ \text{large} \end{vmatrix}$

$\Rightarrow \begin{vmatrix} \text{fitted model has to force to fit close to } y_i \\ \text{in this } \mathbf{x}_i, \text{ model cannot fit so well} \end{vmatrix}$

- h_i roughly determines how close (\mathbf{x}_i, y_i) to the surface $(\mathbf{x}_i, \hat{y}_i)$
 - h_i corresponds to Mahalanobis distance.
- What h_i are too large?

$$\sum_{i=1}^n h_i = p, \quad \frac{1}{n} \leq h_i \leq 1 \quad \forall i \Rightarrow \text{large leverage} \gg \frac{p}{n}$$

\Rightarrow rule of thumb: if $h_i > \frac{2p}{n}$

Residuals - 1

- (Internally) studentized residual r_i
 - Because $\text{Var}(\hat{\varepsilon}_i) = (1 - h_i)\sigma^2$ then

$$r_i = \frac{\hat{\varepsilon}_i}{\hat{\sigma}\sqrt{1 - h_i}}, \quad \text{Var}(r_i) \approx 1$$

- (Externally) jackknife residual t_i
 - Exclude i^{th} observation and re-compute the estimates to get $\hat{\beta}_{(i)}$ and $\hat{\sigma}_{(i)}$, where (i) denotes that i^{th} case has been excluded.
 - Since

$$\text{Var}(y_i - \hat{y}_{(i)}) = \sigma^2 \left[1 + \mathbf{x}_i^\top (\mathbf{X}_{(i)}^\top \mathbf{X}_{(i)})^{-1} \mathbf{x}_i \right]$$

$$\begin{aligned} t_i &= \frac{y_i - \hat{y}_{(i)}}{\hat{\sigma}_{(i)} \sqrt{1 + \mathbf{x}_i^\top (\mathbf{X}_{(i)}^\top \mathbf{X}_{(i)})^{-1} \mathbf{x}_i}} \\ &= r_i \sqrt{\frac{n - p - 1}{n - p - r_i^2}} \end{aligned}$$

Residuals - 2

- Test for outliers

- Given a specific case i , conclude an outlier if

$$|t_i| > T_{n-p-1} \left(\frac{\alpha}{2} \right)$$

- In practice, a few (or all) t_i will be tested \Rightarrow Multiple testing

- Bonferroni Correction

- Test need to adjust the significant level of the test accordingly.
- Consider the test

H_0 : no outlier in the n obs **vs** H_1 : at least one outlier

$$\alpha^* = P \left(\bigcup_{i=1}^n R_i \middle| H_0 \right) \leq \sum_{i=1}^n P(R_i | H_0) = n\alpha$$

- Conclude an outlier if

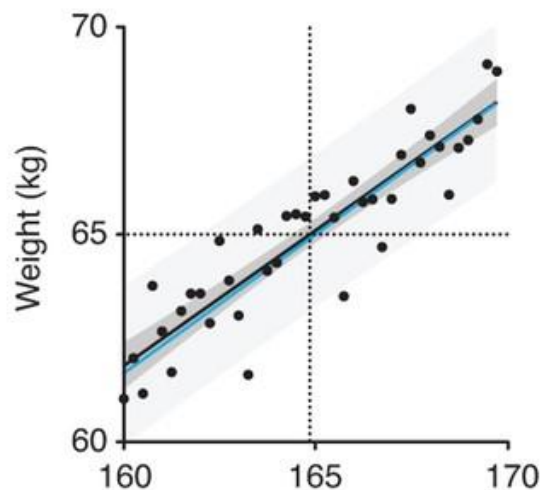
$$|t_i| > t_{n-p-1} \left(\frac{\alpha}{2n} \right)$$

- It is conservative, tends not to label points as outlier.

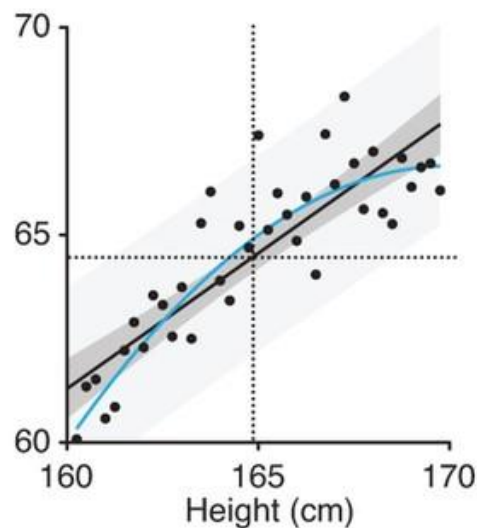
Residual Plot

a

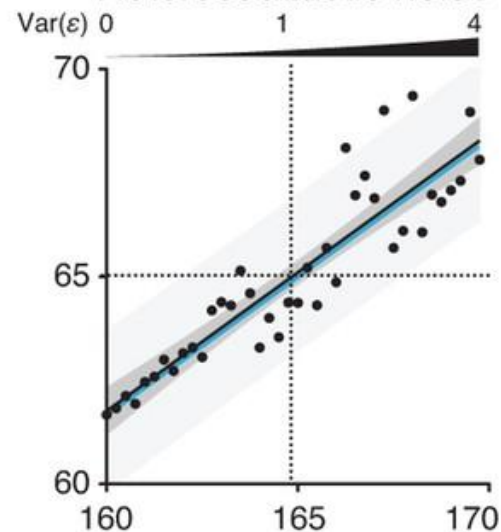
All assumptions satisfied



Nonlinear term in model

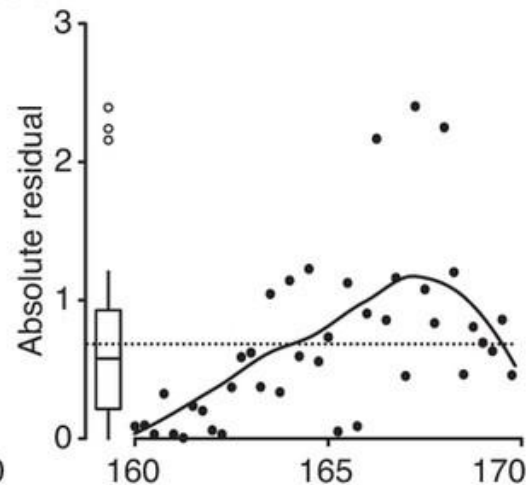
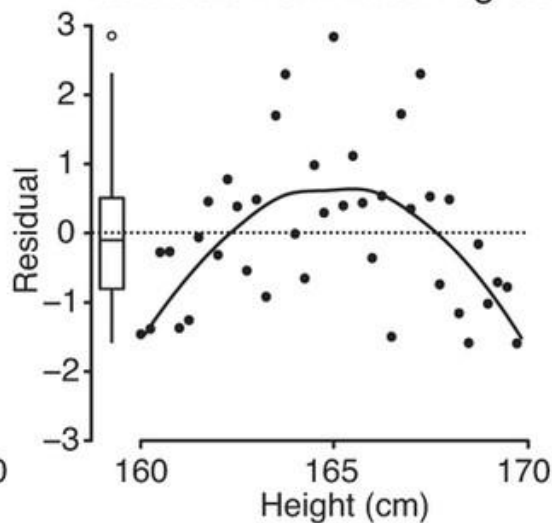
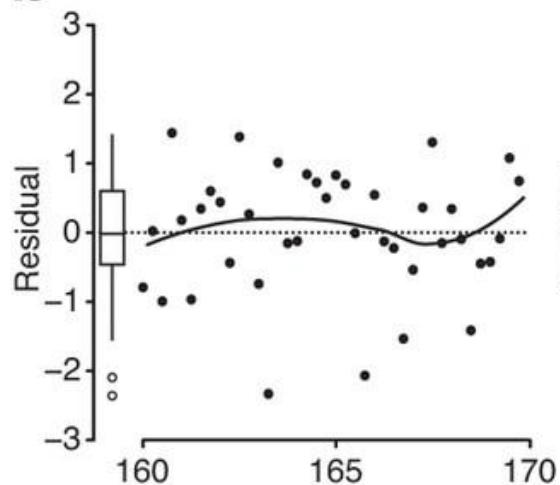


Heteroscedastic noise



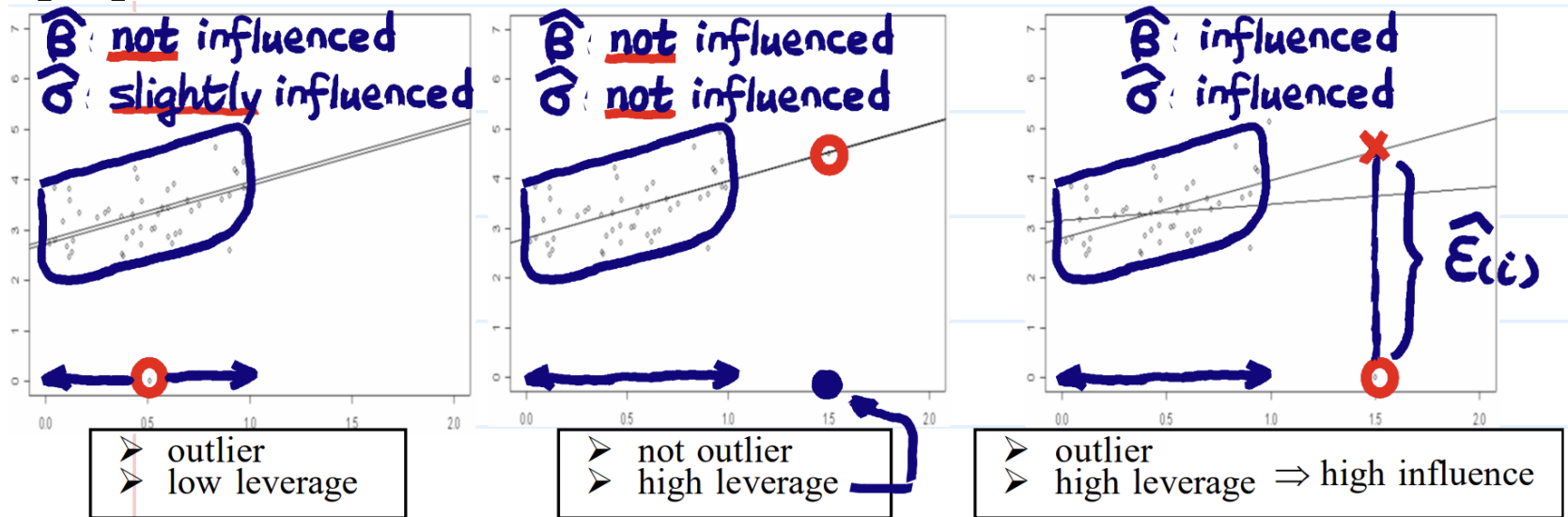
b

Residuals from fitted regression



Influential observation

- An influential point may or may not be an outlier and may or may not have large leverage but it will tend to have at least one of those two properties.



- Cook's distance (scale and unit free)

$$D_i = \frac{(\hat{\beta} - \hat{\beta}_{(i)})^T (X^T X) (\hat{\beta} - \hat{\beta}_{(i)})}{p \hat{\sigma}^2} = \left(\frac{1}{p} \right) r_i^2 \left(\frac{h_i}{1 - h_i} \right)$$

- It is a combination of **residual** and **leverage**. If assume \mathbf{X} is normal, can do a test on D_i .

Variance Stabilizing transformation - 1

- It's better try to understand the cause of non-constant variance before taking any remedies.
 - Large response have more room to vary.
 - Response constrained to lie between a maximum and a minimum.
 - Response from Poisson or binomial distribution,...
- Remedies for non-constant variance
 - Weighted least squares
 - Find a transformation h such that $\text{Var}(h(y))$ is a constant.

$$h(y) \approx h(E(y)) + h'(E(y))(y - E(y))$$

$$\text{Var}(h(y)) \approx \{h'(E(y))\}^2 \text{Var}(y) = C \Rightarrow h'(E(y)) \propto \frac{1}{\sqrt{\text{Var}(y)}}$$

$$h(E(y)) = \int \frac{1}{\sqrt{\text{Var}(y)}} d(E(y)) = \int \frac{1}{\sqrt{g(\mu)}} d(\mu)$$

- Ex. $\text{Var}(y) \propto [E(y)]^2 \Rightarrow g(\mu) = \mu^2$

$$h(\mu) = \int \frac{1}{\sqrt{\mu^2}} d\mu = \ln(\mu) \Rightarrow \text{suggest } h(y) = \ln(y)$$

Variance Stabilizing transformation - 2

TABLE 5.1 Useful Variance-Stabilizing Transformations

Relationship of σ^2 to $E(y)$	Transformation
$\sigma^2 \propto \text{constant}$	$y' = y$ (no transformation)
$\sigma^2 \propto E(y)$	$y' = \sqrt{y}$ (square root; Poisson data)
$\sigma^2 \propto E(y)[1 - E(y)]$	$y' = \sin^{-1}(\sqrt{y})$ (arcsin; binomial proportions $0 \leq y_i \leq 1$)
$\sigma^2 \propto [E(y)]^2$	$y' = \ln(y)$ (log)
$\sigma^2 \propto [E(y)]^3$	$y' = y^{-1/2}$ (reciprocal square root)
$\sigma^2 \propto [E(y)]^4$	$y' = y^{-1}$ (reciprocal)

- Other transformation

- Box-Cox transformation (make the errors as nearly like normal)

$$t_{\lambda}(y) = \begin{cases} \frac{y^{\lambda} - 1}{\lambda}, & \text{if } \lambda \neq 0 \\ \ln(y), & \text{if } \lambda = 0 \end{cases}$$

choose λ to fit data well using maximum likelihood

$$H_0: \lambda = \lambda_0 \text{ vs } H_1: \lambda \neq \lambda_0$$

the likelihood ratio test

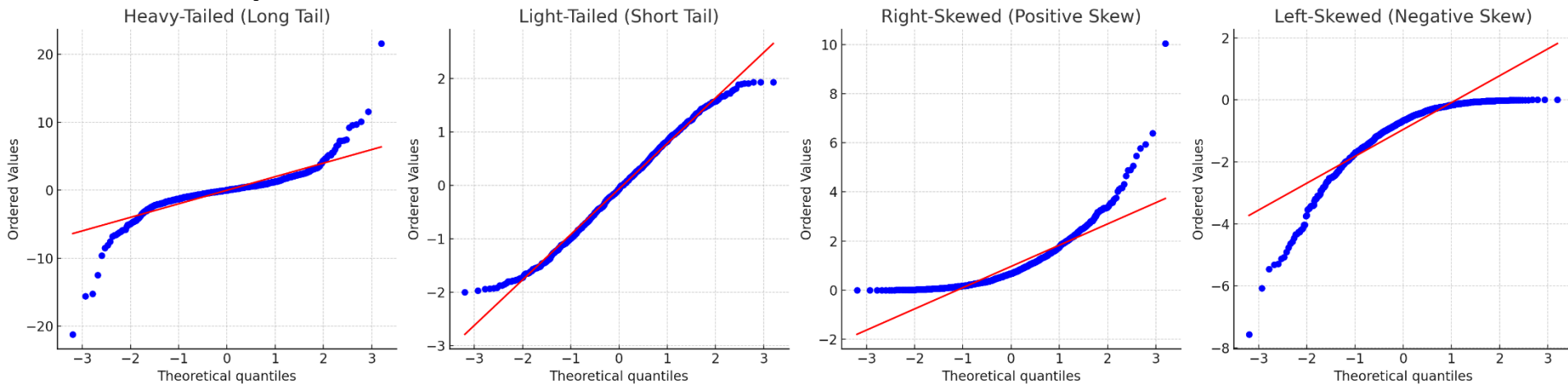
$$-2[\mathcal{L}(\lambda_0) - \mathcal{L}(\hat{\lambda})] \sim \chi_{(1)}^2 \text{ under } H_0$$

Normality Test

- It is used to determine if data is well-modeled by normal distribution and to compute how likely it is for a random variable underlying the data to be normal.
 - Bell curve shape.
 - Symmetric around the $x = \mu$.
 - Log-concave.
- Graphical methods
 - Normal probability plot.
 - Quantile-Quantile plot (Q-Q plot).
- Numerical methods
 - D'Agostino's K-squared Test.
 - Jarque-Bera Test.
 - Anderson-Darling Test.
 - Kolmogorov-Smirnov Test.
 - Shapiro-Wilk Test.

Q-Q plot

- Assessing normality assumption of ε
 - Sort the data $\hat{\varepsilon}_{(1)} \leq \hat{\varepsilon}_{(2)} \leq \dots \leq \hat{\varepsilon}_{(n)}$
 - Plot $\hat{\varepsilon}_{(i)}$ against $\Phi^{-1}\left(\frac{i}{n+1}\right)$ where $\Phi(\cdot)$ is cdf of $N(0,1)$
- If the residuals are normal distributed, an approximately straight-line relationship will be observed.
- Non-normality:
 - long-tail
 - short-tail
 - asymmetric



Collinearity

- A challenge arises when $\mathbf{X}^\top \mathbf{X}$ is close to singular but not exactly so. It leads to imprecise estimates of $\boldsymbol{\beta}$. The standard errors are inflated so that T test may fail to reveal significant factors.
- It can be detected in several ways:

- VIF (Variance Inflation factor)

$$\text{Var}(\hat{\beta}_j) = \sigma^2 \left(\frac{1}{1 - R_j^2} \right) \left(\frac{1}{\sum_i (g_{ij} - \bar{g}_j)^2} \right)$$

- Let $\lambda_1 \geq \dots \geq \lambda_p \geq 0$ be eigenvalues of $\mathbf{X}^\top \mathbf{X}$. Zero eigenvalues denote exact collinearity, while the presence of some small eigenvalues indicates collinearity. The condition number κ measures the relative sizes of the eigenvalues and is defined

$$\kappa = \sqrt{\lambda_1 / \lambda_p}$$

where $\kappa > 30$ is consider large. Other condition numbers, $\sqrt{\lambda_1 / \lambda_i}$ are also worth considering.

PCA - 1

- Let Σ be the covariance matrix associated with the $\mathbf{X} = (X_1, \dots, X_p)$, and Σ have the eigenvalue-eigenvector pairs $(\lambda_1, e_1), \dots, (\lambda_p, e_p)$ where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$. Then the i^{th} principal component (PC) is given by

$$Z_i = e_i^\top \mathbf{X} = e_{i,1}X_1 + \dots + e_{i,p}X_p, \quad i = 1, \dots, p$$

with these choices

$$\text{Var}(Z_i) = e_i^\top \Sigma e_i = \lambda_i, \quad i = 1, \dots, p$$

$$\text{Cov}(Z_i, Z_j) = e_i^\top \Sigma e_j = 0, \quad i \neq j$$

Then

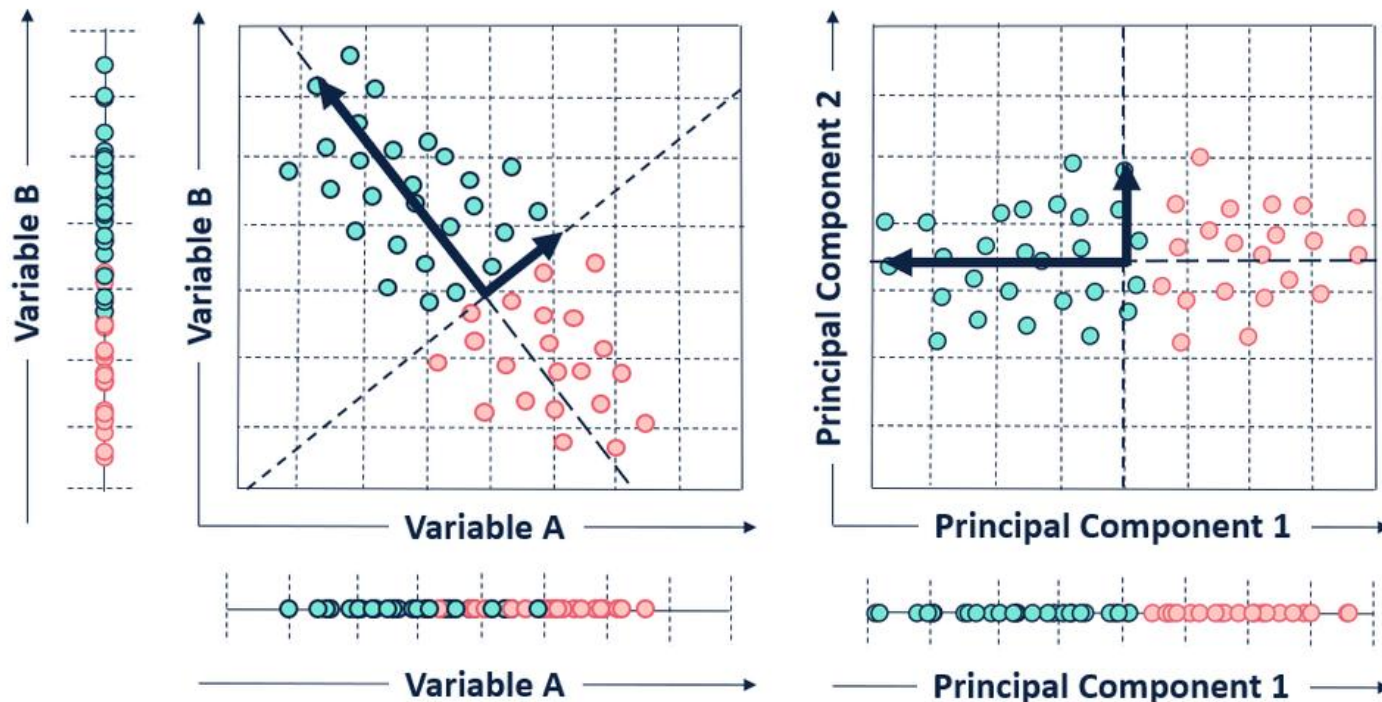
$$\sigma_{1,1} + \dots + \sigma_{p,p} = \sum_{i=1}^p \text{Var}(X_i) = \lambda_1 + \dots + \lambda_p = \sum_{i=1}^p \text{Var}(Z_i)$$

- Proportion of total population variance due to i^{th} PC is

$$\frac{\lambda_i}{\lambda_1 + \dots + \lambda_p}, \quad i = 1, \dots, p$$

PCA - 2

- For a geometric interpretation of PCA, suppose we have two features A and B. The 1st PC is a line through the widest part; the 2nd PC is the line at right angles to the 1st PC.
- In other words, the 1st PC goes through the fattest part of the "football" and the 2nd PC through the next fattest part of the "football" and orthogonal to the 1st PC; and so on.



PCA - 3

- How many PCs to retain?
 - Screen plot.
 - Amount of total variance explained (say 70% ~ 90%).
 - Cutoff point 1 or 0.7.
- What are the distributions of the $(\hat{\lambda}_1, \dots, \hat{\lambda}_p)$ and $(\hat{e}_1, \dots, \hat{e}_p)$?
 - Let $\Lambda = \text{diag}\{\lambda_1, \dots, \lambda_p\}$ of Σ , then $\sqrt{n}(\hat{\lambda} - \lambda) \sim N(\mathbf{0}, 2\Lambda^2)$
 - Let

$$E_i = \lambda_i \sum_{\substack{k=1 \\ k \neq i}}^p \frac{\lambda_k}{(\lambda_k - \lambda_i)^2} e_k e_k^\top$$

then $\sqrt{n}(\hat{e}_i - e_i) \sim N(\mathbf{0}, E_i)$

- A large sample $100(1 - \alpha)\%$ confidence interval for λ_i is thus

$$\frac{\hat{\lambda}_i}{\left(1 + Z_{\alpha/2} \sqrt{2/n}\right)} \leq \lambda_i \leq \frac{\hat{\lambda}_i}{\left(1 - Z_{\alpha/2} \sqrt{2/n}\right)}$$

Serial correlation

- When a time order is available
 - Plot $\hat{\varepsilon}$ against time
 - Plot $\hat{\varepsilon}_{t+1}$ against $\hat{\varepsilon}_t$, when t related to time
 - Use formal tests (Ex. Durbin-Watson)

$$DW = \frac{\sum_{t=2}^T (\hat{\varepsilon}_t - \hat{\varepsilon}_{t-1})^2}{\sum_{t=1}^T \hat{\varepsilon}_t^2} \in [0, 4]$$

- positive correlated $\Rightarrow DW \rightarrow 0$
- negative correlated $\Rightarrow DW \rightarrow 4$
- under null $\Rightarrow DW \approx 2$
- Null distribution depends on X
- Rewrite the test statistic by lag one sample autocorrelation

$$DW \approx 2(1 - \hat{\rho}_1), \quad \hat{\rho}_1 = \frac{\sum_{t=1}^T \hat{\varepsilon}_t \hat{\varepsilon}_{t+1}}{\sum_{t=1}^T \hat{\varepsilon}_t^2}$$

The DW statistic gets smaller as the serial correlation increase.

Dummy

- Categorical (qualitative) predictors
 - Nominal vs ordinal.
 - What properties can we explore for qualitative predictor?

$$\text{category } i \rightarrow y_{i,j}, \quad \mu_i = E(y_{i,j})$$

can only study difference between μ_i .

- One dichotomous predictor (Two categories)

$$d(C) = \begin{cases} 1, & \text{if } C = c_1 \\ 0, & \text{if } C = c_2 \end{cases}$$

$$\text{Ex. } y = \beta_0 + \beta_1 d + \varepsilon$$

$$\begin{cases} \mu_1 = E(y|d = 0) = \beta_0 \\ \mu_2 = E(y|d = 1) = \beta_0 + \beta_1 \end{cases} \Rightarrow \begin{cases} \beta_0 = \mu_1 \\ \beta_1 = \mu_2 - \mu_1 \end{cases}$$

- One polytomous predictor (More than two categories)
 - For k categories, dummy variables d_1, \dots, d_{k-1} are needed to depict the difference between categories.

Coding Scheme

- Consider $y = \beta_0 + \beta_1 d_1 + \beta_2 d_2 + \beta_3 d_3 + \varepsilon$ with 4 categories c_1, c_2, c_3, c_4 .

- Treatment coding (c_1 as a reference)

$$\begin{cases} \mu_1 = E(y|d_1 = 0, d_2 = 0, d_3 = 0) = \beta_0 \\ \mu_2 = E(y|d_1 = 1, d_2 = 0, d_3 = 0) = \beta_0 + \beta_1 \\ \mu_3 = E(y|d_1 = 0, d_2 = 1, d_3 = 0) = \beta_0 + \beta_2 \\ \mu_4 = E(y|d_1 = 0, d_2 = 0, d_3 = 1) = \beta_0 + \beta_3 \end{cases} \Rightarrow \begin{cases} \beta_0 = \mu_1 \\ \beta_1 = \mu_2 - \mu_1 \\ \beta_2 = \mu_3 - \mu_1 \\ \beta_3 = \mu_4 - \mu_1 \end{cases}$$

- Sum coding (β_0 represent overall mean)

$$\begin{cases} \mu_1 = E(y|d_1 = -1, d_2 = -1, d_3 = -1) = \beta_0 - \beta_1 - \beta_2 - \beta_3 \\ \mu_2 = E(y|d_1 = 1, d_2 = 0, d_3 = 0) = \beta_0 + \beta_1 \\ \mu_3 = E(y|d_1 = 0, d_2 = 1, d_3 = 0) = \beta_0 + \beta_2 \\ \mu_4 = E(y|d_1 = 0, d_2 = 0, d_3 = 1) = \beta_0 + \beta_3 \end{cases}$$

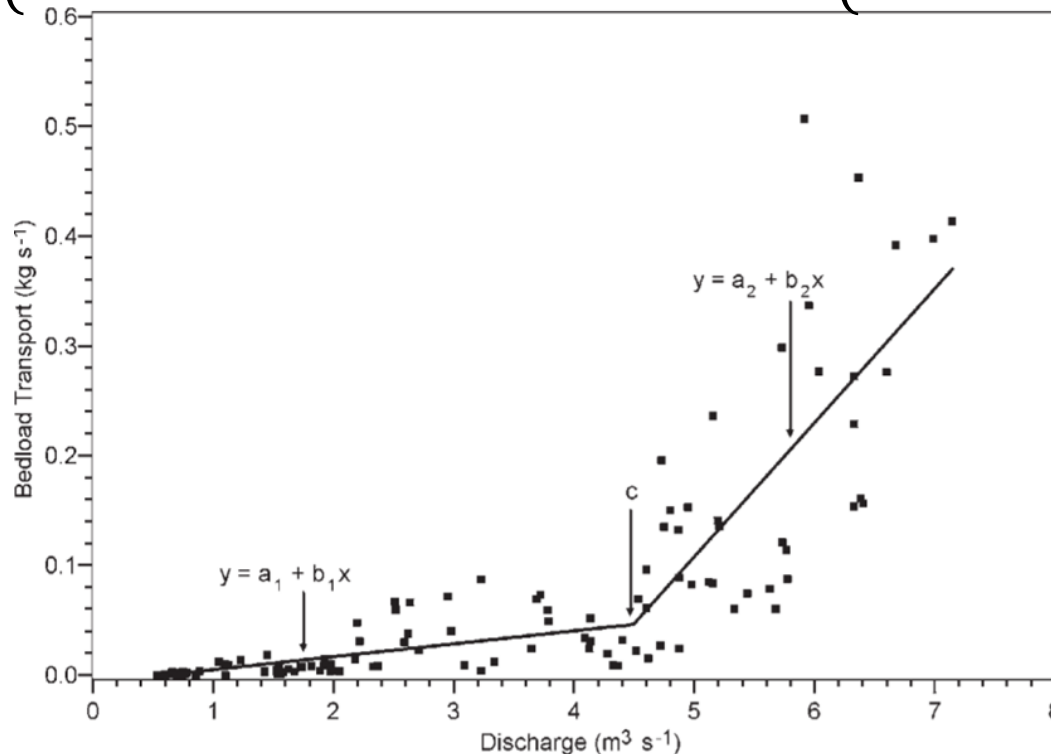
$$\Rightarrow \begin{cases} \beta_0 = \frac{\mu_1 + \mu_2 + \mu_3 + \mu_4}{4} = \bar{\mu} \\ \beta_1 = \mu_2 - \bar{\mu} \\ \beta_2 = \mu_3 - \bar{\mu} \\ \beta_3 = \mu_4 - \bar{\mu} \end{cases}$$

Piecewise Regression

- Sometimes we have reason to believe that different linear regression models apply in different region of data.
- Consider a simple case

$$y = \beta_0 + \beta_1 b_L(x) + \beta_2 b_R(x) + \varepsilon$$

$$b_L(x) = \begin{cases} c - x, & \text{if } x < c \\ 0, & \text{ow} \end{cases}, \quad b_R(x) = \begin{cases} x - c, & \text{if } x > c \\ 0, & \text{ow} \end{cases}$$



Polynomial regression

- Polynomial regression: The relationship between response and predictors is smooth, but not a straight line.

- 1 predictor with general degree d

$$y = \beta_0 + \beta_1 X + \beta_2 X^2 + \cdots \beta_d X^d + \varepsilon$$

- 2 predictors with degree $d = 2$

$$y = \beta_0 + (\beta_1 X_1 + \beta_2 X_2) + (\beta_{1,1} X_1^2 + \beta_{2,2} X_2^2) + \beta_{1,2} X_1 X_2 + \varepsilon$$

- m predictors with degree $d = 2$

$$y = \beta_0 + \left(\sum_{i=1}^m \beta_i X_i \right) + \left(\sum_{i=1}^m \beta_{i,i} X_i^2 \right) + \left(\sum_{1 \leq i < j \leq m} \beta_{i,j} X_i X_j \right) + \varepsilon$$

- Increasing degree $d \Rightarrow$ model may have too many parameters

$$1 + m + m + \binom{m}{2} = \frac{(m+1)(m+2)}{2}$$

- Polynomial terms can cause numerical instability and collinearity. It is hard to fit jump function.

Inner Product Space

- An inner product space is a vector space \mathcal{V} over the field \mathcal{F} together with an inner product, that is, a map

$$\langle \cdot, \cdot \rangle : \mathcal{V} \times \mathcal{V} \rightarrow \mathcal{F}$$

that satisfies the three properties $\forall x, y, z \in \mathcal{V}$ and $a, b \in \mathcal{F}$

- Symmetry $\langle x, y \rangle = \langle y, x \rangle$
 - Linearity $\langle ax + by, z \rangle = a\langle x, z \rangle + b\langle y, z \rangle$
 - Positive definiteness $\langle x, x \rangle > 0$ if $x \neq 0$
- Every inner product space induces a norm, that is defined by

$$\|x\| = \sqrt{\langle x, x \rangle}$$

with this norm, every inner product space becomes a normed vector space. Every general property of normed vector spaces applies to inner product spaces.

- Ex. (Euclidean Vector Space)

$\langle \cdot, \cdot \rangle : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ is a inner product on \mathbb{R}^n if and only if exist a symmetric positive-definite matrix M such that

$$\langle x, y \rangle = x^T M y, \quad \forall x, y \in \mathbb{R}^n$$

if M is identity matrix then $\langle \cdot, \cdot \rangle$ is the dot product.

Orthogonal Polynomial

- An orthogonal polynomial sequence $\{P_n\}_{n=0}^{\infty}$ is a family of polynomials such that

$$\langle P_n, P_m \rangle = 0, \quad m \neq n$$

Usually the sequence is required to be orthonormal, namely

$$\langle P_n, P_n \rangle = 1$$

- It is common to define the inner product of functions f and g w.r.t a nonnegative weight function σ over an interval $[a, b]$

$$\langle f, g \rangle_{\sigma} = \int_a^b f(x)g(x)\sigma(x)dx, \quad \|f\| = \sqrt{\langle f, f \rangle_{\sigma}}$$

Various polynomial sequences

- Hermite polynomials $\Leftrightarrow \langle f, g \rangle_{\sigma} = \int_{-\infty}^{\infty} f(x)g(x)e^{-x^2}dx$
- Legendre polynomials $\Leftrightarrow \langle f, g \rangle_{\sigma} = \int_{-1}^1 f(x)g(x)dx$
- Chebyshev polynomials $\Leftrightarrow \langle f, g \rangle_{\sigma} = \int_{-1}^1 f(x)g(x)\frac{1}{\sqrt{1-x^2}}dx$
- Laguerre polynomials $\Leftrightarrow \langle f, g \rangle_{\sigma} = \int_0^{\infty} f(x)g(x)e^{-x}dx$

Legendre Polynomials

- Let $\{B_n(x), n \in \mathbb{N}\}$ be a linearly independent sequence of $C([a, b])$. Then an orthogonal basis of functions, $\{P_n(x), n \in \mathbb{N}\}$ can be found and is given by

$$P_0(x) = B_0(x)$$

$$P_n(x) = B_n(x) - \sum_{j=0}^{n-1} \frac{\langle B_n, P_j \rangle_\sigma}{\langle P_j, P_j \rangle_\sigma} P_j(x), \quad n = 1, 2, \dots$$

- Ex. Consider $B_n(x) = x^n, x \in (-1, 1)$ and $\sigma(x) = 1$

$$P_0(x) = B_0(x) = 1$$

$$P_1(x) = B_1(x) - \frac{\langle B_1, P_0 \rangle_\sigma}{\langle P_0, P_0 \rangle_\sigma} P_0(x) = x - \frac{\langle x, 1 \rangle_\sigma}{2} 1 = x$$

$$P_2(x) = B_2(x) - \frac{\langle B_2, P_0 \rangle_\sigma}{\langle P_0, P_0 \rangle_\sigma} P_0(x) - \frac{\langle B_2, P_1 \rangle_\sigma}{\langle P_1, P_1 \rangle_\sigma} P_1(x)$$

$$= x^2 - \frac{\langle x^2, 1 \rangle_w}{2} 1 - \frac{\langle x^2, x \rangle_w}{2/3} x = x^2 - \frac{1}{3}$$

$$\Rightarrow \text{orthogonal set } \left\{ 1, x, x^2 - \frac{1}{3} \right\}$$

Relation to moments

- The orthogonal polynomials P_n can be expressed in terms of moments

$$\mu_n = \int_a^b x^n \sigma(x) dx$$

as follow

$$P_n(x) = c_n \det \begin{bmatrix} \mu_0 & \mu_1 & \cdots & \mu_n \\ \mu_1 & \mu_2 & \cdots & \mu_{n+1} \\ \vdots & \vdots & \ddots & \vdots \\ \mu_{n-1} & \mu_n & \cdots & \mu_{2n-1} \\ 1 & x & \cdots & x^n \end{bmatrix}$$

where the c_n are arbitrary (depend on the normalization of P_n).

- Ex. Consider the Legendre Polynomial

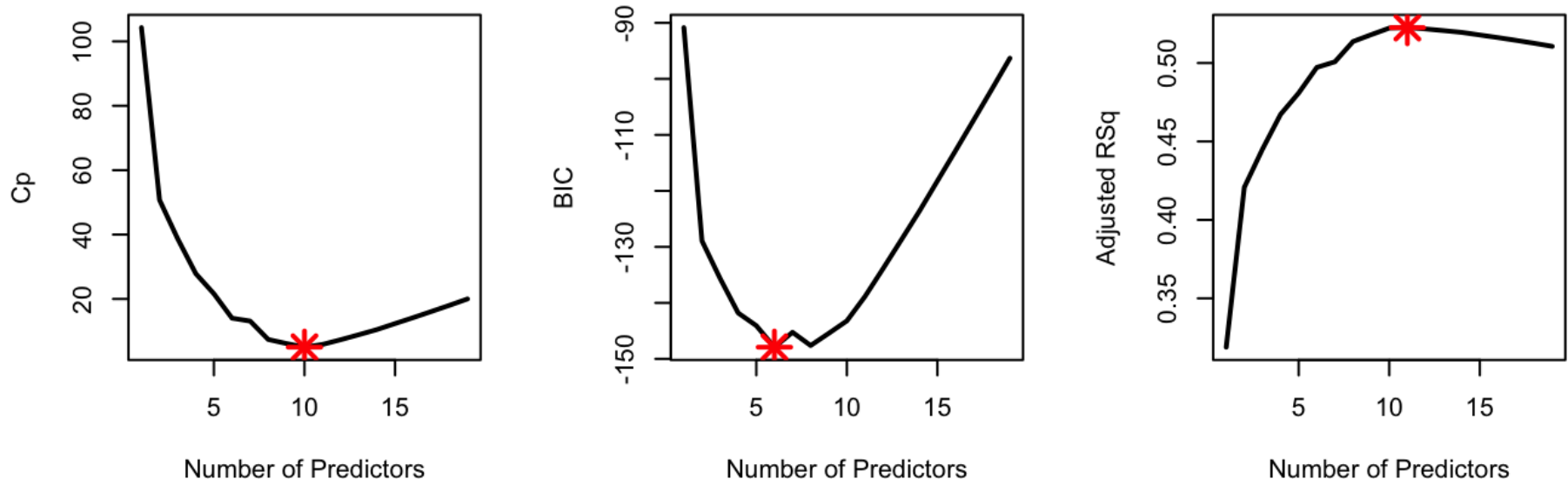
$$\begin{aligned} P_2(x) &= c_n \det \begin{bmatrix} \mu_0 & \mu_1 & \mu_2 \\ \mu_1 & \mu_2 & \mu_3 \\ 1 & x & x^2 \end{bmatrix} = c_n \det \begin{bmatrix} 2 & 0 & 2/3 \\ 0 & 2/3 & 0 \\ 1 & x & x^2 \end{bmatrix} \\ &= c_n \frac{4}{3} \left(x^2 - \frac{1}{3} \right) \end{aligned}$$

Model Selection - 1

- Testing-Based Procedures
 - May miss optimal model because of its one-at-a-time adding/dropping.
 - α -values (α -to-enter and α -to-remove) should not be treated too literally: because multiple testing occurring.
 - Removal of less significant terms tend to increase the significance of the remaining terms \Rightarrow it may lead to overstate the importance of the remaining terms.
 - For prediction purpose, testing-based procedure tends to pick smaller models than desired.
- Criterion-Based Procedures
 - AIC and BIC
$$\text{AIC} = -2(\text{maximized loglikelihood}) + 2p$$
$$\text{BIC} = -2(\text{maximized loglikelihood}) + \log(n)p$$
 - For regression,
$$-2(\text{maximized loglikelihood}) = n \log \left(\frac{\text{RSS}_{\mathcal{M}_p}}{n} \right) + \text{const}$$

Model Selection - 2

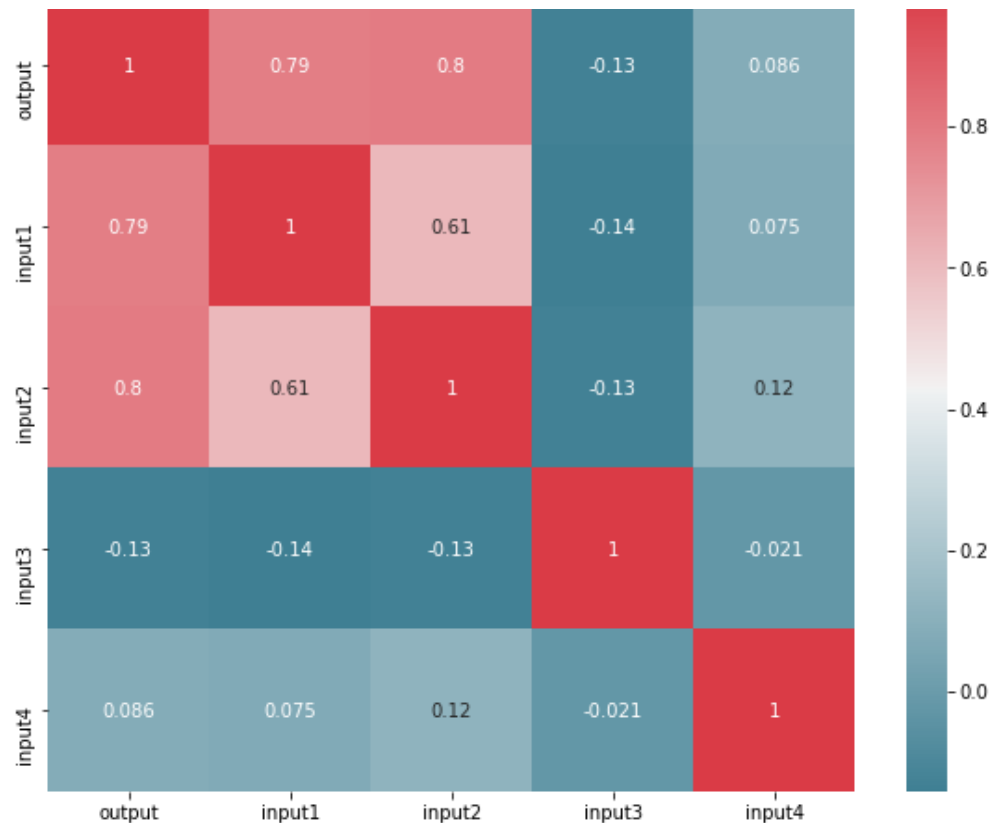
- Best Subset Selection:
 - Let \mathcal{M}_0 denote the which contains no predictor.
 - For $k = 1, \dots, p$
 - Fit all $\binom{p}{k}$ models that contain exactly k predictors.
 - Pick the best among these $\binom{p}{k}$ models and call it \mathcal{M}_k .
 - Select a single best model from among $\mathcal{M}_0, \mathcal{M}_1, \dots, \mathcal{M}_k$.



Model Selection - 3

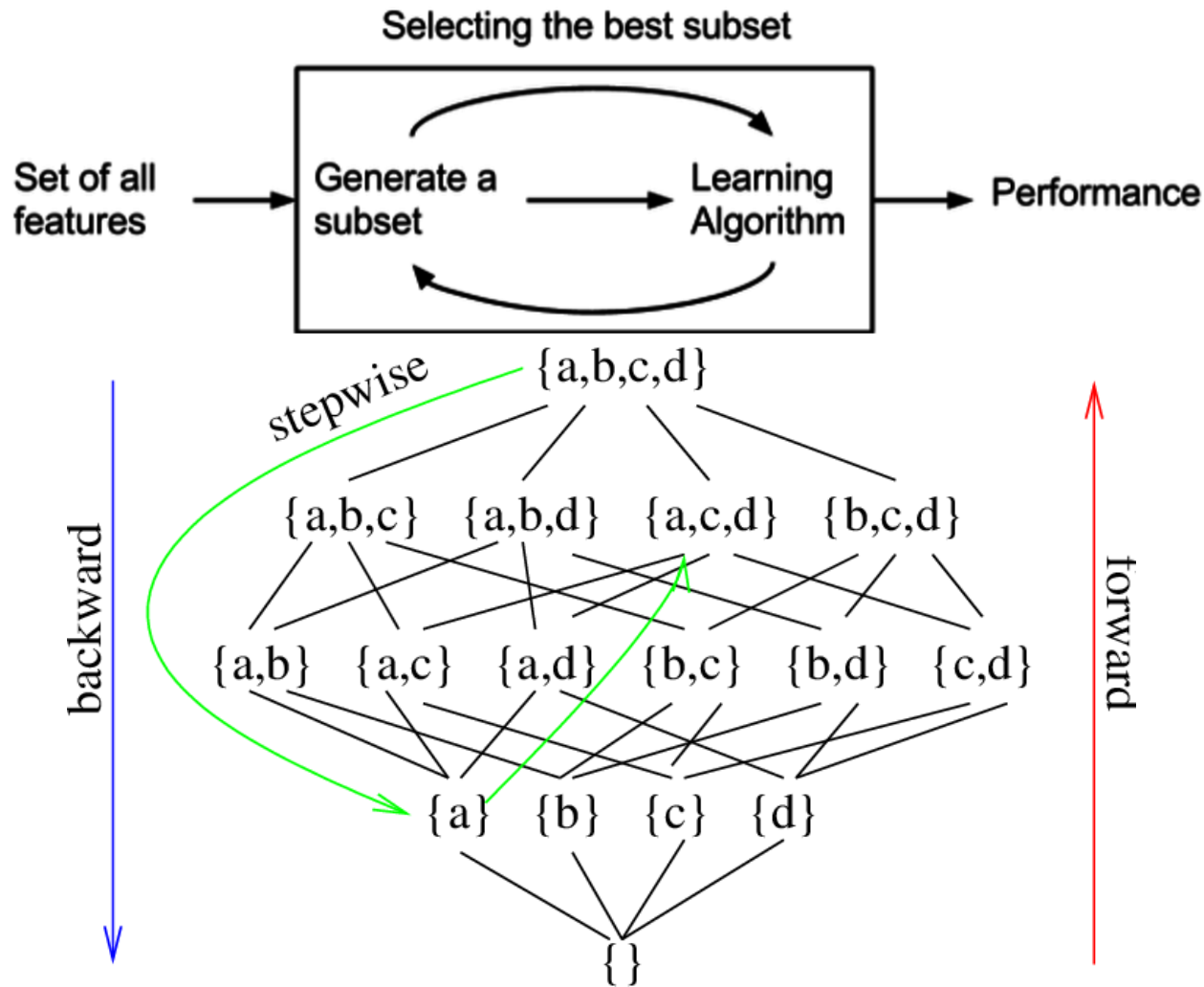
- Three types of selection procedure
 - Filter Mode (Select Best K): Apply a measure to assign a scoring to each feature. The features are ranked by the score and either selected to be kept or removed from the dataset.

Set of all Features → **Selecting the Best Subset** → **Learning Algorithm** → **Performance**



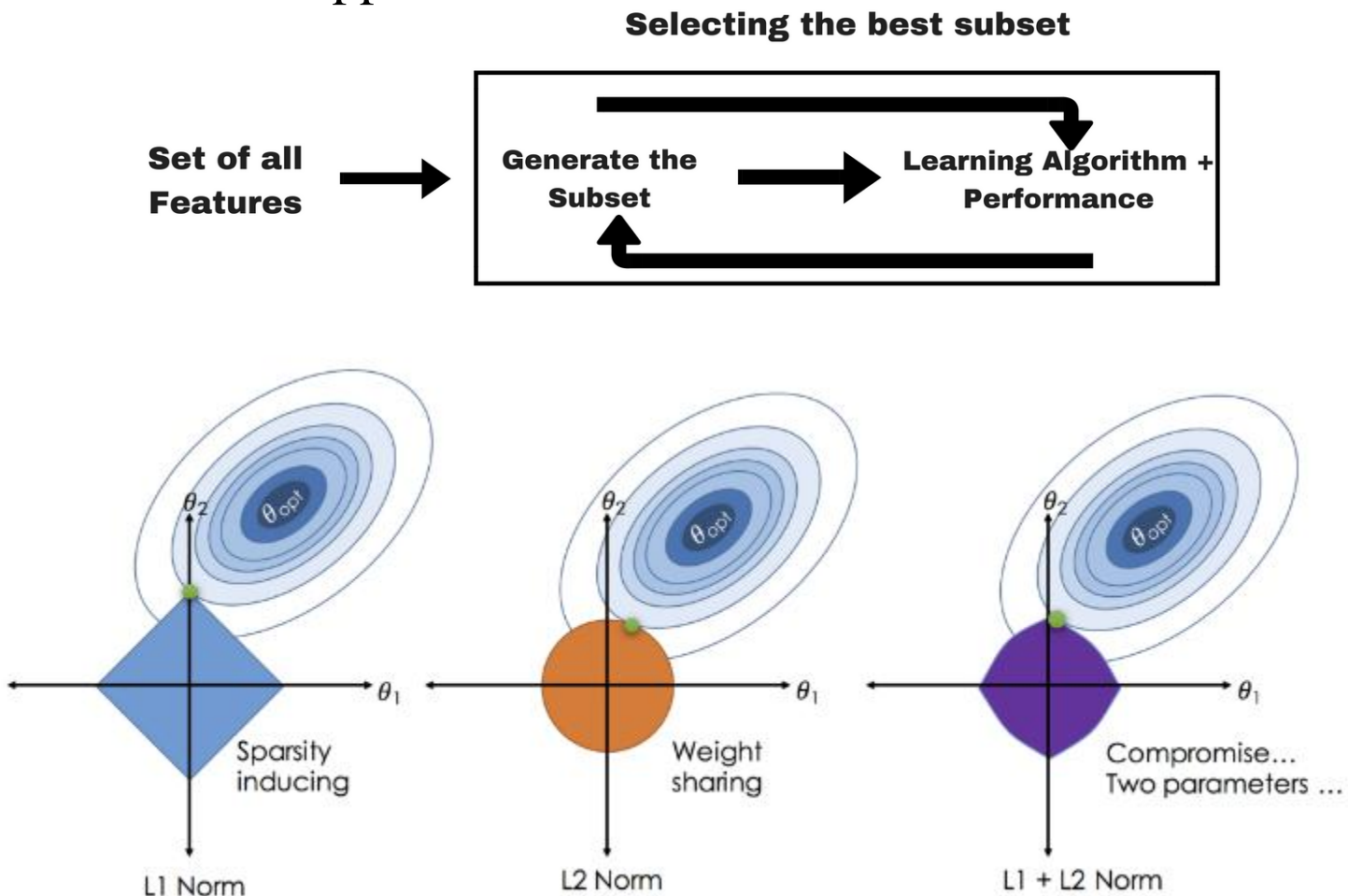
Model Selection - 4

- Wrap Mode (Select from Model): Consider the selection of a set of features as a **search problem**, where different combinations are prepared, evaluated and compared to other combinations.



Model Selection - 5

- ▶ **Embedded Mode:** A catch-all group of techniques which perform feature selection as part of the model construction process. Embedded methods combines the advantageous aspects of both Filter and Wrapper methods.



Shrinkage

- Choose the β that minimize

$$\mathcal{S}(\beta, \lambda) = (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) + \lambda \|\beta\|^2$$

for some choice of $\lambda > 0$.

- Ridge regression minimize \mathcal{S} under ℓ^2 norm

$$\|\beta\|_{\ell^2} = \sum_{i=1}^p \beta_i^2$$

Lasso regression minimize \mathcal{S} under ℓ^1 norm

$$\|\beta\|_{\ell^1} = \sum_{i=1}^p |\beta_i|$$

- One can show Ridge regression has analytic solution of β

$$\hat{\beta}_{\ell^2} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

Lasso regression has no close form formula.

- As $\lambda \searrow 0$, the Ridge and Lasso estimate become the OLS estimate
- As $\lambda \nearrow \infty$, the Ridge and Lasso estimate shrink to 0.