

Copyright © 2013 John Smith

PUBLISHED BY PUBLISHER

BOOK-WEBSITE.COM

Licensed under the Creative Commons Attribution-NonCommercial 3.0 Unported License (the "License"). You may not use this file except in compliance with the License. You may obtain a copy of the License at http://creativecommons.org/licenses/by-nc/3.0. Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License.

First printing, March 2013



- 1	Part One				
1	Linear Regression	. 7			
1.1	Projection in Euclidean Space	7			
1.2	Cochran's Theorem	14			
1.3	Gaussian Linear Regresson Model	15			
1.4	Statistical Inference for eta , σ^2	18			
1.5	Delete One Prediciton	19			
1.6	Residuals	20			
1.7	Influence and Cook's Distance	21			
1.8	Orthogonal Decomposition	22			
1.9	Lack of Fit Test	23			
2	ANOVA (1-way)	25			
2.1	Overview	25			
3	Mutiway ANOVA	27			
3.1	Overview	27			
4	Nonorthogonal Design	29			
4.1	Overview	29			

5 5 .1	Random Effects Model Overview	. <mark>31</mark> 31
Ш	Part Two	
6 6.1	Basic Concepts Overview	. 35 35
7 7.1	Estimation	. <mark>37</mark> 37
8 8.1	Inference	. 39 39
9 9.1	Residuals	. 41 41
10 10.1	Cetegorical Prediction	. 43 43
11 11.1	Some Important GLM Overview	. 45 45
12 12.1	Multivariate GLM Overview	. 47 47
Ш	Part Three	
13 13.1	Principle Componant Analysis	. <mark>51</mark> 51
14 14.1	Canonical Correlation Analysis	. 53
15 15.1	Independent Componant Analysis	. 55 55
	Index	. 57

Part One

1	Linear Regression 7
1.1	Projection in Euclidean Space
1.2	Cochran's Theorem
1.3	Gaussian Linear Regresson Model
1.4	Statistical Inference for β , σ^2
1.5	Delete One Prediciton
1.6	Residuals
1.7	Influence and Cook's Distance
1.8	Orthogonal Decomposition
1.9	Lack of Fit Test
2	ANOVA (1-way)
2.1	Overview
۷. ۱	CVCIVICW
2	M !! ANOVA
3	Mutiway ANOVA
3.1	Overview
4	Nonorthogonal Design
4.1	Overview
5	Random Effects Model
5.1	Overview

1. Linear Regression

- projection
- orthongonal decomposition
- Gaussian Linear Regression
- prediction (generally of \hat{y})
- different types of errors
- influence
- lack of fit
- \bullet R^2
- Multicollinearity

1.1 Projection in Euclidean Space

Monday August 22

Definition 1.1.1 — Euclidean Space. One way to think of the Euclidean plane is as a set of points satisfying certain relationships, expressible in terms of distance and angle. **Euclidean space** is an abstraction detached from actual physical locations, specific reference frames, measurement instruments, and so on.

Let Euclidian Space be denoted by \mathbb{R}^{P} .

$$\mathbb{R}X \dots X\mathbb{R} = \{(x_1, \dots, x_p) : x_1 \in \mathbb{R} \dots, x_p \in \mathbb{R}^P\}$$

Definition 1.1.2 — **Inner Product.** In linear algebra, an inner product space is a vector space with an additional structure called an inner product. This additional structure associates each pair of vectors in the space with a scalar quantity known as the inner product of the vectors. **Inner products** allow the rigorous introduction of intuitive geometrical notions such as the length of a vector or the angle between two vectors. They also provide the means of defining orthogonality between vectors (zero inner product).

Let $a \in \mathbb{R}^P, b \in \mathbb{R}^P$

$$a^T b = \sum_{i=1}^P a_i b_i$$

$$a^T b = \langle a, b \rangle$$

Definition 1.1.3 — **Hilbert Space**. The mathematical concept of a Hilbert space generalizes the notion of Euclidean space. It extends the methods of vector algebra and calculus from the two-dimensional Euclidean plane and three-dimensional space to spaces with any finite or infinite number of dimensions. A Hilbert space is an abstract vector space possessing the structure of an inner product that allows length and angle to be measured. Furthermore, Hilbert spaces are complete: there are enough limits in the space to allow the techniques of calculus to be used.

Hilbert Inner Product Space $\{\mathbb{R}^P, \langle a, b \rangle\}$

General Inner Product

Let $\Sigma \in \mathbb{R}^{P_X P}$ set of all $P_X P$ matrices. Assume Σ is a positive definite matrix.

$$x^T \Sigma x < 0$$
$$\forall x \in \mathbb{R}^P$$

 $x \neq 0$

Then $a^T \Sigma b$ also satisfies the conditions for inner product.

$$a^T \Sigma b = \langle a, b \rangle_{\Sigma}$$

$$a^T b = a^T I b = \langle a, b \rangle_I$$

 $\{\mathbb{R}^P, <, >_{\Sigma}\}$ is a more general inner product space.

Linear Transformation

A matrix, $A \in \mathbb{R}^{PxP}$ can be viewed as linear transformation $T_A : \mathbb{R}^P \to \mathbb{R}^P, x \mapsto Ax$



Bing Li will denote T_A as A.

- \rightarrow means maps to for a domain.
- \mapsto means maps to for a value.
- \Rightarrow means implies.

If $A: \mathbb{R}^P \to \mathbb{R}^P$.

$$ker(A) = \{x \in \mathbb{R}^P, Ax = 0\}$$

 $ran(A) = \{Ax : x \in \mathbb{R}^P\}$

Definition 1.1.4 — Kernel. In linear algebra, the kernel, or sometimes the null space, is the set of all elements v of V for which L(v) = 0, where 0 denotes the zero vector in W.

In coordinate plane, think of a function that crosses the x-axis. The kernel would be all points on x where y = 0.

Definition 1.1.5 — Range. In coordinate plane, how much of the y axis is reached with the function? Now extend this idea to more dimensions.

A linear transformation is **idempotent** if

$$A = A^{2}$$
$$Ax = A(A(x))$$
$$\forall x \in \mathbb{R}^{P}$$

If A were a number it could only be 1 or 0.

Wednesday August 24

Let $T \in \mathbb{R}^{PxP}$ then there exists a unique operator $R \in \mathbb{R}^{PxP}$ such that $\forall x, y \in \mathbb{R}^{P}$,

$$\langle x, Ty \rangle = \langle Rx, y \rangle$$

(general inner product, $a^T \Sigma b$). Aside: What this states is that if you give me any operator in the first you can find one in the second.

R is called the **adjoint operator** of T. Written as T^* , that is,

$$\langle x, Ty \rangle = \langle T^*x, y \rangle$$

Derived Facts

$$< x, Ty > = < T^*, y >$$

= $< y, T^*x >$
= $< (T^*)^*y, x >$
= $< x, (T^*)^*y >$

(by the definition)
(inner products the order doesn't matter)
(Use the definition again)
(swap order)

So,
$$T = (T^*)^*$$
.

It is easy to see in our case

$$\langle x, Ty \rangle_{\Sigma} = x^{T} \Sigma Ty$$

$$= x^{T} \Sigma T \Sigma^{-1} \Sigma y$$

$$= (\Sigma^{-1} T^{T} \Sigma x)^{T} \Sigma y$$

$$= \langle \Sigma^{-1} T^{T} \Sigma x, y \rangle_{\Sigma}$$

So, $T^* = \Sigma^{-1}T^T\Sigma$ when $\Sigma = I_P$ (identity) and $T^* = T^T$.

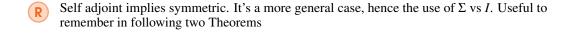
Derived Facts

An operator is **self adjoint** if its adjoint is itself. (i.e. if $T = T^*$ or $\langle x, Ty \rangle = \langle Tx, y \rangle$). In the case of $<,>_{\Sigma}$,

$$T = \Sigma^{-1} T^T \Sigma$$

if

$$\Sigma = I_P$$
, $T = T^T$



Theorem 1.1.1 If $A \in \mathbb{R}^{PxP}$ is symmetric, then there exists **eigenvalue-eigenvector pairs**. $(\lambda_1, v_1), \dots (\lambda_P, v_P)$ such that $v_1 \perp \dots \perp v_P$. Orthoginal basis (ONB) such that

$$A = \sum_{i=1}^{P} \lambda_i v_i v_i^T \text{(spectral decomposition)}$$

More generally, if A is a linear operator in \mathcal{H} (finite dimential inner product such as $(\mathbb{R}^P,<,>_{\Sigma})$). its eigen pair (linear operator now) (λ,ν) is defined by

$$\begin{cases} Av = \lambda v \\ \langle v, v \rangle = 1 \end{cases}$$

Definition 1.1.6 — Orthogonal Basis. In the following, $(\mathbb{R}^P, <, >_{\Sigma}) = \mathcal{H}$ (H for Hilbert) ONB is defined by:

- 1. $v_i \perp v_j, \langle v_i, v_j \rangle = 0$ 2. $||v_i|| = 1$ 3. $span\{v_1, \dots, v_P\} = \mathcal{H}$

Theorem 1.1.2 Suppose $A: \mathcal{H} \to \mathcal{H}$ is a self adjoint linear operator. Then A has eigen pairs: $(\lambda_1, \nu_1, \dots, (\lambda_P, \nu_P))$ where $\{\nu_1, \dots, \nu_P\}$ is ONB of \mathbb{R} such that

$$A = \sum_{i=1}^{P} \lambda_i v_i v_i^T \Sigma$$

Proof. (λ, v) is eigen pair of A, which means

$$Av = \lambda v$$

$$< v, v > = 1$$

$$v^T \Sigma v = 1$$

Let $u = \sum_{i=1}^{n} v_i$.

Aside: $\Sigma^{\alpha} = \Sigma \lambda_i^{\alpha} v_i v_i^T$

Let
$$v = \Sigma^{-\frac{1}{2}}u$$
.

$$A\Sigma^{-\frac{1}{2}}u = \lambda\Sigma^{-\frac{1}{2}}u$$
$$\Sigma^{-\frac{1}{2}}u = \lambda u$$

So, (λ, ν) is an eigen pair of A in $(\mathbb{R}, <, >_{\Sigma}) \Leftrightarrow (\lambda, u)$ '...' of $\Sigma^{\frac{1}{2}} A \Sigma^{-\frac{1}{2}}$ in $(\mathbb{R}, <, >_{I})$. Note that, A is self adjoint in $(\mathbb{R}, <, >_{\Sigma})$. So, $A = \Sigma^{-1} A^{T} \Sigma$

$$\begin{split} \Sigma^{\frac{1}{2}} A \Sigma^{-\frac{1}{2}} &= \Sigma^{\frac{1}{2}} A^T \Sigma \Sigma^{-\frac{1}{2}} \\ &= \Sigma^{-\frac{1}{2}} A \Sigma^{\frac{1}{2}} \\ &= (\Sigma^{\frac{1}{2}} A \Sigma^{-\frac{1}{2}})^T \end{split}$$

Note: $\Sigma^{\frac{1}{2}}A\Sigma^{-\frac{1}{2}}$ is symmetric!! So by Theorem 1.1, $\Sigma^{\frac{1}{2}}A\Sigma^{-\frac{1}{2}} = \sum \lambda_i v_i v_i^T$ where (λ_i, v_i) eigenpairs of $\Sigma^{\frac{1}{2}}A\Sigma^{-\frac{1}{2}}$.

That means $(\lambda_i, \Sigma^{\frac{1}{2}} v_i)$ are eigen pairs of A.

So,
$$\Sigma^{\frac{1}{2}}A\Sigma^{-\frac{1}{2}} = \sum_{i=1}^{P} \Sigma^{\frac{1}{2}} u_i u_i^T \Sigma^{\frac{1}{2}} \Rightarrow A = \sum_{i=1}^{P} \lambda u_i u_i^T \Sigma$$

Definition 1.1.7 — Projection. If P is an operator in $(\mathbb{R}^P, <, >)$ then P is called a **projection** if it is both idempotent $(P = P^2)$ and self adjoint $(P = P^*)$.

Preposition 1.1 If *A* is a linear operator then $ker(A) = ran(A^*)^{\perp}$

Proof. Take
$$x \in ker(A) (\Rightarrow Ax = 0)$$
.
 $\forall y \in ran(A^*), x \perp y$
 $\Rightarrow x \perp y \forall y = A^*z, z \in \mathbb{R}^P$
Hence,

$$\langle x, y \rangle = \langle x, A^*z \rangle$$

$$= \langle Ax, z \rangle$$

$$= \langle 0, z \rangle$$

$$= 0$$

$$\Rightarrow x \perp y$$

$$\Rightarrow x \in ran(A^*)^{\perp}$$

Or vice versa.

Friday August 26

ightharpoonup \perp means orthogonal complement.

$$\mathscr{S}^{\perp} = \{ v \in \mathbb{R}^P, v \perp \mathscr{S} \}$$

$$v \perp w \forall w \in \mathscr{S}$$

$$< v, w > = 0 \forall w \in \mathcal{S}$$

= $\{v \in \mathbb{R}^P, < v, w > = 0 \forall w \in \mathcal{S}\}$

Recall,
$$ker(A) = ran(A^*)^{\perp}$$

So, if A is slef adjoint then this is true and ran(A) is also span(A) which is the subspace spanned all columns of A.

Theorem 1.1.3 If P is a projection, then

- 1. $Pv = v, \forall v \in ran(P)$
- 2. Pv = 0, $\forall v \perp ran(P)$
- 3. If Q is another projections such that the ran(Q) = ran(P) then Q = P. (The range determines the operator, because it is what decomposes the operator.)

Asside: P acts like one on some spaces, and zero on orthogonal space.

Proof. 1. Let
$$v \in ran(P)$$
. Since $P^2 = P$ (idempotent) then
$$P^2v = Pv$$

$$\Rightarrow P^2v - PV = 0$$

$$\Rightarrow P(Pv - v) = 0$$

$$\Rightarrow Pv - v \in ker(P)$$

$$\Rightarrow Pv - v \perp ran(P)$$

$$\Rightarrow < Pv - v, Pv - v >= 0$$

$$\Rightarrow ||Pv - v|| = 0$$

$$\Rightarrow Pv - v = 0$$

$$\Rightarrow Pv = v$$
2. If
$$v \perp ran(P)$$

$$\Rightarrow v \in ker(P)$$

$$\Rightarrow Pv = 0$$
3. If Q is another operator with $ran(Q) = ran(P) = \mathscr{S}$ then $\forall v \in \mathscr{S}$

$$Qv = v = Pf(\forall v \perp \mathscr{S})$$

$$Qv = 0 = Pv$$

$$Qv = Pv \forall, v \in \mathscr{S}$$

$$Q = P$$

Theorem 1.1.4 Suppose \mathscr{S} is a subspace of \mathbb{R}^P , R v_1, \ldots, v_m is a basis of \mathscr{S} .

Let
$$v = (v_1, \ldots, v + m) \in \mathbb{R}^{xM}$$
.

Then,

1. $A = v(v^T \Sigma v)^{-1} v^T \Sigma$ is a projection.

2. $ran(A) = \mathcal{S}$

Proof. 1. idempotent.
$$A^{2} = v(v^{T}\Sigma v)^{-1}v^{t}\Sigma v(v^{T}\Sigma v)^{-1}v^{T}\Sigma$$

$$= v(v^{T}\Sigma v)^{-1}v^{T}\Sigma$$

$$= A$$

2. Self adjoint.

Let
$$x, y \in \mathbb{R}^P$$

 $\langle x, Ay \rangle = x^T \sum v (v^T \sum v)^{-1} v^{\Sigma} y$
 $= (v(v^T \sum v)^{-1} v^T \sum x)^T \sum y$
 $= \langle Ax, y \rangle$

3. $ran(A) = \mathcal{S}$?

Let $x \in \mathbb{R}^P$.

$$Ax = v(v^T \Sigma v)^{-1} v^T \Sigma x \in span(v) = \mathscr{S}$$

So let $x \in \mathcal{S}$,

$$x \in ran(v)$$

$$x = vy$$

for some $y \in \mathbb{R}^P$

$$= v(v^T \Sigma v)^{-1} v^T \Sigma v y$$

 $\in ran(A)$

So, $\mathscr{S} \subseteq ran(A)$ and then $\mathscr{S} = ran(A)$.

We write *A* as $P_{\mathscr{S}}(\Sigma)$ (orthogonal projection on to \mathscr{S} with respect to Σ - product).

In the following, let $I : \mathbb{R}^P \to \mathbb{R}^P$ be the identity mapping. $(x \mapsto x)$ Let \mathscr{S} be a subspace in \mathbb{R}^P .

Let
$$Q_{\mathcal{S}}(\Sigma) = I - P_{\mathcal{S}}(\Sigma)$$

Proprosition 1.2
$$Q_{\mathscr{S}}(\Sigma) = P_{\mathscr{S}^{\perp}}(\Sigma)$$

Proof. Show $Q_{\mathcal{S}}(\Sigma)$ is projection.

1. Idempotent

$$\begin{aligned} Q_{\mathscr{S}}^{2}(\Sigma) &= Q_{\mathscr{S}}(\Sigma)Q_{\mathscr{S}}(\Sigma) \\ &= (I - P_{\mathscr{S}}(\Sigma))(I - P_{\mathscr{S}}(\Sigma)) \\ &= I - P_{\mathscr{S}}(\Sigma) - P_{\mathscr{S}}(\Sigma) + P_{\mathscr{S}}P_{\mathscr{S}} \\ &= Q_{\mathscr{S}}(\Sigma) \end{aligned}$$

2. Self-adjoint

$$x, y \in \mathbb{R}^P$$

3. Range

$$ran(Q_{\mathscr{S}}(\Sigma)) = \mathscr{S}^{\perp}$$
. Take $x \perp \mathscr{S} = ran(P_{\mathscr{S}}(\Sigma))^{\perp} = ker(P_{\mathscr{S}}(\Sigma))$.

$$\Rightarrow P_{\mathscr{S}}(\Sigma) = 0$$

$$\Rightarrow Q_{\mathscr{S}}(\Sigma)x = x - P_{\mathscr{S}}(\Sigma)x = x$$

$$X \in ran(Q_{\mathscr{S}}(\Sigma))$$

$$\Rightarrow \mathscr{S}^{\perp} \subseteq ran(Q_{\mathscr{S}}(\Sigma))$$
Take $x \in ran(Q_{\mathscr{S}}(\Sigma))$, $\forall y \in \mathscr{S} = ran(P_{\mathscr{S}}(\Sigma))$

$$y = P_{\mathscr{S}}(\Sigma)z \text{ for some } z \in \mathbb{R}^{P}$$

$$< x, y > = < x, P_{\mathscr{S}}(\Sigma)z > = < P_{\mathscr{S}}(\Sigma)x, z > = 0$$

$$\Rightarrow x \in \mathscr{S}^{\perp}$$

$$\Rightarrow ran(Q_{\mathscr{S}}(\Sigma)) = \mathscr{S}^{\perp}$$

1.2 Cochran's Theorem

This section will be about the distribution of the squared norm of a projection of a Gaussian random vector.

Preposition 1.3 If A is idempotent, then its eigenvalues are either 0 or 1.

Proof. λ is eigenvalue of A.

$$\Rightarrow Av = \lambda v(||v|| = 1)$$

$$\lambda = Av = A^2v = \lambda Av = \lambda^2$$

So, λ is 0 or 1.

Monday August 29

Lemma 1.1 Suppose $V \sim N(0, \sigma^2 I_P)$.

P is projection with I_P - inner product. Then $V^T P V \sim \sigma^2 \chi_S^2$ where df = rank(P).

Proof. P is symmetric, and it has spectral decomposisition,

$$ARA^{T}$$

where the A's are orthogonal and R is diagonal with diagonal entries 0 or 1.

Then,

$$A^T V \sim N_P(0, A^T(\sigma^2 I_P)A) = N_P(0, \sigma I_P)$$

Let,

$$Z = RA^T V$$

then,

$$Z \sim N_P(0, \sigma^2 R^2) = N_P(0, \sigma^2 R)$$

That means among the components of Z, some are distributied as N(0, 1) and the rest are zero and they are independant. So,

$$Z^T Z \sim \chi_S^2 = V^T P V$$

Corollary 1.1 Suppose $X \sim N(0, \Sigma)$. Consider the Hilbert space $(\mathbb{R}^P, <, >_{\Sigma^{-1}})$.

$$\langle a,b\rangle_{\Sigma^{-1}}=a^T\Sigma^{-1}b$$

Let $\mathscr S$ be a subspace of $\mathbb R^P$ and $P_{\mathscr S}(\sigma^{-1})$ be the projection onto $\mathscr S$ with respect to $<,>_{\Sigma}^{-1}$ (special case of Fisher information inner product)

Then,

$$||P_{\mathscr{S}}(\Sigma^{-1})x||_{\Sigma^{-1}}^2 \sim \chi_r^2$$

where $r = dim(\mathcal{S})$.

Proof. Let V be a basis matrix of \mathscr{S} (i.e. the col of V form basis in \mathscr{S}).

$$\begin{aligned} ||P_{\mathscr{S}}(\Sigma^{-1})x||_{\Sigma^{-1}}^{2} &= < P_{\mathscr{S}}(\Sigma^{-1})x, P_{\mathscr{S}}(\Sigma^{-1})x > \\ &= x^{T} P_{\mathscr{S}}(\Sigma^{-1}) \Sigma^{-1} P_{\mathscr{S}}(\Sigma^{-1})x \\ &= x^{T} (V(V^{T}\Sigma^{-1}v)^{-1}v^{T}\Sigma^{-1})^{T} \Sigma^{-1} (V(V^{T}\Sigma^{-1}v)^{-1}v^{T}\Sigma^{-1})^{T}x \\ &= x^{T} \Sigma^{-1} V(V^{T}\Sigma^{-1}v)^{-1}v^{T}\Sigma^{-1}V(V^{T}\Sigma^{-1}v)^{-1}v^{T}\Sigma^{-1})^{T}x \\ &= (\Sigma^{-\frac{1}{2}}x)^{T} || \end{aligned}$$

But,

$$\Sigma^{-\frac{1}{2}}x \sim N(0, I_P)$$

$$\Sigma^{-\frac{1}{2}}V(V^T\Sigma^{-1}v)^{-1}v^T\Sigma^{-\frac{1}{2}})$$

is a projection with repect to I_P -inner producted (idempotent, self adjoint, YES). By Lemme 1.1, $\sim \chi_r^2$.

It is then easy to derive Cocharan's Theorem. (see proof in Homework 1)

Theorem 1.2.1 Let $X \sim N(0, \Sigma)$ and $\mathcal{H} = \{\mathbb{R}^P, <, >_{\Sigma^{-1}}\}$. Let $\mathcal{S}_1, dots, \mathcal{S}_k$ be linear subspaces of \mathbb{R}^P such that $\mathcal{S}_i \perp \mathcal{S}_j$ in $<, >_{\Sigma^{-1}}$

Let $r_i = dim(\mathcal{S}_i)$.

Let
$$w_i = ||P_{\mathcal{S}_i}(\Sigma^{-1})X||_{\Sigma^{-1}}^2$$

Then,

- 1. $W_i \sim \chi_{r_i}^2$
- 2. $W_1 \perp \!\!\! \perp, \dots, \perp \!\!\! \perp W_k$ where $\perp \!\!\! \perp$ indicates independence.

1.3 Gaussian Linear Regresson Model

$$Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

$$X = \begin{pmatrix} x_{11} & \dots & x_{1P} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{np} \end{pmatrix} \in \mathbb{R}^{nxp}$$

Consider the linear model,

$$y = X\beta + \varepsilon$$

$$\varepsilon \sim N(0, \sigma^2 I_n)$$

where X has full comlumn rank $(n \ge p)$.

Here X is treated as fixed

• Maximum Likelihood Estimator

$$E(y) = X\beta \in \mathbb{R}^n$$

$$Var(y) = \sigma^2 I_n$$

$$y \sim N_p(X\beta, \sigma^2 I_n)$$

Multivariate normal density

$$y \sim N(\mu, \Sigma)$$

$$f_Y(y) = \frac{1}{(2\pi)^{\frac{n}{2}} [det(\Sigma)]^{\frac{1}{2}}} \exp(-\frac{1}{2}(y-\mu)^T \Sigma^{-1}(y-\mu))$$

In our case,

$$\Sigma = \sigma I_n$$

$$det(\Sigma) = det(\sigma^2 I_n) = \sigma^2 det(I_n) = \sigma^{2n}$$

50,

$$f_Y(y) = \frac{1}{(2\pi)^{\frac{n}{2}} \sqrt{\sigma^{2n}}} \exp(-\frac{1}{2\sigma^2} ||y - \mu||^2)$$

$$\log(f_{y}(\eta)) = \frac{n}{2}\log(\sigma^{2}) - \frac{1}{2\sigma^{2}}||y - \mu||^{2} = \ell(\beta, \sigma^{2}, y)$$

$$\frac{\partial}{\partial \beta} = \dots = -\frac{1}{2\sigma^2} 2X^T (y - X\beta) = 0$$

$$\hat{\boldsymbol{\beta}} = (\boldsymbol{X}^T \boldsymbol{x})^{-1} \boldsymbol{X}^T \boldsymbol{Y} \in \mathbb{R}^P$$

$$\frac{\partial}{\partial \sigma^2} l(\beta, \sigma^2, y) = \dots = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} ||y - X\beta||^2 = 0$$

$$\hat{\sigma^2} = \frac{1}{n} ||y - X\hat{\beta}||^2$$

In summary, the MLE for (β, σ^2) in Gaussian Linear Model are

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

$$\hat{\sigma^2} = \frac{1}{n} ||y - X\hat{\beta}||^2$$

Note that

$$X\hat{\beta} = X(X^TX)^{-1}X^Ty = \hat{y}$$

So,
$$\hat{y} = P_{span(x)}(I_P) = P_X$$
.
Now,

$$\hat{\sigma}^2 = \frac{1}{n} ||(I_n - P_X)y||^2 = \frac{1}{n} ||Q_X y||^2$$

where $(I_n - P_X)$ is projection on to $span(X)^{\perp}$.

It turns out that $(X^T y, y^T y)$ is complete, sufficient statistic for this Gaussian linear model.

Wednesday August 31

Recall,

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

$$\hat{\sigma^2} = \frac{1}{n} ||y - X \hat{\beta}||^2$$

$$Q_x = I_n - P_x$$

$$P_X + X (X^T X)^{-1} X^T$$

Several properties,

$$E(\hat{\boldsymbol{\beta}}) = B \text{ (unbiased)}$$

$$Var(\hat{\boldsymbol{\beta}}) = (X^T X)^{-1} X^T \sigma^2 I_n X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1}$$

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

Because P_x has rank p and Q_x has rank (n-p), then

$$||Q_x y||^2 \sim \chi^2_{(n-p)}$$

Let's find an unbiased estimator for σ^2

$$E(\hat{\sigma^2}) = E(\frac{1}{n}||Q_x y||^2)$$
$$= \frac{n-p}{n}\sigma^2$$
$$E(\frac{n}{n-p}\hat{\sigma^2}) = \sigma^2$$

Moreover, $\hat{\beta}$ has one-to-one transformation with

$$(X^T X)^{-1} X^t y \leftrightarrow X (X^T X)^{-1} X^t y = P_{xy}$$

$$Cov(P_{Xy}, Q_{Xy}) = P_X \sigma^2 I_n Q_X$$

= $\sigma^2 P_X Q_X$
= 0

$$P_{Xy} \perp \!\!\!\perp Q_{Xy}$$
 (due to normality)

$$\hat{\beta} \leftrightarrow P_{Xy}$$

 $\hat{\sigma}^2$ is a funciton of Q_{Xy} , so $\hat{\beta} \perp \!\!\! \perp \hat{\sigma}^2$

In your homework, $\hat{\beta}$, $\hat{sigma}^2 \leftrightarrow$ complete sufficient.

 $\hat{\beta}$, $\tilde{sigma^2}$ is UMVUE (Lehmann-Sheffe).

Theorem 1.3.1 — Gaussian Regression Model. Under this model:

1.
$$\hat{\beta}$$
, $\tilde{\sigma}^2$ UMVUE for β , σ^2
2. $\hat{\beta} \sim N(\beta, \sigma^2(X^TX)^{-1})$
3. $(n-p)\tilde{\sigma}^2 \sim \sigma^2 \chi^2_{(n-p)}$

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

3.
$$(n-p)\tilde{\sigma}^2 \sim \sigma^2 \chi^2_{(n-p)}$$

4.
$$\hat{\beta} \perp \tilde{\sigma}^2$$

Statistical Inference for β , σ^2

Suppose we want to test

$$H_0: \beta_1 = \beta_{i0}$$

Let $M = (X^T X)^{-1}$.

Then,

$$\hat{\beta} \sim N(\beta_i 0, \sigma^2 M_{ii})$$

where, $M_{ii} \leftarrow (i, i)^{th}$ entry of M

Also,
$$\frac{(n-p)\tilde{\sigma^2}}{\sigma^2} \sim \chi^2_{(n-p)}$$

$$\hat{eta}$$
 \perp $\tilde{\sigma^2}$

$$\frac{\frac{\hat{\beta_i} - \beta_{i0}}{\sqrt{\sigma^2 M_{ii}}} \sim N(0,1)}{\sqrt{\frac{(n-p)\sigma^2 \tilde{/}\sigma^2}{n-p}}} \sim t_{(n-p)}$$

$$T = \frac{\hat{\beta}_i - \beta_{i0}}{\sigma \tilde{\sqrt{M}_{ii}}} \sim t_{(n-p)} = (*)$$

Reject H_0 if

$$|\frac{\hat{eta}_i - eta_{i0}}{\sigma \sqrt{\tilde{M}_{ii}}}| > t_{\frac{lpha}{2}(n-p)}$$

$$X \sim N(\mu, 1) \ y \sim \chi_r^2 \ X \perp \!\!\! \perp y$$

 $\frac{X}{\sqrt{\frac{y}{r}}} \sim t_n(\mu)$

Power at β_{i1}

$$\hat{\beta}_i \sim N(\beta_{i1}, \sigma^2 M_{i1})$$

So,

$$rac{\hat{eta}_i - eta_{i0}}{ ilde{\sigma}\sqrt{M_{ii}}} \sim t_{(n-p)} (rac{eta_{i1} - eta_{i0}}{\sigma\sqrt{M_{ii}}})$$

(alternative distrabution of T)

By this (*),

$$P(\in (-t_{\frac{\alpha}{2}(n-p)}, t_{\frac{\alpha}{2}(n-p)}))$$

Convert this to put β_{i0} in between $(1 - \alpha)100$ percent C.I. for $\beta_{i.}$.

$$(\hat{eta_1} - t_{rac{n}{2}(n-p)}\hat{oldsymbol{\sigma}}\sqrt{M_{ii}},\hat{eta_1} + t_{rac{n}{2}(n-p)}\hat{oldsymbol{\sigma}}\sqrt{M_{ii}})$$

1.5 Delete One Prediciton

Very useful in variable selection, cross validation, diagnostics.

Prediction: $\hat{y} = X\hat{\beta} = P_x y$

But this has a drawback as it favors overfitting. Projectioning onto larger spaces will always decrease the norm, $||Q_Xy||^2$. (This can decrease errors which would cause you to think it's better, even though it's not.)

To prevent overfitting, try to be objective, withhold y_i when predicting y_i (inverse of a matrix, rank 1 perpendicular)

Theorem 1.5.1 Suppose $A \in \mathbb{R}^{PxP}$ is a symmetric, nonsingular matrix. and $v \in \mathbb{R}^{P}$. Then,

$$(A \pm vv^T)^{-1} = A^{-1} \pm \frac{A^{-1}vv^tA^{-1}}{1 \pm v^TA^{-1}v}$$

Use what is left to compute $\hat{\beta}_{-i}$.

$$\hat{\beta}_{-i} = (X_{-1}^T X_{-i})^{-1} X_{-i}^T y_{-i}$$

This can be expanded in simple sum, so that you don't have to do n regressions.

$$(X_{-i}^{T}X_{-i})^{-1} = (X^{T}X - X_{i}X_{i}^{T})^{-1}$$

$$= A^{-1} + \frac{A^{-1}vv^{T}A^{-1}}{1 - v^{t}A^{-1}v}$$

$$= (X^{T}X)^{-1} + \frac{(X^{T}X)^{-1}X_{i}X_{i}^{T}(X^{T}X)^{-1}}{1 - X_{i}^{T}MX_{i}}$$

$$X_{i}^{T}MX_{i} = X_{i}^{T}(X^{T}X)^{-1}$$

$$= (P_{x})_{ii}$$

$$= P_{i}$$

$$\hat{\beta}_{i} = (X^{T}X - X_{i}X_{i}^{T})^{-1}(X^{T}y - X_{i}y_{i})$$

$$= [M + \frac{MX_{i}X_{i}^{T}M}{1 - P_{i}}](X^{T}y - X_{i}y_{i})$$

$$= MX^{T}y + \frac{MX_{i}X_{i}^{T}MX^{T}y}{1 - P_{i}} - MX_{i}y_{i} - \frac{MX_{i}X_{i}^{T}MX_{i}y_{i}}{1 - P_{i}}$$

$$= \dots$$

$$= \hat{\beta} - \frac{MX_{i}}{1 - P_{i}}(y_{i} - X_{i}^{T}\hat{\beta})$$

Delete-one regression.

$$X_i \hat{\beta_{-i}} = \hat{y}_i - \frac{P_i}{1 - P_i} (y_i - \hat{y}_i)$$

Friday September 2

Delete- one error

$$y_i - \hat{y}_i^{(-i)}$$

Recall, you want to leave out y^i so you don't overfit.

The above is equivalent to

$$\begin{aligned} y_{i} - X_{i}^{T} \hat{\beta}_{-i} \\ y_{i} - \hat{y}_{i} - \frac{P_{i}}{1 - P_{i}} (y_{i} - \hat{y}_{i}) \\ (y_{i} - \hat{y}_{i}) (1 - \frac{P_{i}}{1 - P_{i}})) \\ \frac{1}{1 - P_{i}} (y_{i} - \hat{y}_{i}) \end{aligned}$$

Delete-one cross validation

$$\sum_{i=1}^{n} (y_i - \hat{y}_i^{(-i)})^2$$

This method is not affected by over fitting.

The following is often used for "tuning" or variable selection (i.e. penalty, bandwidth, regularization, etc). $\sum_{i=1}^{n} \frac{1}{(1-P_i)^2} (y_i - \hat{y}_i)^2$ Note: we will come back to variable selection later.

$$eta = egin{pmatrix} eta_1 \ dots \ eta_n \end{pmatrix} \ A \subseteq \{1,\dots,P\}$$

Cross validation of *A* minimizes over $A \in 2^{\{1,\dots,P\}}$. Best cross validation set.

1.6 Residuals

• Residual

$$\hat{e}_i = y_i - \hat{y}_i$$

· Standardized Residual

$$Var(\hat{e}_i) = Var(y_i - \hat{y}_i) = Var((Q_X)_{ii}y_i)$$

$$= ((Q_X)_{ii}y_i)\sigma^2$$

$$= (1 - P_i)\sigma^2$$

$$sd(\hat{e}_i) = \sqrt{1 - P_i}\sigma$$

$$\hat{sd}(\hat{e}_i) = \sqrt{1 - P_i}\tilde{\sigma}$$

$$\tilde{\sigma} = \frac{\sum_{i=1}^{n} (y_i - \hat{y}_i^{(-i)})^2}{n - p}$$

· Standardized residual

$$E_i^* = rac{\hat{e}_i}{ ilde{\sigma}\sqrt{1-P_i}}$$

• Prediction Error Sum of Squares (PRESS) Residual

$$y_i - \tilde{y}_i^{(-i)} = \frac{1}{1 - P_i} \hat{e}_i = \hat{e}_{iP}$$

$$\hat{e}_{iP} \sim N(0, \frac{\sigma^2}{1 - P_i})$$

• Standardized PRESS Error

$$\frac{\hat{e}_{iP}}{\tilde{\sigma}/\sqrt{1-i}} = \frac{\frac{1}{1-P_i}\hat{e}_i}{\tilde{\sigma}(\sqrt{1-P_i})} = \frac{\hat{e}_i}{\tilde{\sigma}(\sqrt{1-P_i})} = e_i^*$$

1.7 Influence and Cook's Distance

Definition 1.7.1 — Influence. The difference between predictions with and without a data point.

$$\hat{y}_i - \hat{y}_i^{(-i)}$$

$$\hat{y}_i - \hat{y}_i^{(-i)}$$

$$X_i\hat{\beta} - X_i\hat{\beta}_{-i}$$

Recall,

$$\hat{\beta}_{-i} - \hat{\beta} = -\frac{MX_i(y_i - \hat{y}_i)}{1 - P_i} = -\frac{MX_i\hat{e}_i}{1 - P_i}$$

$$||X_i\hat{\beta} - X_i\hat{\beta}_{-i}||^2 =$$

Cook's Distance (Technometrics, 1976?)

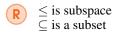
$$||\frac{\hat{y} - \hat{y}^{(-i)}||^2}{\tilde{\sigma}^2} = \frac{|i\hat{e}_i^2}{(1 - P_i)^2 \tilde{\sigma}^2}$$

Definition 1.7.2 — Cook's Distance. Cook's distance measures the influence of the i^{th} deservation.

Orthogonal Decomposition

Recall, \mathbb{R}^n is Euclidean Space.

 \mathcal{S} is a subspace $(\mathcal{S} < \mathbb{R}^n)$



$$\mathcal{S}_1 \leq \mathcal{S}_1 \mathcal{S}_2 \leq \mathcal{S}$$

$$\mathscr{S}_1 + \mathscr{S}_2 = \{x + y : x \in \mathscr{S}_1, y \in \mathscr{S}_2\}$$

Suppose
$$\mathcal{S}_1, \mathcal{S}_2 \leq \mathcal{S}$$
, $\mathcal{S}_1 + \mathcal{S}_2 = \mathcal{S}$, $\mathcal{S}_1 \perp \mathcal{S}_2$ then.

$$\{\mathscr{S}_1,\mathscr{S}_2\}$$

is called an orthogonal decomposition of $\mathscr S$ In this case,

$$\mathcal{S}_1 \oplus \mathcal{S}_2 = \mathcal{S}$$

More generally,

Definition 1.8.1 — Orthogonal Decomposition (O.D.). Let $\mathscr{S}_1, \ldots, \mathscr{S}_k$ be subspaces of \mathscr{S}

1.
$$\mathscr{S}_1, \ldots, \mathscr{S}_k = \{v_1 + \cdots + v_k : v_1 \in \mathscr{S}_1, \ldots, v_k \in \mathscr{S}_k\}$$

2.
$$\mathscr{S}_i \perp \mathscr{S}_j \quad \forall i \neq j$$

1. $\mathscr{S}_1,\ldots,\mathscr{S}_k=\{v_1+\cdots+v_k:v_1\in\mathscr{S}_1,\ldots,v_k\in\mathscr{S}_k\}$ 2. $\mathscr{S}_i\perp\mathscr{S}_j\quad \forall i\neq j$ Then, $\{\mathscr{S}_1,\mathscr{S}_2,\ldots,\mathscr{S}_k\}$ is an **orthogonal decomposition** of \mathscr{S} . We may write $\mathscr{S}=\mathscr{S}_1\oplus\mathscr{S}_2\oplus\cdots\oplus\mathscr{S}_k$.

Proposition 1.5 If $\mathcal{S}_1, \dots, \mathcal{S}_k$ is an O.D. of \mathcal{S} , then any $v \in \mathcal{S}$ can be uniquely written as

$$v_1 + \cdots + v_k$$

, where $v_1 \in \mathcal{S}_1, \dots v_k \in \mathcal{S}_k$.

Wednesday September 7

Definition 1.8.2 — Direct Difference. Let $\mathscr{S}_1 \leq \mathscr{S}_2 \leq \mathbb{R}^n$. Then,

$$\mathscr{S}_2 \cap \mathscr{S}_1^{\perp} \equiv \mathscr{S}_2 \ominus \mathscr{S}_1$$

is called direct difference. This is almost the same as orthogonal complement, except it is within \mathcal{S}_2 .

Proposition 1.6 If $\mathcal{S}_1 \leq \mathcal{S}_2$, then

$$\mathscr{S}_2 = \mathscr{S}_1 \oplus (\mathscr{S}_2 \ominus \mathscr{S}_1)$$

Proposition 1.7 - Orthogonal Decomposition and Projection Consider a Hilbert Space, $\mathcal{H} = \{\mathbb{R}^n, <, >_A\},$

1.9 Lack of Fit Test 23

1. If $\mathscr{S} \leq \mathscr{S}_1 \perp \mathscr{S}_2$ in \mathscr{H} , then

$$P_{\mathcal{S}_1}(A)P_{\mathcal{S}_2}(A)=0$$

2. If $\mathscr{S} \leq \mathscr{H}, \dots, \mathscr{S}_k \leq \mathscr{H}$, and $\mathscr{S}_1 \perp \dots \perp \mathscr{S}_k$, then

$$P_{\mathcal{S}_1,\oplus\cdots\oplus\mathcal{S}_k}(A) = P_{\mathcal{S}_1}(A) + \cdots + P_{\mathcal{S}_k}(A)$$

3. If $\mathcal{S}_1 \leq \mathcal{S}_2 \leq \mathbb{R}^n$, then

$$P_{\mathscr{S}_2 \ominus \mathscr{S}_1}(A) = P_{\mathscr{S}_2}(A) - P_{\mathscr{S}_1}(A)$$

Theorem 1.8.1 — Generalization of the earlier Cochran's Theorem. Suppose $X \sim N(0, \Sigma)$ where $\Sigma \in \mathbb{R}^{n \times n}$ is positive definite.

Let
$$\mathcal{H} = \{<,>_{\Sigma^{-1}}\}$$
. Suppose $\mathcal{S}_1,\ldots\mathcal{S}_k,\mathcal{S} \leq \mathcal{H}$ such that $\mathcal{S} = \mathcal{S}_1 \oplus \cdots \oplus \mathcal{S}_k$.

Let

$$w_i = ||P_{\mathcal{S}_i}(\Sigma^{-1})X||_{\Sigma^{-1}}^2$$
$$w = ||P_{\mathcal{S}}(\Sigma^{-1})X||_{\Sigma^{-1}}^2$$

Then,

- 1. $w = w_1 + \cdots + w_k$
- 2. $w_1 \!\!\perp \!\!\!\perp \dots \!\!\!\perp \!\!\!\!\perp \!\!\!\! w_k$
- 3. $w_i \sim \chi_{r_i}^2$ $w \sim \chi_r^2$

where r_i is the $dim(\mathcal{S}_i)$, r is the $dim(\mathcal{S})$, and $r = r_1 + \cdots + r_k$.

Notation 1.1. We use \oplus for spaces. We can also use \oplus function to stack up matrices. Let A_1, \ldots, A_k be matrices with arbitrary dimensions.

$$A_1 \oplus \cdots \oplus A_k = \begin{pmatrix} A_1 & \dots & 0 \\ & \ddots & \\ 0 & \dots & A_k \end{pmatrix}$$

1.9 Lack of Fit Test

Goodness of Fit

At each x_i you have multiple observations, say y_{i1}, \ldots, y_{im_i} . In this case, you may test to see if a linear model, $y_i = x_i^T \beta + \varepsilon_i$, is the correct choice for fitting the data. In general, lack of fit refers to testing whether any (linear, generalized, etc) model is adequately describing the data.

Denote

$$y_{i} = \begin{pmatrix} y_{i1} \\ \vdots \\ y_{im_{i}} \end{pmatrix}$$

$$1_{m_{i}} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

$$X = \begin{pmatrix} X_1^T \\ \vdots \\ X_m^T \end{pmatrix}$$
Assume

$$y_{ij} = X_i^T \beta + \varepsilon_{ij}$$

where $\varepsilon \sim^{iid} N(0, \sigma^2)$.

The point is that you have $y_{i1} \dots y_{jm}$ for each X_i .

In matrix form,

$$(1_{m_1} \oplus \cdots \oplus 1_{m_n})X\beta + \varepsilon$$

So, let *N* denote a full sample size.

$$N = m_1 + \cdots + m_n$$

this is a special case of linear model, except the design matrix is structured $(1_{m_1} \oplus \cdots \oplus 1_{m_n})X$ instead of X. So the formula for MLE (and so on) is the same.

$$X \leftrightarrow (1_{m_1} \oplus \cdots \oplus 1_{m_n})X$$

So,

$$\hat{\beta} = ([(1_{m_1} \oplus \cdots \oplus 1_{m_n})X])^T ([(1_{m_1} \oplus \cdots \oplus 1_{m_n})X])^{-1} [(1_{m_1} \oplus \cdots \oplus 1_{m_n})X]^T y$$

$$\hat{y} = (1_{m_1} \oplus \cdots \oplus 1_{m_n}) X \hat{\beta}
= (1_{m_1} \oplus \cdots \oplus 1_{m_n}) X ([(1_{m_1} \oplus \cdots \oplus 1_{m_n}) X])^T ([(1_{m_1} \oplus \cdots \oplus 1_{m_n}) X])^{-1} [(1_{m_1} \oplus \cdots \oplus 1_{m_n}) X]^T y
= (1_{m_1} \oplus \cdots \oplus 1_{m_n}) X [X^T \begin{pmatrix} m_1 & \dots & 0 \\ & \ddots & \\ 0 & \dots & m_n \end{pmatrix} X]^{-1} X^T (1_{m_1} \oplus \cdots \oplus 1_{m_n})$$

So, in linear model with replication we have our hypotheses for lack of fit test,

$$H_O: E(y_i) = 1_{m_i} X_i^T \beta$$

$$H_1: E(y_i) = 1_{m_i} \mu_i$$

We are testing whether the arbitrary means, $\mu_1, \dots \mu_n$ sit on the same line.



- General linear models
- Scheffe's simulteaneous confidence
- Singular decomposition
- Non Gaussian error



- Orthogonal design
- Additive 2 way ANOVA
- simultaneous intervals
- nonadditive
- decomposition of sum of squares
- Latin square
- nested design



$$\bullet \ \ \bar{X}_{\dot{i}} - \bar{X}_{\dot{i}}$$



Part Two

6 6.1	Basic Concepts Overview	35
7 7.1	Estimation	37
8 8.1	Inference Overview	39
9 9.1	Residuals Overview	41
10 10.1	Cetegorical Prediction Overview	43
11 11.1	Some Important GLM Overview	45
12 12.1	Multivariate GLM Overview	47







• deviance <-> sum of squares





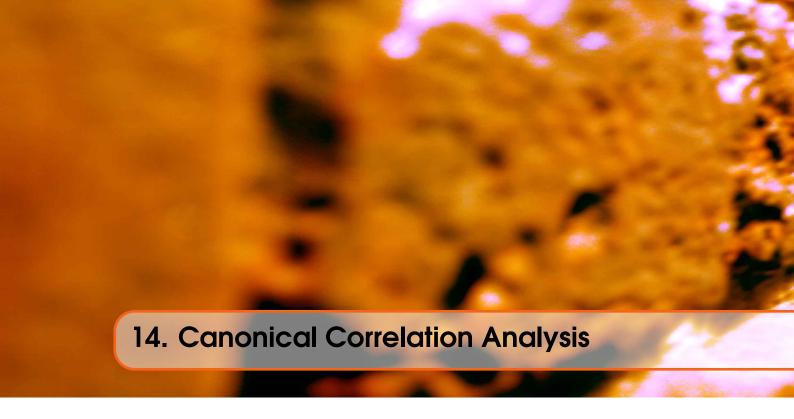




Part Three

13 13.1	Principle Componant Analysis Overview	51
14 14.1	Canonical Correlation Analysis Overview	53
15 15.1	Independent Componant Analysis Overview	55
	Index	57









Cochran's Theorem, 14

Delete One Prediciton, 19

Gaussian Linear Regresson Model, 15

Influence, Cook's Distance, 21

Lack of Fit Test, 23

Orthogonal Decomposition, 22 Overview, 25, 27, 29, 31, 35, 37, 39, 41, 43, 45, 47, 51, 53, 55

Projection, 7

Residuals, 20

Statistical Inference for β , σ^2 , 18