

## 1<sup>o</sup> Linear Basis Model.

Basic Linear model:

$$y(x, w) = w_0 + w_1 x_1 + \dots + w_D x_D = w^T \cdot x$$

Basis function  $\phi(\cdot)$ :

Basis function is some sort of "fixed preprocessing" of the input data

By using nonlinear basis functions, the linear model can contain some linear information on the input space  $x$ .

II

Linear basis model:

$$y(x, w) = w_0 + \sum_{i=1}^{M+1} w_i \phi_i(x) = w^T \Phi(x)$$

$$x = \begin{bmatrix} x_1 & x_2 & x_3 & \dots & x_N \\ \vdots & & & & \\ x_m & x_n & x_m & \dots & x_m \end{bmatrix} \in \mathbb{R}^{M \times N}$$

$$\text{let } x_i := \begin{bmatrix} x_{1i} \\ x_{2i} \\ \vdots \\ x_{Mi} \end{bmatrix}$$

$\Rightarrow \Phi(\cdot): X \rightarrow \mathbb{R}^J$

$$\Rightarrow \Phi(x) = \begin{bmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_{J-1}(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \dots & \phi_{J-1}(x_2) \\ \vdots & & & \\ \phi_0(x_N) & \phi_1(x_N) & \dots & \phi_{J-1}(x_N) \end{bmatrix}$$

Examples of Basis functions:

1<sup>o</sup> Polynomial basis:  $\phi_j(x) = x^j$

2<sup>o</sup> Fourier basis: periodic data, or data with known boundaries.

$$\phi_0(x) = 1 \quad \phi_j(x) = \cos(\omega_j x + \phi_j)$$

3<sup>o</sup> Radial basis function.

A function that is symmetric around a center and typically decay to 0 when near the center

Gaussian Radial basis:

$$\phi_j(x) = \exp\left\{-\frac{1}{\rho} \|x - c_j\|_2^2\right\} = \exp\left\{-\frac{1}{2\rho} \sum_{d=1}^D (x_d - c_{j,d})^2\right\}$$

4<sup>o</sup> Sigmoid basis: (Parametric Soft Thresholds)

$$\phi_j(x) = \frac{1}{1 + \exp\{-x^T w_j\}} \Rightarrow \text{close to zero when } x^T w_j \text{ is large negative.}$$

2<sup>o</sup> How to solve parameters:

2) Maximum Likelihood.

t: the true target variable

$y(x, w)$ : our model.

Assumption 1.

$$\Rightarrow t = y(x, w) + \varepsilon. \quad \varepsilon \sim N(0, \sigma^2) \sim N(0, \beta^{-1}), \beta \text{ is the precision.}$$

$$\Rightarrow P\{t | x, w, \beta\} = N\{t | y(x, w), \beta^{-1}\}$$

Given data  $X = \{x_1, \dots, x_N\}$

Assumption 2: different samples are i.i.d.

$$\Rightarrow P\{t | X, w, \beta\} = \prod_{n=1}^N N\{t_n | w^T \phi(x_n), \beta^{-1}\}$$

$$\begin{aligned}\Rightarrow \ln p_{\theta}(x, w, \beta) &= \sum_{n=1}^N \ln N\left\{t_n \mid w^T \phi(x_n), \beta\right\} \\ &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln (2 \pi) - \beta \cdot \frac{1}{2} \sum_{n=1}^N\left\{t_n - w^T \phi(x_n)\right\}^2 \\ &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln (2 \pi) - \underbrace{\beta \cdot E_D(w)}_{\text{square loss function.}}\end{aligned}$$

$\Rightarrow$  Maximizing MLE  $\Leftrightarrow$  Minimizing square loss function.

Thus:  $\nabla \ln p_{\theta}(w, \beta) = \sum_{n=1}^N\left\{t_n - w^T \phi(x_n)\right\} \phi(x_n)^T = 0$

$$\Rightarrow w_{ML} = (\Phi^T \Phi)^{-1} \Phi^T t.$$

b) Least Squares:

objection: minimize square loss.

$$\Rightarrow E_D(w) = \frac{1}{2} \sum_{n=1}^N\left\{t_n - w^T \phi(x_n)\right\}^2$$

$$=\frac{1}{2}(t - w^T \Phi)^T(t - w^T \Phi)$$

$$\Rightarrow \nabla_w E_D(w) = -\frac{1}{2} \Phi^T t + \frac{1}{2} \Phi^T \Phi w = 0$$

$$\Rightarrow \hat{w} = (\Phi^T \Phi)^{-1} \Phi^T t.$$

$$\text{Var}(\hat{w}) = E\left\{\left[\hat{w} - E(w)\right]\left[\hat{w} - E(w)\right]^T\right\}$$

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

2<sup>o</sup> Regularized LS,

General way: modify the objective function by adding a penalty term

$$E_D(w) + \lambda E_W(w)$$

Collinearity: the event of two (or multiple) covariates being strongly related

- ⇒ The space spanned by super-collinear covariates is a lower-dimensional subspace of the parameter space.
- ⇒ The design matrix  $X$  is rank deficient.

Why collinearity is a problem

First consider the case that two (or multiple) covariates being perfectly dependent (super-collinearity)

⇒  $X$  is column rank deficient, column rank  $< p$ .

∴ There exists a non-trivial  $U \in \mathbb{R}^p$  s.t.  $XU = 0_p$

⇒  $X^T X U = X^T 0_p = 0_p$  ⇒  $\det(X^T X) = 0 \Rightarrow$  one or more of the eigenvalues of  $A$  is zero.

∴  $X^T X$  is not invertible.

By spectral decomposition,

$X^T X = \sum_{j=1}^p \lambda_j U_j U_j^T$ ,  $U_j$  is the eigenvector and  $\lambda_j$  is the corresponding eigenvalue

⇒  $X^T X$  is singular

⇒  $(X^T X)^{-1} = \sum_{j=1}^p \lambda_j^{-1} U_j U_j^T$

⇒ We can't have  $U = (X^T X)^{-1} X^T b$ .

⇒ one or more of  $\lambda_j$  must be 0.

Next we consider the case that there isn't an exact linear relationship among the predictors, but the correlations are close to one.

⇒  $X^T X$  is invertible

However:

let  $\Sigma \equiv X^T X$

$$\therefore \Sigma^{-1} = \frac{1}{\det(\Sigma)} \text{adj}(\Sigma)$$

$\Sigma$  is close to singular  $\Rightarrow \det(\Sigma) \rightarrow 0$

$\therefore \text{Var}(w) = \sigma^2(X^T X)$  is large.

when  $P = n$ ,

the system of equations does not have a unique solution

Conclusion: when  $X$  is rank deficient, there are three consequent:

1°  $w$  cannot be estimated

2° the variance of the estimation is large

3° the solution of  $w$  is not unique.

To solve the above problem, we introduce minimum least squares estimator.

Defn. The minimum least squares estimator of the regression parameter minimizes the square loss and is of minimum length.

$$\Rightarrow \hat{w}_{MLS} = \underset{w \in \mathbb{R}^P}{\text{arg min}} \|t - \Phi w\|_2^2$$

$$\text{s.t. } \|\hat{w}_{MLS}\|_2^2 < \|w\|_2^2 \quad \text{for all } w \text{ that minimizes}$$

the square loss.

$\Rightarrow \boxed{\text{Ridge Regression}}$

Objective function:

$$\min \sum_{n=1}^N \{t_n - w^\top \phi(x_n)\}^2 + \frac{\lambda}{2} w^\top w.$$

Ridge estimator:  $\hat{w} = (\Phi^\top \Phi + \lambda I_p)^{-1} \Phi^\top t$ .

How can Ridge estimator solve the above problems?

By singular value decomposition:

$$\Phi = U_\phi D_\phi V_\phi^\top$$

For OLS:

$$\begin{aligned}\hat{w} &= (\Phi^\top \Phi)^{-1} \Phi^\top t \\ &= (U_\phi D_\phi^\top D_\phi V_\phi^\top)^{-1} U_\phi D_\phi^\top U_\phi^\top + \\ &= U_\phi (D_\phi^\top D_\phi)^{-1} U_\phi^\top U_\phi D_\phi^\top U_\phi^\top + \\ &= U_\phi (D_\phi^\top D_\phi)^{-1} D_\phi^\top U_\phi^\top +\end{aligned}$$

For Ridge:

$$\begin{aligned}\hat{w} &= (\Phi^\top \Phi + \lambda I_p)^{-1} \Phi^\top t \\ &= (U_\phi D_\phi^\top U_\phi^\top U_\phi D_\phi U_\phi^\top + \lambda I_p)^{-1} U_\phi D_\phi^\top U_\phi^\top + \\ &= (U_\phi D_\phi^\top D_\phi U_\phi^\top + \lambda U_\phi U_\phi^\top)^{-1} U_\phi D_\phi^\top U_\phi^\top + \\ &= U_\phi (D_\phi^\top D_\phi + \lambda I_p)^{-1} U_\phi^\top U_\phi D_\phi^\top U_\phi^\top + \\ &= U_\phi (D_\phi^\top D_\phi + \lambda I_p)^{-1} D_\phi^\top U_\phi^\top +\end{aligned}$$

For OLS: we have  $(D_\phi^T D_\phi)^{-1} D_\phi^T$  in the estimation

$$(D_\phi^T D_\phi)^{-1} D_\phi^T \in \mathbb{R}^{P \times n}$$

$$= \begin{bmatrix} \frac{1}{s_1} & & & 0 \\ & \frac{1}{s_2} & & 0 \\ & & \ddots & \\ & & & \frac{1}{s_p} & 0 \end{bmatrix} \Rightarrow \text{the reciprocal of the non-zero singular values on the diagonal of the left prep matrix}$$

$$(D_\phi^T D_\phi + \lambda I_{pp})^{-1} D_\phi^T \in \mathbb{R}^{P \times n}$$

$$= \begin{bmatrix} \frac{s_1}{s_1^2 + \lambda} & & & 0 \\ & \frac{s_2}{s_2^2 + \lambda} & & 0 \\ & & \ddots & \\ & & & \frac{s_p}{s_p^2 + \lambda} & 0 \end{bmatrix}$$

let  $d_{jj}$  denotes the non-zero singular values on the left diagonal:

$$d_{jj}^{-1} = d_{jj} (d_{jj}^2 + \lambda)$$

$\Rightarrow$  Ridge penalty shrinks the singular values

Result that in the case of super-collinearity,  $(\Phi^T \Phi)$  is not invertable

because  $(\Phi^T \Phi)^{-1} = \sum_{j=1}^P (d_{jj}^2)^{-1} \cdot v_j v_j^T$ , where one or more of  $d_{jj}$  is 0

By Ridge penalty:

$$(\Phi^T \Phi + \lambda I_{pp})^{-1} = \sum_{j=1}^P (d_{jj}^2 + \lambda)^{-1} v_j v_j^T$$

$$\text{If } d_{jj} = 0 \Rightarrow d_{jj}^2 + \lambda > 0$$

$\Rightarrow \Phi^T \Phi + \lambda I_{pp}$  is invertable.

We could transform the Ridge regression into an optimization problem:

$$w^* = \underset{\|w\|_2^2 \leq \lambda}{\operatorname{arg\min}} \|t - \Phi w\|^2$$

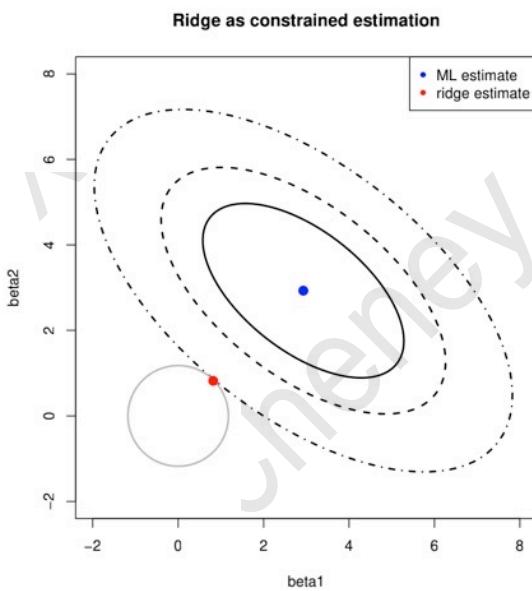
This could be solved by KKT condition.

$$\text{Lagrangean: } L = \|t - \mathbf{x}w\|_2^2 + \nu(\|w\|_2^2 - \lambda)$$

$$\text{1st KKT: } \nabla_w L = 0$$

$$\text{2nd KKT: } \nu(\|w(\nu)\|_2^2 - \lambda) = 0$$

$\Rightarrow$  Both have solution when  $\nu = \|w(\nu)\|_2^2$



} Lasso

objective function:

$$\|t - \mathbf{x}w\|_2^2 + \lambda\|w\|_1$$

Proposition. Lasso does not have a closed form solution, except when  $\Phi$  is an orthogonal matrix.

Proof. Let  $w_{OLS} \equiv (\Phi^T \Phi)^{-1} \Phi^T t = \Phi^T t$

$$\min_w \|t - \Phi w\|_2^2 + \lambda \|w\|_1 = \min_w t^T t - t^T \Phi w - w^T \Phi^T t + w^T \Phi^T \Phi w$$

$$+ \lambda \sum_{j=1}^P |w_j|$$

$$\propto \min_w -w_{OLS}^T w - w^T w_{OLS} + w^T w + \lambda |w|$$

$$= \min_{w_1, \dots, w_p} \sum_{j=1}^P (-2 w_{j, OLS} \cdot w_j + w_j^2 + \lambda |w_j|)$$

$$= \sum_{j=1}^P \min_{w_1, \dots, w_p} (-2 w_{j, OLS} \cdot w_j + w_j^2 + \lambda |w_j|)$$

$\Rightarrow$  The function becomes fixed  $w_j$  for each:

$$\min_{w_j} -2 w_{j, OLS} w_j + w_j^2 + \lambda |w_j|$$

$$= \min \begin{cases} -2 w_{j, OLS} w_j + w_j^2 + \lambda w_j & \text{if } w_j > 0 \\ -2 w_{j, OLS} w_j + w_j^2 - \lambda w_j & \text{if } w_j < 0 \end{cases}$$

$$\Rightarrow w_j^*(\lambda) = \begin{cases} w_{j, OLS} - \frac{1}{2}\lambda & \text{if } w_{j, OLS} > \frac{1}{2}\lambda \\ w_{j, OLS} + \frac{1}{2}\lambda & \text{if } w_{j, OLS} < -\frac{1}{2}\lambda \\ 0 & \text{otherwise} \end{cases}$$

$$\Rightarrow w_j^*(\lambda) = \text{sgn}(w_{j, OLS}) (|w_{j, OLS}| - \frac{1}{2}\lambda)_+$$

How to estimate  $w$  in the general case

Proposition 1: Lasso estimator may not be unique.

Proof: Considering the case of super-collinearity.

$$t_i = \Phi_{i+} w + \varepsilon_i$$

If  $\Phi$  only has two columns, with

$$\Phi_{+,1} = \Phi_{+,2}$$

$$u := w_1 + w_2, v := w_1 - w_2$$

$$\Rightarrow \|t - \Phi w\|_2^2 + \lambda \|w\|_1 = \|t - \Phi_{+,1} \cdot u\|_2^2 + \frac{1}{2} \lambda (|u+v| + |u-v|)$$

$|u+v| + |u-v|$  is minimized with respect to  $v$  for any  $v$  such that  $|v| < |u|$

$\Rightarrow$  if  $|v| < |u|$ :

$$|u+v| + |u-v| = |u|$$

$$\Rightarrow \perp = \|t - \Phi_{+,1} \cdot u\|_2^2 + \lambda |u|$$

If  $\lambda \rightarrow 0$ ,  $\perp \rightarrow \|t - \Phi_{+,1} \cdot u\|_2^2$

$\Rightarrow$  For any  $v = w_1 - w_2$ ,  $\perp$  will be the same.

$\Rightarrow w^*(\lambda)$  is not uniquely defined.

$$\text{Since } w_1^* = \frac{1}{2} (\hat{u}(w) + \hat{v}(w))$$

$$w_2^* = \frac{1}{2} (\hat{u}(w) - \hat{v}(w))$$

Proposition 2: Though  $w^*(w)$  is not uniquely defined,  $\Phi \cdot w^*$  is uniquely defined.

Proof: Suppose, for contradiction,

$$\exists w_\alpha^*, w_\beta^* \text{ s.t.}$$

$$\Phi w_\alpha^* \neq \Phi w_\beta^*$$

$$\|t - \Phi w_\alpha^*\|_2^2 + \lambda \|w_\alpha^*\|_1 = \perp = \|t - \Phi w_\beta^*\|_2^2 + \lambda \|w_\beta^*\|_1$$

$\Rightarrow \exists \theta \in (0, 1)$ :

$$\begin{aligned} & \|t - \underline{\Phi}[(1-\theta)w_\alpha^* + \theta w_\beta^*]\|_2^2 + \lambda \| (1-\theta)w_\alpha^* + \theta w_\beta^* \|_1 \\ &= \| (1-\theta)[t - \underline{\Phi}w_\alpha^*] + \theta[t - \underline{\Phi}w_\beta^*] \|_2^2 + \lambda \| (1-\theta)w_\alpha^* + \theta w_\beta^* \|_1 \\ &< (1-\theta) \|t - \underline{\Phi}w_\alpha^*\|_2^2 + \theta \|t - \underline{\Phi}w_\beta^*\|_2^2 + \lambda (1-\theta) \|w_\alpha^*\|_1 + \theta \|w_\beta^*\|_1 \\ &= (1-\theta)L + \theta L = L \end{aligned}$$

which means,  $\exists \theta \in (0, 1)$ , s.t. the cost function of  $w = (1-\theta)w_\alpha^* + \theta w_\beta^*$  is lower than  $L$ .

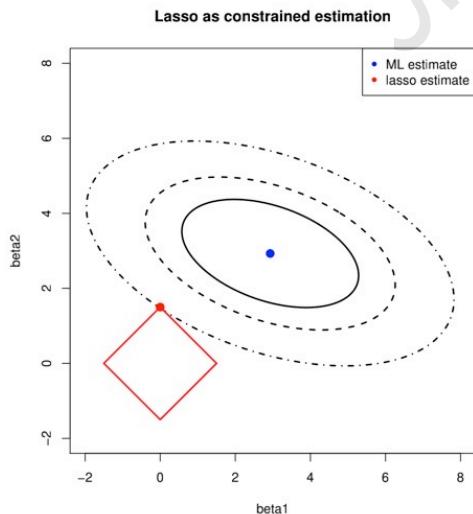
$\Rightarrow$  contradiction!

$\Rightarrow \underline{\Phi}w^* = \underline{\Phi}w_\beta^*$

$\Rightarrow \underline{\Phi}w^*$  is unique.

Lasso as an optimization problem:

$$w^* = \underset{\|w\|_1 \leq \lambda}{\text{argmin}} \|t - \underline{\Phi}w\|$$



$\Rightarrow$  the lasso estimate is sparse.

How to estimate?

## 1<sup>o</sup> Quadratic Planning.

$$\min_{\mathbf{w} \geq 0} \frac{1}{2} (\mathbf{t} - \mathbf{\Phi w})^T (\mathbf{t} - \mathbf{\Phi w})$$

$\mathbf{Rw} \geq 0$

$\mathbf{R} \in \mathbb{R}^{P \times P}$ : the matrix that specifies the linear constraints on the parameter  $w$

constraints:  $\{ \mathbf{w} \in \mathbb{R}^P : \| \mathbf{w} \|_1 < \lambda_1 \}$

i.e. for  $P=2$ :

$$\{ \mathbf{w} \in \mathbb{R}^2 : w_1 + w_2 < \lambda_1 \} \cap \{ \mathbf{w} \in \mathbb{R}^2 : w_1 + w_2 > -\lambda_1 \} \cap \{ \mathbf{w} \in \mathbb{R}^2 : w_1 - w_2 > -\lambda_1 \} \cap \{ \mathbf{w} \in \mathbb{R}^2 : w_1 - w_2 < \lambda_1 \}$$

$$\Rightarrow \mathbf{R} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \\ -1 & 1 \end{bmatrix} \Rightarrow \mathbf{Rw} \geq -\lambda_1$$

By Lagrangian:

$$L(\mathbf{w}, \nu) = \frac{1}{2} (\mathbf{t} - \mathbf{\Phi w})^T (\mathbf{t} - \mathbf{\Phi w}) + \nu^T \mathbf{R} \mathbf{w}$$

$$\Rightarrow \mathbf{w}^* = (\mathbf{\Phi}^T \mathbf{\Phi})^{-1} \mathbf{\Phi}^T \mathbf{t} + (\mathbf{\Phi}^T \mathbf{\Phi})^{-1} \mathbf{R}^T \nu$$

Substituting  $\mathbf{w}^*$  to  $\mathbf{w}$

$$\min_{\nu} \frac{1}{2} \nu^T \mathbf{R} (\mathbf{\Phi}^T \mathbf{\Phi})^{-1} \mathbf{\Phi}^T \mathbf{t} + \nu^T \mathbf{R} (\mathbf{\Phi}^T \mathbf{\Phi})^{-1} \mathbf{\Phi}^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{\Phi} (\mathbf{\Phi}^T \mathbf{\Phi})^{-1} \mathbf{\Phi}^T \mathbf{t}$$

s.t.  $\nu \geq 0$

We could find the optimum  $\nu^*$  of this quadratic planning problem and thus find the optimum  $w$ .

## 2<sup>o</sup> Gradient Descent.

For the function  $f: \mathbb{R}^P \rightarrow \mathbb{R}$  at  $x \in \mathbb{R}^P$  in the direction of  $v$ ,

$$f'(x) = \lim_{\tau \rightarrow 0} \frac{[f(x + \tau v) - f(x)]}{\tau}$$

for  $x \in \mathbb{R}$ ,  $f(x) = |x|$ ,  $v \in \mathbb{R} \setminus \{0\}$ , define:

$$f'(x) = \begin{cases} v \frac{x}{|x|} & \text{if } x \neq 0 \\ v & \text{if } x=0 \end{cases}$$

limit  $v \in \mathbb{R}^+$  to be:  $\|v\|=1$  and  $v$  is the direction of the steepest descent.

$$\nabla f(x) = \begin{cases} f'(x) \cdot v_{\text{opt}} & \text{if } f'(x) \geq 0 \\ 0_p & \text{if } f'(x) < 0 \end{cases}$$

$$v_{\text{opt}} \in \mathbb{R}^p = \operatorname{argmax}_{\{v: \|v\|=1\}} f'(x)$$

$$\text{Define: } L_{\text{OLS}} = (\hat{t} - \Phi w)^T (\hat{t} - \Phi w)$$

$$\Rightarrow \frac{\partial}{\partial w_j} L_{\text{OLS}}(\hat{t}, \Phi, w) = \begin{cases} \frac{\partial}{\partial w_j} L_{\text{OLS}}(\hat{t}, \Phi, w) - \lambda \operatorname{sign}(w_j) & \text{if } w_j \neq 0 \\ \frac{\partial}{\partial w_j} L_{\text{OLS}}(\hat{t}, \Phi, w) - \lambda_1 \operatorname{sign}\left[\frac{\partial}{\partial w_j} L_{\text{OLS}}(\hat{t}, \Phi, w)\right] & \text{if } w_j = 0 \text{ and } \frac{\partial L_{\text{OLS}}}{\partial w_j} > \lambda_1 \\ 0 & \text{otherwise.} \end{cases}$$

How to select the Hyperparameters  $\alpha$ ?

10 The degree of freedom (dof)

Degrees of freedom denotes how many independent things are there.

The dof of regression:

$$\text{OLS: } \hat{t} = \Phi(\Phi^T \Phi)^{-1} \Phi^T \tau$$

$$\text{Hat matrix: } H = \Phi(\Phi^T \Phi)^{-1} \Phi^T$$

The dof of regression is:  $\text{tr}(H)$

For Ridge Regression,

$$H(\lambda) = \Phi (\Phi^T \Phi + \lambda I_P)^{-1} \Phi^T$$

$$\text{let } D_\phi := \text{diag}(\Phi)$$

$$\begin{aligned}\Rightarrow \text{tr}(H(\lambda)) &= \text{tr} [\Phi (\Phi^T \Phi + \lambda I_P)^{-1} \Phi^T] \\ &= \sum_{j=1}^P (D_\phi^T D_\phi)_{jj} [(D_\phi^T D_\phi)_{jj} + \lambda]^{-1}\end{aligned}$$

1<sup>o</sup> The range of the value of  $\lambda$  can be determined by the expected dof

To avoid overfitting, one could limit the dof of the model.

If we want the maximum dof of the model to be  $v$ :

$$\text{tr}(H(\lambda)) \leq v$$

$$\Rightarrow \sum_{j=1}^P (D_\phi^T D_\phi)_{jj} [(D_\phi^T D_\phi)_{jj} + \lambda]^{-1} \leq v$$

2<sup>o</sup> Find the optimum  $\lambda$  by Cross Validation.