

Practical Session 1.

AMMI , 18/05/2020

Exercise 1.

$$(a) \quad f(x) = x A \quad \in \mathbb{R}^{m \times n}$$

$$\nabla f \in \mathbb{R}^{m \times n} \quad \nabla f_{ij} = f'_{ij}(x) \quad f'_{ij}(x) = f(x)_{ij} = A_{ij} x$$

$$\hookrightarrow \nabla f_{ij} = A_{ij} \quad \boxed{\nabla f = A}$$

$$(b) \quad f(x) = Ax \quad \in \mathbb{R}^m$$

m dimensions, n variables : $\nabla f \in \mathbb{R}^{m \times n}$.

$$\nabla f_{ij} = \frac{\partial f_i}{\partial x_j} \quad f_i(x) = A_i^T x = \sum_k A_{ik} x_k.$$

$$\frac{\partial f_i}{\partial x_j} = A_{ij} \quad \boxed{\nabla f = A}$$

$$(c) \quad f(x) = \text{Tr}(A^T X) \quad \in \mathbb{R} \quad X \in \mathbb{R}^{m \times n}$$

1 dimension, $m \times n$ dimension, variables : $\nabla f \in \mathbb{R}^{m \times n}$

$$f(x) = \sum A_{ij} x_{ij}$$

$$\nabla f_{ij} = \frac{\partial f}{\partial x_{ij}} = A_{ij} \quad \boxed{\nabla f = A}$$

Exercise 2.

$$X \in \mathbb{R}^{n \times p}, \lambda > 0 \quad M = X^T X + \lambda I_p$$

$\rightarrow \det M' := X^T X \in \mathbb{R}^{p \times p}$, we will show that M' has non-negative eigenvalues.
Let μ be an eigenvalue of M' , and y an eigenvector

$$M'y = \mu y \Rightarrow y^T M'y = \mu y^T y = \mu \|y\|^2$$

$$y^T M'y = y^T X^T X y = (Xy)^T X y = \|Xy\|^2$$

$$\text{Therefore } \mu = \frac{\|Xy\|^2}{\|y\|^2} \geq 0$$

\rightarrow The eigenvalues of M are equal to λ + the eigenvalues of M'
 $\mu \in \text{Eig}(M') \Leftrightarrow \mu + \lambda \in \text{Eig}(M) \quad (1)$

Therefore the eigenvalues of M are $\geq \lambda$ and > 0 ($\neq 0$)

so M is invertible

(M is invertible \Leftrightarrow its eigenvalues are $\neq 0$)

Proof of (1):

. M' is symmetric ($M'^T = M'$) so there exists $U \in \mathbb{R}^{p \times p}$ orthogonal ($U^T U = I$)
such that $M' = U^T D U$ $D = \text{Diag}(\lambda_1, \dots, \lambda_n)$

The λ_i 's are the eigenvalues of M'

$$\begin{aligned} M &= M' + \lambda I \\ &= U^T D U + \lambda U^T U \\ &= U^T (D + \lambda I) U \\ &= U^T \begin{pmatrix} \lambda_1 + \lambda & & \\ & \ddots & 0 \\ 0 & \cdots & \lambda_n + \lambda \end{pmatrix} U \end{aligned}$$

\hookrightarrow The eigenvalues of M are $\lambda_i + \lambda$

Exercise 3

(a), (b), (c) : seen in class

(d) \rightarrow pseudo-inverse

\rightarrow regularize with small λ : $(X^T X)^{-1}$ becomes $(X^T X + \lambda u I)^{-1}$
(Ridge Regression)

which is invertible (Ex2)

Exercise 4.

For $A \in \mathbb{R}^P$, $\text{Var}(A) \in \mathbb{R}^{P \times P}$ is the covariance matrix.

$$\text{Var}(A) = \mathbb{E}((A - \mathbb{E}(A))(A - \mathbb{E}(A))^T)$$

$$\begin{aligned}\text{Var}(A)_{ij} &= \text{Cov}(A_i, A_j) \\ &= \mathbb{E}(A_i A_j) - \mathbb{E}(A_i) \mathbb{E}(A_j)\end{aligned}$$

\rightarrow In this exercise, $\mathbb{E}\varepsilon = 0$ and $\text{Var}(\varepsilon) = \sigma^2 I$

so $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0$ if $i \neq j$

and $\text{Var}(\varepsilon_i) = \sigma^2$

\rightarrow Define $P = (X^T X)^{-1} X^T$.

$$\begin{aligned}\hat{\beta}^{\text{OLS}} &= PY \\ &= (X^T X)^{-1} X^T X \beta^* + P\varepsilon \quad \downarrow Y = X\beta^* + \varepsilon \quad (\text{assumption}) \\ &= \beta^* + P\varepsilon.\end{aligned}$$

$$\begin{aligned}\mathbb{E}(\hat{\beta}^{\text{OLS}}) &= \beta^* + \mathbb{E}(P\varepsilon) \quad \downarrow \mathbb{E}(P\varepsilon) = P\mathbb{E}(\varepsilon) = 0 \\ &= \beta^*\end{aligned}$$

\mathbb{E} is linear

$\rightarrow \hat{\beta}^{\text{OLS}}$ is unbiased

Exercise 4 (Continued)

$$\rightarrow \text{Var}(\hat{\beta}^{\text{OLS}}) = \text{Var}(\hat{\beta}^{\text{OLS}} - \beta^*) \\ = \text{Var}(P\epsilon) \\ = P \underbrace{\text{Var}(\epsilon)}_{= \sigma^2 I} P^T$$

$$\text{Var}(\hat{\beta}^{\text{OLS}}) = \sigma^2 P P^T \quad \begin{array}{l} \text{plugging back } P = (X^T X)^{-1} X^T \\ \Rightarrow P^T = X (X^T X)^{-1} \\ ((M^{-1})^T = (M^T)^{-1}) \end{array}$$

$$= \sigma^2 (X^T X)^{-1}$$

Exercise 5. (slide 17)

Gauss Markov Theorem: Least Squares estimator is Best Linear Unbiased Estimator

Assumptions: $Y = X\beta^* + \epsilon$

$$\left\{ \begin{array}{l} E\epsilon = 0 \\ \text{Var}(\epsilon) = \sigma^2 I \end{array} \right.$$

Among linear unbiased estimators, $\hat{\beta}^{\text{OLS}}$ has the lowest variance

In mathematical terms:

If $\tilde{\beta} = CY$ is an unbiased estimator,

$$\text{then } \text{Var}(\hat{\beta}^{\text{OLS}}) \leq \text{Var}(\tilde{\beta})$$

(Δ this condition does not mean $\text{Var}(\hat{\beta}^{\text{OLS}})_{ij} \leq \text{Var}(\tilde{\beta})_{ij}$ for i, j)

i.e. $M = \text{Var}(\tilde{\beta}) - \text{Var}(\hat{\beta}^{\text{OLS}})$ is a positive semidefinite matrix

i.e. $\forall x_0, x_0^T M x_0 \geq 0$

PROOF: Let $\tilde{\beta} = CY$ be an unbiased estimator, $C \in \mathbb{R}^{p \times n}$

Define $D := C - (X^T X)^{-1} X^T$

We are going to show that $\text{Var}(\hat{\beta}^{\text{OLS}}) \leq \text{Var}(\tilde{\beta})$

→ Step 1: show that $\mathbf{D}\mathbf{X} = \mathbf{0}$

$$\tilde{\beta} \text{ is unbiased : } \mathbb{E}(\tilde{\beta}) = \beta^* \quad (1)$$

$$\begin{aligned}\mathbb{E}(\tilde{\beta}) &= \mathbb{E}(CY) && \Downarrow \mathbb{E} \text{ is linear} \\ &= C\mathbb{E}(Y) && \Downarrow Y = X\beta^* + \varepsilon \\ &= CX\beta^* + C\mathbb{E}(\varepsilon) && \Downarrow \mathbb{E}(\varepsilon) = 0 \\ &= CX\beta^* \quad (2) && \Downarrow \end{aligned}$$

Combining (1) and (2) : $\beta^* = CX\beta^*$

$$\Rightarrow (CX - I)\beta^* = 0 \qquad \mathbf{D}\mathbf{X} = CX - I$$

$$\therefore \mathbf{D}\mathbf{X}\beta^* = 0 \qquad \text{This is true for all } \beta^*$$

$$\hookrightarrow \mathbf{D}\mathbf{X} = \mathbf{0} \quad \square$$

→ Step 2: show that $\text{Var}(\tilde{\beta}) = \text{Var}(\hat{\beta}^{\text{OLS}}) + \sigma^2 DD^T$

→ We reuse the notation $P = (X^T X)^{-1} X^T$
since $\mathbf{D}\mathbf{X} = \mathbf{0}$, $\underline{DP^T} = P D^T = \mathbf{0}$

$$\begin{aligned}\tilde{\beta} &= \beta^* CY \\ &= (D+P)(X\beta^* + \varepsilon) \\ &= \underbrace{DX\beta^*}_{=0} + D\varepsilon + \underbrace{PX\beta^*}_{=I} + P\varepsilon \\ &= \beta^* + (D+P)\varepsilon\end{aligned}$$

$$\begin{aligned}\text{Var}(\tilde{\beta}) &= \text{Var}((D+P)\varepsilon) \\ &= (D+P) \underbrace{\text{Var}(\varepsilon)}_{=\sigma^2 I} (D+P)^T \\ &= \sigma^2 \left(DD^T + \underbrace{PD^T}_{=0} + \underbrace{DP^T}_{=0} + \underbrace{PP^T}_{=(X^T X)^{-1}} \right) \\ &= \sigma^2 DD^T + \underbrace{\sigma^2 (X^T X)^{-1}}_{= \text{Var}(\hat{\beta}^{\text{OLS}})} \quad \square\end{aligned}$$

note: $Y \in \mathbb{R}^p$ random vector, $C \in \mathbb{R}^{p \times p}$ fixed.

$$\text{Then } \mathbb{E}(CY) = C\mathbb{E}(Y) \quad \text{and} \quad \text{Var}(CY) = C \text{Var}(Y) C^T$$

→ Step 3 (final) $M = \text{Var}(\tilde{\beta}) - \text{Var}(\hat{\beta}^{\text{OLS}})$ is positive semi-definite (pd)

$$\begin{aligned} M &= \text{Var}(\tilde{\beta}) - \text{Var}(\hat{\beta}^{\text{OLS}}) \\ &= \sigma^2 D D^T \end{aligned}$$

$$\begin{aligned} \text{Let } x_0 \in \mathbb{R}^n : \quad x_0^T M x_0 &= \sigma^2 x_0^T D D^T x_0 \\ &= \sigma^2 (D^T x_0)^T (D^T x_0) \\ &= \sigma^2 \|D^T x_0\|^2 \geq 0 \end{aligned}$$

This is true $\forall x_0$: M is p.d.

Finally: $\underline{\text{Var}(\tilde{\beta}) \geq \text{Var}(\hat{\beta}^{\text{OLS}})}$

$\tilde{\beta}$ has a "larger" variance than $\hat{\beta}^{\text{OLS}}$.

WHY DOES THIS MATTER?

Suppose we observe measurements $y = \mu + \underbrace{\varepsilon}_{\text{noise}}$ what we want to estimate

n measurements $y_i = \mu + \varepsilon_i$ $E(\varepsilon_i) = 0$
 $\text{Var}(\varepsilon) = \sigma^2 I$

We want an estimator $\hat{\mu}$ of μ .

Which is the best estimation?

① $\hat{\mu}_1 = \frac{1}{n} \sum_{i=1}^n y_i$

② $\hat{\mu}_2 = y_1$

③ $\hat{\mu}_3 = \frac{1}{n} \sum y_i + \frac{1}{n}$.

→ BIAS: $E(\hat{\mu}_1) = E(\hat{\mu}_2) = \Theta \mu$ $\hat{\mu}_1$ and $\hat{\mu}_2$ are unbiased.
• $E(\hat{\mu}_3) = \mu + \frac{1}{n} \neq \mu$ $\hat{\mu}_3$ is biased

→ VARIANCE: $\text{Var}(\hat{\mu}_1) = \text{Var}(\hat{\mu}_3) = \frac{\sigma^2}{n}$
• $\text{Var}(\hat{\mu}_2) = \sigma^2$

→ ERROR: $E(\hat{\mu}) = \text{bias}^2(\hat{\mu}) + \text{Var}(\hat{\mu})$

$$\left\{ \begin{array}{l} E(\hat{\mu}_1) = \frac{\sigma^2}{n} \\ E(\hat{\mu}_2) = \sigma^2 \\ E(\hat{\mu}_3) = \frac{1}{n^2} + \frac{\sigma^2}{n} \end{array} \right. \quad (\text{much})$$

$\rightarrow \hat{\mu}_1$ and $\hat{\mu}_2$ are both unbiased, but $\hat{\mu}_1$ is better than $\hat{\mu}_2$

$\rightarrow \hat{\mu}_3$ is biased, and $\hat{\mu}_2$ is not, but $\hat{\mu}_3$ is (much) better than $\hat{\mu}_2$

\hookrightarrow Variance matters just as much as bias.

Exercise 6

slide 33

RIDGE REGRESSION

$$\text{loss } J(\beta) = \text{MSE}(\beta) + \lambda \|\beta\|^2. \quad \lambda > 0$$

$$= \frac{1}{n} (Y - X\beta)^T (Y - X\beta) + \lambda \|\beta\|^2.$$

gradient

$$\nabla J(\beta) = \frac{2}{n} X^T (X\beta - Y) + 2\lambda \beta.$$

$$\nabla J(\hat{\beta}_\lambda^{\text{ridge}}) = 0 \quad \Leftrightarrow \quad \hat{\beta}_\lambda^{\text{ridge}} = (X^T X + \lambda n I)^{-1} X^T Y$$

NB: when $\lambda > 0$, $(X^T X + \lambda n I)^{-1}$ always exists (exercise 2)

Assumption: $X^T X = n I$ (orthogonal vectors x_i)

$$\hookrightarrow \hat{\beta}_\lambda^{\text{ridge}} = \frac{1}{n(1+\lambda)} X^T Y \quad \downarrow Y = X\beta^* + \varepsilon$$

$$= \frac{1}{n(1+\lambda)} \underbrace{(X^T X \beta^* + X^T \varepsilon)}_{= n I}$$

$$\hat{\beta}_\lambda^{\text{ridge}} = \frac{1}{n(1+\lambda)} \beta^* + \frac{1}{n(1+\lambda)} X^T \varepsilon$$

Taking expectations:

$$\begin{aligned} E(\hat{\beta}_\lambda^{\text{ridge}}) &= \frac{1}{1+\lambda} \beta^* + \frac{1}{n(1+\lambda)} X^T E(\varepsilon) \\ &= \frac{1}{1+\lambda} \beta^* \end{aligned}$$

$$\rightarrow \text{Hence, } \text{bias}(\hat{\beta}_\lambda^{\text{ridge}}) = \mathbb{E}(\hat{\beta}_\lambda^{\text{ridge}}) - \beta^*$$

$$= -\frac{\lambda}{1+\lambda} \beta^* \quad \begin{array}{l} \rightarrow \text{large } \lambda = \text{large bias} \\ \rightarrow \text{small } \lambda = \text{small bias.} \end{array}$$

$$\rightarrow \text{Var}(\hat{\beta}_\lambda^{\text{ridge}}) = \text{Var}\left(\frac{1}{n(1+\lambda)} X^T \varepsilon\right)$$

$$= \frac{1}{n^2(1+\lambda)^2} \underbrace{X^T \text{Var}(\varepsilon) X}_{\sigma^2 I}$$

$$= \frac{\sigma^2}{n^2(1+\lambda)^2} \underbrace{X^T X}_{nI}$$

$$= \frac{\sigma^2}{n(1+\lambda)^2} I \quad \begin{array}{l} \rightarrow \text{large } \lambda = \text{small variance} \\ \rightarrow \text{small } \lambda = \text{large variance.} \end{array}$$

under our assumptions, $\text{Var}(\hat{\beta}^{\text{OLS}}) = \sigma^2 (X^T X)^{-1} = \frac{\sigma^2}{n} I$

$$\therefore \text{Var}(\hat{\beta}_\lambda^{\text{ridge}}) = \frac{1}{(1+\lambda)^2} \text{Var}(\hat{\beta}^{\text{OLS}})$$

$\lambda \rightarrow 0$: $\hat{\beta}_\lambda^{\text{ridge}} \rightarrow \hat{\beta}^{\text{OLS}}$: no bias, large variance
 $\lambda \rightarrow +\infty$: $\hat{\beta}_\lambda^{\text{ridge}} \rightarrow 0$: large bias, no variance.

Exercise 7 : Ridge performance

$$f(\lambda) = \underset{S, X_0}{\mathbb{E}} \left[\text{bias}^2(X_0^\top \hat{\beta}_\lambda^{\text{ridge}}) + \text{Var}(X_0^\top \hat{\beta}_\lambda^{\text{ridge}}) \right] \quad \begin{cases} \mathbb{E} X_0 = 0 \\ \mathbb{E} X_0^\top X_0 = I \end{cases}$$

= expectation of total error.

→ Let us compute the two terms:

$$\begin{aligned} - \text{bias}^2(X_0^\top \hat{\beta}_\lambda) &= (X_0^\top \text{bias}(\hat{\beta}_\lambda))^2 \\ &= \frac{\lambda^2}{(1+\lambda)^2} (X_0^\top \beta^*)^2 \quad \text{bias}(\hat{\beta}) = \beta^* \\ &= \frac{\lambda^2}{(1+\lambda)^2} \beta^{*\top} X_0 X_0^\top \beta^* \quad (X_0^\top \beta^*)^2 = \sum_{j=1}^p X_{0,j} X_{0,j} \beta_j^* \beta_j^* \\ &= \beta^{*\top} X_0 X_0^\top \beta^* \\ - \text{Var}(X_0^\top \hat{\beta}_\lambda^{\text{ridge}}) &= X_0^\top \text{Var}(\hat{\beta}_\lambda^{\text{ridge}}) X_0 \\ &= \frac{\sigma^2}{n(1+\lambda)} X_0^\top X_0 \end{aligned}$$

→ Taking expectations, with $\mathbb{E}(X_0 X_0^\top) = I$, we get:

$$f(\lambda) = \frac{\lambda^2}{(1+\lambda)^2} \|\beta^*\|^2 + \frac{\sigma^2 p}{n(1+\lambda)^2}. \quad \mathbb{E}(X_0^\top X_0) = \sum_{i=1}^p \mathbb{E}(X_{0,i} X_{0,i}^\top) = p$$

→ Finding the optimum λ^* :

$$\text{Define } a = \|\beta^*\|^2 = f(+\infty) \quad b = \frac{\sigma^2 p}{n} = f(0)$$

$$f(\lambda) = \frac{1}{(1+\lambda)^2} (a\lambda^2 + b)$$

$$f'(\lambda) = \frac{2}{(1+\lambda)^3} (a\lambda + b)$$

$$f'(\lambda^*) = 0 \quad (\Rightarrow) \quad \lambda^* = \frac{b}{a} = \frac{\sigma^2 p}{n \|\beta^*\|^2}$$

keeping our notation: $f(\lambda^*) = \frac{1}{(1+\frac{b}{a})^2} \left(a \left(\frac{b}{a} \right)^2 + b \right) = \frac{ab}{a+b}$

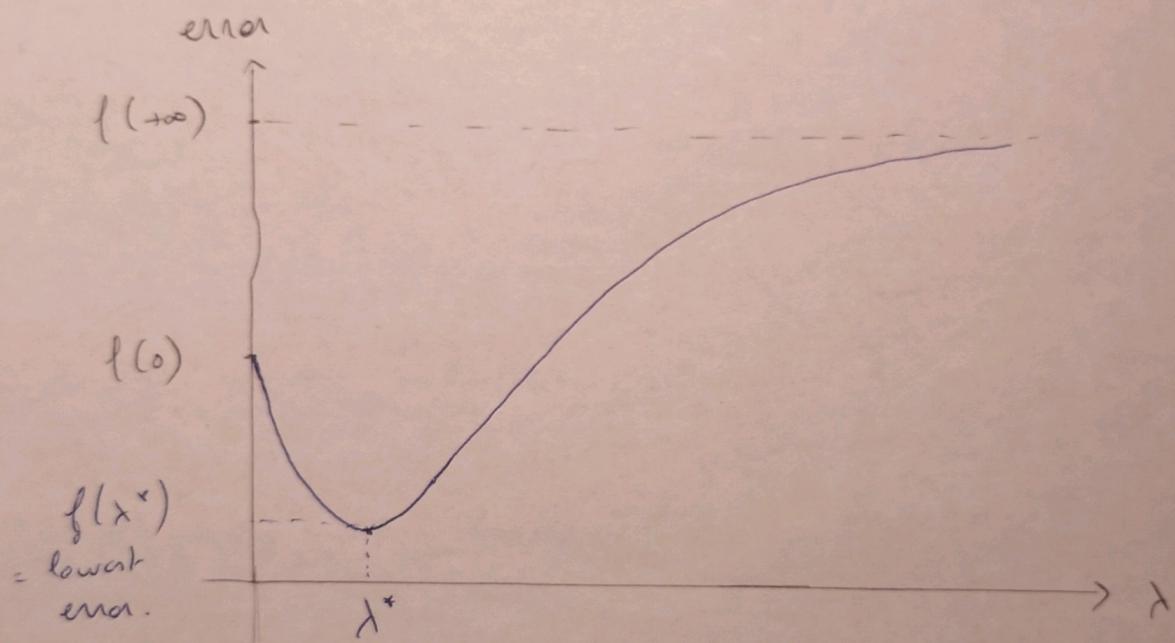
$$f(\lambda^*) = \frac{ab}{a+b} = \frac{f(0) f(+\infty)}{f(0) + f(+\infty)}$$

$$a \geq 0 \Rightarrow \frac{ab}{a+b} \leq \frac{a(a+b)}{a+b} = a$$

$$b \geq 0 \Rightarrow \frac{ab}{a+b} \leq \frac{(a+b)b}{a+b} = b$$

$$\left. \begin{array}{l} \frac{ab}{a+b} \leq \min(a, b) \end{array} \right\}$$

So finally, $f(\lambda^*) \leq \min(f(0), f(+\infty))$



(can be)

CONCLUSION : Some regularization \checkmark is better than none
and better than the estimator O

Exercise 8.

In notebook of practical session 2 in more detail

Exercise 8

$$\beta^{\text{new}} = \beta^{\text{old}} - \left[\nabla_{\beta}^2 J(\beta^{\text{old}}) \right]^{-1} \nabla_{\beta} J(\beta^{\text{old}})$$

lemma:

β^{new} is also the solution to $\min_{\beta} J_q(\beta)$

where $J_q(\beta) := J(\beta^{\text{old}}) + \nabla J(\beta^{\text{old}})^T (\beta - \beta^{\text{old}})$
 $+ \frac{1}{2} (\beta - \beta^{\text{old}})^T \nabla_{\beta}^2 J(\beta^{\text{old}}) (\beta - \beta^{\text{old}})$

J_q is the quadratic approximation to J in β^{old} .

(analogy in 1 dimension: $J_q(x) = J(x_0) + J'(x_0)(x-x_0) + \frac{1}{2} J''(x_0)(x-x_0)^2$)

proof: $\nabla J_q(\beta) = \nabla J(\beta^{\text{old}}) + \nabla_{\beta}^2 J(\beta^{\text{old}}) (\beta - \beta^{\text{old}})$
 $= 0 \Leftrightarrow \beta = \beta^{\text{old}} - \left[\nabla_{\beta}^2 J(\beta^{\text{old}}) \right]^{-1} \nabla_{\beta} J(\beta^{\text{old}})$
 $= \beta^{\text{new}}$

Writing all the terms of J_q that depend on β :

$$\rightarrow \beta^T \nabla J(\beta^{\text{old}}) = \frac{1}{n} \beta^T X^T P y + 2\lambda \beta^T \beta^{\text{old}}$$

$$\rightarrow \frac{1}{2} \beta^T \nabla_{\beta}^2 J(\beta^{\text{old}}) \beta = \frac{1}{2n} \beta^T X^T W X \beta + \lambda \beta^T \beta$$

$$\rightarrow -\beta^T \nabla_{\beta}^2 J(\beta^{\text{old}}) \beta^{\text{old}} = -\frac{1}{n} \beta^T X^T W X \beta^{\text{old}} - 2\lambda \beta^T \beta^{\text{old}}.$$

Putting it all together:

$$J_q(\beta) = -\frac{1}{n} \beta^T X^T W X \beta^{\text{old}} + \frac{1}{n} \beta^T X^T P y + \frac{1}{2n} \beta^T X^T W X \beta + \lambda \beta^T \beta + C$$

$$= -\frac{1}{n} \beta^T X^T W \underbrace{(\beta^{\text{old}} - W^{-1} P y)}_{:= g} + \frac{1}{2n} \beta^T X^T W X \beta + \lambda \beta^T \beta + C$$

$$+ \frac{1}{2n} (X \beta - g)^T W (X \beta - g) + \lambda \|\beta\|^2 + C$$