$$\Rightarrow \beta | X \sim N_p\left(\left(X^TX\right)^{-1}X^T\right)\left(X\beta\right), \left(X^TX\right)^{-1}X^T\sigma^2 I\left(\left(X^TX\right)^{-1}X^T\right)^T\right)$$

$$\Sigma = \sigma^2 \cdot \left(\left(X^TX\right)^{-1}X^T X \left(X^TX\right)^{-1}\right)$$

$$= N_p\left(\beta, \sigma^2\left(X^TX\right)^{-1}\right) := \frac{1}{\sqrt{(2\pi)^p |\sigma^2(X^TX)^{-1}|}} e^{-\frac{1}{2}(\beta-\beta)^T\left(\sigma^2(X^TX)^{-1}\right)^{-1}(\beta-\beta)}$$

$$\propto e^{-\frac{1}{2}\left((X^TX)^{-1}X^TY-\beta\right)^T \frac{1}{\sigma^2}(X^TX)\left((X^TX)^{-1}X^TY-\beta\right)}$$

≠h

why??

One of the crown achievements
of statistics ... the OLS
estimator is normal, unbiased, + a
whole bunch of other properties
which I would go into if I was teaching
a course on this...

But we care about something different ...

$$P\left(\beta | X, Y\right) \propto P\left(Y | \beta, X\right) P\left(\beta | X\right)$$

Bayesian estimation of linear
model with the 5 OLS assumptions...

Assume $\sigma^2$ known (nuisance parameter anyway)

$$P\left(\beta | X, y, \sigma^2\right) \propto P\left(Y | \beta, X, \sigma^2\right) P\left(\beta | X, \sigma^2\right)$$

let $P\left(\beta | X, \sigma^2\right) \propto 1$
no information

$$\underbrace{\frac{1}{\sqrt{(2\pi)^n |\sigma^2 I|}} e^{-\frac{1}{2}(y-X\beta)^T(\sigma^2 I)^{-1}(y-X\beta)}}(1)$$

$$\propto e^{-\frac{1}{2\sigma^2}(y^T-(X\beta)^T)(y-X\beta)}$$

$$\underset{\beta^T X^T}{\downarrow}$$

$$= e^{-\frac{1}{2\sigma^2}} \left( Y^T X (X^T X)^{-1} - \beta^T \right) \left( X^T X \right) \left( (X^T X)^{-1} X^T Y - \beta \right)$$

$$= e^{-\frac{1}{2\sigma^2}} \left( Y^T X \quad - \quad \beta^T X^T X \right) \left( (X^T X)^{-1} X^T Y - \beta \right)$$

FOIL
THIS
Now

$$= e^{-\frac{1}{2\sigma^2}} \left( Y^T X (X^T X)^{-1} X^T Y - \beta^T X^T Y - Y^T X \beta + \beta^T X^T X \beta \right)$$

$\underbrace{\quad}$ 1xp p×n n×1  by $n\times p$, $p\times n$
1x1      1x1

$\left( \beta^T X^T Y \right)^T = Y^T X \beta$

Since $(17.3)^T = 17.3$

i.e. the transpose of
a scalar equals itself
these two terms are equal

$$= e^{-\frac{1}{2\sigma^2}} \left( Y^T X (X^T X)^{-1} X^T Y - 2 Y^T X \beta + \beta^T X^T X \beta \right)$$

go back to place...

$$\propto e^{-\frac{1}{2\sigma^2}\left(y^Ty - \beta^TX^Ty - y^TX\beta + \beta^TX^TX\beta\right)} \propto e^{-\frac{1}{2\sigma^2}\left(-2y^TX\beta + \beta^TX^TX\beta\right)}$$

$\underbrace{\beta^TX^Ty}_{\substack{\text{depends on all} \\ 1 \times 1}}$  $\underbrace{y^TX\beta}_{\substack{\text{does not depend} \\ 1 \times 1}}$

Proposme since $y$ is given

$$\cancel{(\beta^TX^Ty)^T \quad y^TX\beta}$$
$$\cancel{\beta^TX^Ty = y^TX\beta} \quad \text{since transpose of a scalar is itself}$$

SAME GAME

$$\propto e^{-\frac{1}{2\sigma^2}\left(-\beta^TX^Ty - y^TX\beta + \beta^TX^TX\beta\right)} \quad e^{-\frac{1}{2\sigma^2} y^TX(X^TX)^{-1}X^Ty}$$

Since ols is a Gaussian

$$\propto \text{kernel of } N_p\left(\cdot, \sigma^2(X^TX)^{-1}\right) \qquad \text{What's the free variable now?}$$
$$\beta \to \beta \Rightarrow \text{mean is}$$

So

$$P(\beta \mid X, y, \sigma^2) = P(\beta \mid X) \quad \text{the OLS } \underset{\wedge}{\overset{dist}{\text{distr.}}} \text{ under non-informative prior}$$
$$(X^TX)^{-1}X^Ty$$

Same as $\quad \theta \mid X, \sigma^2 \sim N\left(\bar{x}, \frac{\sigma^2}{n}\right) \quad$ under $P(\theta) \propto 1$

$P(x \mid \theta, \sigma^2) \propto e^{-\frac{1}{2\sigma^2}(x-\theta)^2} \propto N(\theta, \sigma^2)$

$P(\theta \mid x, \sigma^2) \propto e^{-\frac{1}{2\sigma^2}(x-\theta)^2} \propto N(x, \sigma^2)$

same game $\qquad$ "L2 best" $\qquad$ Both ass's from above

Derivation of Ridge Estimator

Try to demonstrate that if

$\beta \perp X, \sigma^2$ a priori

$$P(\beta | X, \sigma^2) = P(\beta) = N_p\left(0, \frac{\sigma^2}{m} I_p\right)$$

no info about var

no dep of $\beta_i, \beta_j$

same variance for all $\beta_i$'s which is var of $\varepsilon$ modified

then we get conjugacy $\Rightarrow P(\beta_j) \sim N(0, \infty) \; \forall j$

if $m \to 0 \Rightarrow P(\beta) \propto 1$  just as before in unknown case

$$P(\beta) = \frac{1}{\sqrt{(2\pi)^p (\sigma_b^2 I_p)}} \; e^{-\frac{1}{2}(\beta - 0)^T \left(\frac{\sigma^2}{m} I_p\right)^{-1}(\beta - 0)}$$

Note $I^{-1} = I$

$$\propto e^{-\frac{1}{2}\beta^T \left(\frac{m}{\sigma^2} I\right)\beta} = e^{-\frac{1}{2\sigma^2}\beta^T (mI)\beta}$$

$$P(\beta | X, y, \sigma^2) \propto P(y | \beta, X, \sigma^2) \underbrace{P(\beta | X, \sigma^2)}_{= P(\beta)}$$

$$\propto e^{-\frac{1}{2\sigma^2}\left(-\beta^T X^T y - y^T X \beta + \beta^T X^T X \beta\right)} e^{-\frac{1}{2\sigma^2}\beta^T (mI)\beta}$$

Note: $a^T A u + a^T B u = a^T(Au + Bu) = a^T (A + B) u$

$$= e^{-\frac{1}{2\sigma^2}\left(-\beta^T X^T y - y^T X \beta + \beta^T (X^T X + mI)\beta\right)}$$

$$\propto e^{-\frac{1}{2\sigma^2}\left(y^T X (X^T X + mI)^{-1} X^T y - \beta^T X^T y - y^T X \beta + \beta^T (X^T X + mI)\beta\right)}$$

$$= e^{-\frac{1}{2\sigma^2}\left(y^T X - \beta^T (X^T X + mI)\right)\left((X^T X + mI)^{-1} X^T y - \beta\right)}$$

$$= e^{-\frac{1}{2\sigma^2}\left(y^T X - \beta^T (X^T X + mI)\right)(X^T X + mI)^{-1}(X^T X + mI)\left((X^T X + mI)^{-1} X^T y - \beta\right)}$$

$$= e^{-\frac{1}{2\sigma^2}\left(\underbrace{y^T X (X^T X + mI)^{-1}}_{\beta_R^T} - \beta^T\right)(X^T X + mI)\left(\underbrace{(X^T X + mI)^{-1} X^T y}_{or} - \beta\right)}$$

the "ridge estimator"

$$= e^{-\frac{1}{2\sigma^2}(\beta_R - \beta)^T (X^T X + mI)(\beta_R - \beta)}$$

$$\propto N\left(\beta_R, \sigma(X^T X + mI)^{-1}\right) \Rightarrow \hat\theta_{MMSE} = \hat\theta_{MAP} = \hat\theta_{MMAE} = \beta_R = (X^T X + mI)^{-1} X^T y$$

$\theta \in \mathbb{R}^p$ $\quad X_1, \ldots, X_n \sim N_p(\theta, \varepsilon)$, $\quad \theta \sim N_p(\mu_0, \varepsilon_0)$

$$P(\theta | X, \varepsilon^2) \propto P(X | \theta, \varepsilon) P(\theta | X)$$

$$= \prod_{i=1}^{n} N(\theta, \varepsilon) N_p(\mu_0, \varepsilon_0)$$

$$\propto N_p\left( \left(\varepsilon_0^{-1} + n\varepsilon^{-1}\right)^{-1} \left(\varepsilon_0^{-1}\mu_0 + n\varepsilon^{-1}\bar{X}\right), \left(\varepsilon_0^{-1} + n\varepsilon^{-1}\right)^{-1} \right)$$

really hard matrix algebra

if $n=1$

$$= N_p\left( \left(\varepsilon_0^{-1} + \varepsilon^{-1}\right)^{-1} \left(\varepsilon_0^{-1}\mu_0 + \varepsilon^{-1}X\right), \left(\varepsilon_0^{-1} + \varepsilon^{-1}\right)^{-1} \right)$$

if $\varepsilon = \sigma^2 I_p$

$$= N_p\left( \left(\varepsilon_0^{-1} + \tfrac{1}{\sigma^2}I\right)^{-1} \left(\varepsilon_0^{-1}\mu_0 + \tfrac{1}{\sigma^2}X\right), \left(\varepsilon_0^{-1} + \tfrac{1}{\sigma^2}I\right)^{-1} \right)$$

if $\mu_0 = 0$

$$= N_p\left( \left(\varepsilon_0^{-1} + \tfrac{1}{\sigma^2}I\right)^{-1} \left(\tfrac{1}{\sigma^2}X\right), \left(\varepsilon_0^{-1} + \tfrac{1}{\sigma^2}I\right)^{-1} \right)$$

if $\varepsilon_0 = \tfrac{\sigma^2}{m} I_p$

$$= N_p\left( \left(\tfrac{m}{\sigma^2}I + \tfrac{1}{\sigma^2}I\right)^{-1} \left(\tfrac{1}{\sigma^2}X\right), \left(\tfrac{m}{\sigma^2}I + \tfrac{1}{\sigma^2}I\right)^{-1} \right) = N\left( \tfrac{\sigma^2}{m+1}\tfrac{1}{\sigma^2}X, \tfrac{\sigma^2}{m+1}I \right)$$

if $M_0 = 0, \; \Sigma_0 = \frac{\sigma^2}{m} I_p$

$$= N_p \left( \left( \frac{m}{\sigma^2} I + \Sigma^{-1} \right)^{-1} \left( \Sigma^{-1} x \right), \; \left( \frac{m}{\sigma^2} I + \Sigma^{-1} \right)^{-1} \right)$$

if $\quad x = (X^T X)^{-1} X^T y, \quad \Sigma = \sigma^2 (X^T X)^{-1} \implies \Sigma^{-1} = \frac{1}{\sigma^2} (X^T X)$

$$= N_p \left( \left( \frac{m}{\sigma^2} I + \frac{1}{\sigma^2} (X^T X) \right)^{-1} \frac{1}{\sigma^2} (X^T X) \left( (X^T X)^{-1} X^T y \right), \; \left( \frac{m}{\sigma^2} I + \frac{1}{\sigma^2} (X^T X) \right)^{-1} \right)$$

$$= N_p \left( \sigma^2 \left( X^T X + m I \right)^{-1} \frac{1}{\sigma^2} (X^T X)(X^T X)^{-1} X^T y, \; \sigma^2 \left( X^T X + m I \right)^{-1} \right)$$

Ridge estimator

derived as

a special

case

of the

general

conjugacy

formula

What does Ridge do?

$m \to 0$  $p \sim N(0, \infty I) \Rightarrow \beta_j \sim N(0, \infty)$  is uniform

$\beta_R \to \beta = (X^T X)^{-1} X^T y$

$m \to \infty$  super informative

$\beta_R \to 0$  $(X^T X + mI) \to \begin{bmatrix} Big & & Insignificant \\ & Big & \ddots \\ Insignificant & \cdots & Big \end{bmatrix}^{-1} X^T y$

$$= \begin{bmatrix} 0 & 0 & Insignificant \\ & \ddots & \\ Insignificant & & 0 \end{bmatrix} X^T y \approx \vec{0}$$

So  $\beta_R \in [\vec{0}, \vec{\beta}]$  with $m \uparrow \Rightarrow$ more shrinkage to $0$.

Abnormal... if the moons are still right

$$(A+B)^{-1} = A^{-1} - \frac{1}{1+t} A^{-1} B A^{-1} \quad s.t. \ t := trace(BA^{-1})$$

$$(\underbrace{X^T X}_{A} + \underbrace{mI}_{B})^{-1} = (X^T X)^{-1} - \frac{1}{1+mT} (X^T X)^{-1} (mI)(X^T X)^{-1}$$

$$trace(mI(X^T X)^{-1}) = m \cdot \underbrace{trace(X^T X)^{-1}}_{T}$$

$$\beta_R = \left((X^T X)^{-1} - \frac{m}{1+mT}(X^T X)^{-1}(X^T X)^{-1}\right) X^T y = \beta - \frac{m}{1+mT}(X^T X)^{-1}\beta = \beta\left(I - \frac{m}{1+mT}(X^T X)^{-1}\right)$$

Before we had → sqd error loss

$\ell(\varepsilon)$

$b := \arg\min \varepsilon^T\varepsilon = \arg\min (y-Xb)^T(y-Xb) = \arg\min y^Ty - 2\beta^TX^Ty + \beta^TX^TX\beta$

What if we have ... $\|\varepsilon\|_2^2 + m\|\beta\|_2^2$ ← norm form $\Rightarrow$ min $\varepsilon$'s s.t $\beta^T\beta < C$

$\boxed{\ell(\varepsilon) + m\,\ell(\beta)}$ → "Ridge loss function"

$b = \arg\min \varepsilon^T\varepsilon + m\beta^T\beta$

$= \arg\min y^Ty - 2\beta^TX^Ty + \beta^TX^TX\beta + m\beta^T\beta$

$= y^Ty - 2\beta^TX^Ty + \beta^T(X^TX + mI)\beta$

$\nabla(\ ) = -2X^Ty + 2(X^TX+mI)\beta \overset{set}{=} 0$

$\Rightarrow (X^TX+mI)\beta = X^Ty \Rightarrow b_R := (X^TX+mI)^{-1}X^Ty$

Ridge derived as a minimum of the ridge loss function which is similar to the sqd.err. loss

Ridge loss is due to: $\boxed{\|x\|_p := \left(\sum_{i=1}^n |x_i|^p\right)^{1/p} \text{ NORM}}$

$\Leftarrow$ min $\|\varepsilon\|_2^2$ s.t. $\|\beta\|_2^2 < C$ "L2-penalized regression"

Setting up a Lagrange... $(\|\beta\|)$

min $f(\beta;X,y)$ s.t. $\underline{g(\beta) < C}$
$g(\beta) - C = 0$

$\nabla := \begin{bmatrix} \frac{\partial}{\partial b_0} \\ \vdots \\ \frac{\partial}{\partial b_p} \end{bmatrix}$

$\mathcal{L}(\beta,X,y,\lambda) = f(\beta;X,y) + \lambda(g(\beta)-c)$

$0 \overset{set}{=} \nabla\mathcal{L} = \nabla\|\varepsilon\|_2^2 + \lambda\nabla\|\beta\|_2^2 = \nabla(\|\varepsilon\|^2 + \|\beta\|^2)$

$\frac{\partial\mathcal{L}}{\partial\lambda} = 0$ to solve for $\lambda$ in terms of $C$.

MIDTERM ?

let $m := \lambda$ the Lagrange multiplier "

↓ FINAL