

# Estimation of conditional mean squared error of prediction

Filip Lindskog

based on work with M. Lindholm and F. Wahl forthcoming in Annals of Actuarial Science

October 16, 2019

# Conditional MSEP

- $X$  is a random variable whose value will be observed in the future.
- Current and past observations generate the  $\sigma$ -algebra  $\mathcal{F}_0$
- $\hat{X}$  is an  $\mathcal{F}_0$ -measurable predictor of  $X$
- How accurate is the approximation  $X \approx \hat{X}$  given  $\mathcal{F}_0$ ?
- $\text{MSEP}(X, \hat{X}) := \mathbb{E}[(X - \hat{X})^2]$  is not so relevant
- $\text{MSEP}_{\mathcal{F}_0}(X, \hat{X}) := \mathbb{E}[(X - \hat{X})^2 | \mathcal{F}_0]$  is relevant
- How should  $\text{MSEP}_{\mathcal{F}_0}(X, \hat{X})$  be estimated with limited data?

# Conditional MSEP with parametric model

- Assume a stochastic model with parameter  $\theta$
- Write  $h(\theta; \mathcal{F}_0) := \mathbb{E}[X | \mathcal{F}_0]$  and assume  $h(\theta; \mathcal{F}_0)$  is computable
- Consider predictor  $\hat{X} := h(\hat{\theta}; \mathcal{F}_0)$  given  $\mathcal{F}_0$ -measurable estimator  $\hat{\theta}$
- Calculate

$$\begin{aligned}\text{MSEP}_{\mathcal{F}_0}(X, \hat{X}) &= \mathbb{E}[(X - \hat{X})^2 | \mathcal{F}_0] \\ &= \mathbb{E}[(X - \mathbb{E}[X | \mathcal{F}_0] + \mathbb{E}[X | \mathcal{F}_0] - \hat{X})^2 | \mathcal{F}_0] \\ &= \text{Var}(X | \mathcal{F}_0) + \mathbb{E}[(\hat{X} - \mathbb{E}[X | \mathcal{F}_0])^2 | \mathcal{F}_0] \\ &\quad + 2\mathbb{E}[(X - \mathbb{E}[X | \mathcal{F}_0])(\mathbb{E}[X | \mathcal{F}_0] - \hat{X}) | \mathcal{F}_0] \\ &= \text{Var}(X | \mathcal{F}_0) + (h(\hat{\theta}; \mathcal{F}_0) - h(\theta; \mathcal{F}_0))^2\end{aligned}$$

# Naive estimation of conditional MSEP

- From last page

$$\text{MSEP}_{\mathcal{F}_0}(X, \hat{X}) = \text{Var}(X | \mathcal{F}_0) + (h(\hat{\boldsymbol{\theta}}; \mathcal{F}_0) - h(\boldsymbol{\theta}; \mathcal{F}_0))^2$$

- The naive plug-in estimator obtained by replacing  $\boldsymbol{\theta}$  by  $\hat{\boldsymbol{\theta}}$  is

$$\widehat{\text{MSEP}}_{\mathcal{F}_0}(X, \hat{X}) = \widehat{\text{Var}}(X | \mathcal{F}_0) + 0$$

which is a bad estimator of  $\text{MSEP}_{\mathcal{F}_0}(X, \hat{X})$

- The naive estimator corresponds to ignoring prediction error due to parameter estimation error

# Akaike's FPE

- $X_1, \dots, X_n$  observations from AR( $p_0$ ),

$$X_{t+1} = \phi_1 X_t + \dots + \phi_{p_0} X_{t-p_0+1} + Z_{t+1}, \quad (Z_t) \sim \text{iid}(0, \sigma^2),$$

with unknown  $p_0$ . Determine  $p_0$ !

- Let  $\hat{\phi}$  be an estimator of  $\phi = (\phi_1, \dots, \phi_p)^T$  based on AR( $p$ ) data  $X_1, \dots, X_n$ . Let  $(X_t^\perp)_{t=1}^n$  be an independent copy of  $(X_t)_{t=1}^n$ .

$$\text{FPE}_p(X_{n+1}, \hat{X}_{n+1}) := \mathbb{E} \left[ \left( X_{n+1}^\perp - \sum_{k=1}^p \hat{\phi}_k X_{n+1-k}^\perp \right)^2 \right]$$

Set  $\hat{p}_0 := \arg \min_p \widehat{\text{FPE}}_p(X_{n+1}, \hat{X}_{n+1})$

- Convergence  $\sqrt{n}(\hat{\phi} - \phi) \xrightarrow{d}$  multivariate normal distribution gives

$$\text{FPE}_p(X_{n+1}, \hat{X}_{n+1}) \approx \hat{\sigma}^2(n, p) \frac{n+p}{n-p} \quad \text{for large } n$$

# Akaike's FPE in a more general conditional setting

- Let  $\mathcal{T} := \{\underline{t}, \underline{t} + 1, \dots, \bar{t}\} \subset \mathbb{Z}$  with  $\underline{t} < 0 < \bar{t}$
- Let  $(\mathbf{S}_t)_{t \in \mathcal{T}}, (\mathbf{S}_t^\perp)_{t \in \mathcal{T}}$  be iid copies
- Let  $(\mathcal{F}_t)_{t \in \mathcal{T}}, (\mathcal{F}_t^\perp)_{t \in \mathcal{T}}$  be filtrations generated by  $(\mathbf{S}_t)_{t \in \mathcal{T}}, (\mathbf{S}_t^\perp)_{t \in \mathcal{T}}$
- Let  $X = f((\mathbf{S}_t)_{t \in \mathcal{T}}), X^\perp = f((\mathbf{S}_t^\perp)_{t \in \mathcal{T}})$  for some  $f$
- Let  $\hat{\boldsymbol{\theta}} = g((\mathbf{S}_t)_{t \in \mathcal{T}, t \leq 0}), \hat{\boldsymbol{\theta}}^\perp = g((\mathbf{S}_t^\perp)_{t \in \mathcal{T}, t \leq 0})$  for some  $g$

$$\begin{aligned}\text{FPE}_{\mathcal{F}_0}(X, \hat{X}) &:= \mathbb{E}[(X^\perp - h(\hat{\boldsymbol{\theta}}; \mathcal{F}_0^\perp))^2 \mid \mathcal{F}_0^\perp] \\ &= \mathbb{E}[(X - h(\hat{\boldsymbol{\theta}}^\perp; \mathcal{F}_0))^2 \mid \mathcal{F}_0] \\ &= \text{Var}(X \mid \mathcal{F}_0) + \mathbb{E}[(h(\hat{\boldsymbol{\theta}}^\perp; \mathcal{F}_0) - h(\boldsymbol{\theta}; \mathcal{F}_0))^2 \mid \mathcal{F}_0]\end{aligned}$$

whose plug-in estimator does not make the second term vanish

# Estimation of conditional FPE

- First idea: write  $H(\boldsymbol{\theta}) := \text{FPE}_{\mathcal{F}_0}(X, \widehat{X})$  and estimate by  $H(\widehat{\boldsymbol{\theta}})$
- Problem: We do not know the distribution of  $\widehat{\boldsymbol{\theta}}$ , therefore  $\mathbb{E}[(h(\widehat{\boldsymbol{\theta}}^\perp; \mathcal{F}_0) - h(\boldsymbol{\theta}; \mathcal{F}_0))^2 | \mathcal{F}_0]$  is not computable
- Problem: Unlikely that independent copies  $\widehat{\boldsymbol{\theta}}_1^\perp, \dots, \widehat{\boldsymbol{\theta}}_n^\perp$  are available enabling the approximation of

$$\mathbb{E}[(h(\widehat{\boldsymbol{\theta}}^\perp; \mathcal{F}_0) - h(\boldsymbol{\theta}; \mathcal{F}_0))^2 | \mathcal{F}_0]$$

by

$$\frac{1}{n} \sum_{k=1}^n (h(\widehat{\boldsymbol{\theta}}_k^\perp; \mathcal{F}_0) - h(\boldsymbol{\theta}; \mathcal{F}_0))^2$$

# Estimation of conditional FPE

- Second idea: approximating

$$h(\hat{\boldsymbol{\theta}}^\perp; \mathcal{F}_0) \approx h(\boldsymbol{\theta}; \mathcal{F}_0) + \nabla h(\boldsymbol{\theta}; \mathcal{F}_0)^T (\hat{\boldsymbol{\theta}}^\perp - \boldsymbol{\theta})$$

approximates  $H(\boldsymbol{\theta}) := \text{FPE}_{\mathcal{F}_0}(X, \hat{X})$  by  $H^\nabla(\boldsymbol{\theta})$  given by

$$\text{Var}(X \mid \mathcal{F}_0) + \nabla h(\boldsymbol{\theta}; \mathcal{F}_0)^T \mathbb{E}[(\hat{\boldsymbol{\theta}}^\perp - \boldsymbol{\theta})(\hat{\boldsymbol{\theta}}^\perp - \boldsymbol{\theta})^T] \nabla h(\boldsymbol{\theta}; \mathcal{F}_0)$$

- Estimate  $\text{FPE}_{\mathcal{F}_0}(X, \hat{X})$  and  $\text{MSEP}_{\mathcal{F}_0}(X, \hat{X})$  by  $H^\nabla(\hat{\boldsymbol{\theta}})$
- If  $\mathbb{E}[\hat{\boldsymbol{\theta}}] = \boldsymbol{\theta}$ , then

$$H^\nabla(\boldsymbol{\theta}) = \text{Var}(X \mid \mathcal{F}_0) + \nabla h(\boldsymbol{\theta}; \mathcal{F}_0)^T \text{Cov}(\hat{\boldsymbol{\theta}}) \nabla h(\boldsymbol{\theta}; \mathcal{F}_0)$$

- The only remaining challenge is to estimate  $\text{Cov}(\hat{\boldsymbol{\theta}})$

# Data

$C_{i,j}$  denotes the total amount paid to policyholders during development periods  $1, \dots, j$  due to accidents in accident period  $i$ .

	1	2	3	4	5
$i_0$	$C_{i_0,1}$	$C_{i_0,2}$	$C_{i_0,3}$	$C_{i_0,4}$	$C_{i_0,5}$
$\vdots$	...	...	...	...	...
-5	$C_{-5,1}$	$C_{-5,2}$	$C_{-5,3}$	$C_{-5,4}$	$C_{-5,5}$
-4	$C_{-4,1}$	$C_{-4,2}$	$C_{-4,3}$	$C_{-4,4}$	$C_{-4,5}$
-3	$C_{-3,1}$	$C_{-3,2}$	$C_{-3,3}$	$C_{-3,4}$	$C_{-3,5}$
-2	$C_{-2,1}$	$C_{-2,2}$	$C_{-2,3}$	$C_{-2,4}$	$C_{-2,5}$
-1	$C_{-1,1}$	$C_{-1,2}$	$C_{-1,3}$	$C_{-1,4}$	$C_{-1,5}$

$$S_0 = (C_{-1,1}, \dots, C_{-5,5})^T, S_{-1} = (C_{-2,1}, \dots, C_{-6,5})^T, \dots, S_{i_0} = C_{i_0,1},$$
$$X = C_{-1,5} + \dots + C_{-4,5}$$

## Development-period (column) dynamics

- For vector-valued functions  $\mathbf{F}_{i,j}$  and scalar-valued functions  $G_{i,j}$ ,

$$C_{i,j+1} = \mathbf{F}_{i,j}\left((C_{i,k})_{k \leq j}\right)^T \boldsymbol{\beta}_j + \sigma_j G_{i,j}\left((C_{i,k})_{k \leq j}\right) e_{i,j+1},$$

$\{C_{i,1}, e_{i,2}, \dots, e_{i,J}\}$  and  $\{C_{k,1}, e_{k,2}, \dots, e_{k,J}\}$  independent if  $i \neq k$ ,

$$\mathbb{E}[e_{i,j+1} | C_1, \dots, C_j] = 0, \quad \mathbb{E}[e_{i,j+1}^2 | C_1, \dots, C_j] = 1$$

Note:  $e_{i,1}, \dots, e_{i,J}$  not necessarily independent

- Development-period dynamics as vector-valued time series:

$$\mathbf{C}_{j+1} = \mathbf{A}_j \boldsymbol{\beta}_j + \sigma_j \mathbf{D}_j \mathbf{e}_{j+1}$$

- Mack's Chain Ladder:  $C_{i,j+1} = C_{i,j} \boldsymbol{\beta}_j + \sigma_j \sqrt{C_{i,j}} e_{i,j+1}$  when  
 $\mathbf{F}_{i,j}\left((C_{i,k})_{k \leq j}\right) = \mathbf{F}\left((C_{i,k})_{k \leq j}\right) = C_{i,j}$  and  
 $G_{i,j}\left((C_{i,k})_{k \leq j}\right) = G\left((C_{i,k})_{k \leq j}\right) = \sqrt{C_{i,j}}$

# Observed data

	1	2	3	4	5
$i_0$	$C_{i_0,1}$	$C_{i_0,2}$	$C_{i_0,3}$	$C_{i_0,4}$	$C_{i_0,5}$
$\vdots$	...	...	...	...	...
-5	$C_{-5,1}$	$C_{-5,2}$	$C_{-5,3}$	$C_{-5,4}$	$C_{-5,5}$
-4	$C_{-4,1}$	$C_{-4,2}$	$C_{-4,3}$	$C_{-4,4}$	$C_{-4,5}$
-3	$C_{-3,1}$	$C_{-3,2}$	$C_{-3,3}$	$C_{-3,4}$	$C_{-3,5}$
-2	$C_{-2,1}$	$C_{-2,2}$	$C_{-2,3}$	$C_{-2,4}$	$C_{-2,5}$
-1	$C_{-1,1}$	$C_{-1,2}$	$C_{-1,3}$	$C_{-1,4}$	$C_{-1,5}$

Example:  $\tilde{\mathbf{C}}_4 = \tilde{\mathbf{A}}_3 \beta_3 + \sigma_3 \tilde{\mathbf{D}}_3 \tilde{\mathbf{e}}_4$

# Weighted least-squares estimation

The estimators

$$\hat{\beta}_j := \left( \tilde{\mathbf{A}}_j^T \tilde{\Sigma}_j^{-1} \tilde{\mathbf{A}}_j \right)^{-1} \tilde{\mathbf{A}}_j^T \tilde{\Sigma}_j^{-1} \tilde{\mathbf{C}}_{j+1}, \quad \tilde{\Sigma}_j := \tilde{\mathbf{D}}_j^2, \quad (\text{WLS})$$

$$\hat{\sigma}_j^2 := \frac{1}{n_j - p_j} \left( \tilde{\mathbf{C}}_{j+1} - \tilde{\mathbf{A}}_j \hat{\beta}_j \right)^T \tilde{\Sigma}_j^{-1} \left( \tilde{\mathbf{C}}_{j+1} - \tilde{\mathbf{A}}_j \hat{\beta}_j \right)$$

lead to

- $\mathbb{E}[\hat{\beta}_j | \tilde{\mathbf{A}}_j, \tilde{\Sigma}_j] = \mathbb{E}[\hat{\beta}_j] = \beta_j$
- $\mathbb{E}[\hat{\sigma}_j^2 | \tilde{\mathbf{A}}_j, \tilde{\Sigma}_j] = \mathbb{E}[\hat{\sigma}_j^2] = \sigma_j^2$
- $\text{Cov}(\hat{\beta}_j, \hat{\beta}_k) = \mathbf{0}$  if  $j \neq k$
- $\text{Cov}(\hat{\beta}_j | \tilde{\mathbf{A}}_j, \tilde{\Sigma}_j) = \sigma_j^2 (\tilde{\mathbf{A}}_j^T \tilde{\Sigma}_j^{-1} \tilde{\mathbf{A}}_j)^{-1}$

# Estimating covariance matrix of regression coefficients

The estimator  $\Lambda(\hat{\sigma}^2; \mathcal{F}_0)$  of  $\text{Cov}(\hat{\beta})$  given by

$$\begin{bmatrix} \hat{\sigma}_1^2 (\tilde{\mathbf{A}}_1^T \tilde{\boldsymbol{\Sigma}}_1^{-1} \tilde{\mathbf{A}}_1)^{-1} & 0 & \dots & 0 \\ 0 & \ddots & & \\ \vdots & & & \\ 0 & & \hat{\sigma}_{J-1}^2 (\tilde{\mathbf{A}}_{J-1}^T \tilde{\boldsymbol{\Sigma}}_{J-1}^{-1} \tilde{\mathbf{A}}_{J-1})^{-1} & \end{bmatrix}$$

satisfies

$$\mathbb{E}[\Lambda(\hat{\sigma}^2; \mathcal{F}_0)] = \text{Cov}(\hat{\beta}), \quad \hat{\beta} := (\hat{\beta}_1^T, \dots, \hat{\beta}_{J-1}^T)^T$$

which can be shown using covariance decomposition

$$\text{Cov}(\hat{\beta}_j) = \text{Cov}(\mathbb{E}[\hat{\beta}_j | \mathcal{H}_j]) + \mathbb{E}[\text{Cov}(\hat{\beta}_j | \mathcal{H}_j)]$$

for suitably chosen  $\mathcal{H}_j$  together with conditional unbiasedness of  $\hat{\beta}_j$ .

## Putting the pieces together

Recall that we want to estimate

$$\text{Var}(X \mid \mathcal{F}_0) + \nabla h(\boldsymbol{\theta}; \mathcal{F}_0)^T \text{Cov}(\hat{\boldsymbol{\theta}}) \nabla h(\boldsymbol{\theta}; \mathcal{F}_0)$$

which, for the considered model class with  $\boldsymbol{\theta} = (\boldsymbol{\beta}, \sigma^2)$  and  $X$  is such that  $h(\boldsymbol{\beta}, \sigma^2; \mathcal{F}_0) = \mathbb{E}[X \mid \mathcal{F}_0]$  does not depend on  $\sigma^2$ , equals

$$\text{Var}(X \mid \mathcal{F}_0) + \nabla_{\boldsymbol{\beta}} h(\boldsymbol{\beta}, \sigma^2; \mathcal{F}_0)^T \text{Cov}(\hat{\boldsymbol{\beta}}) \nabla_{\boldsymbol{\beta}} h(\boldsymbol{\beta}, \sigma^2; \mathcal{F}_0)$$

which is estimated by

$$\widehat{\text{Var}(X \mid \mathcal{F}_0)} + \nabla_{\boldsymbol{\beta}} h(\hat{\boldsymbol{\beta}}, \hat{\sigma}^2; \mathcal{F}_0)^T \boldsymbol{\Lambda}(\hat{\sigma}^2; \mathcal{F}_0) \nabla_{\boldsymbol{\beta}} h(\hat{\boldsymbol{\beta}}, \hat{\sigma}^2; \mathcal{F}_0)$$

Tedious calculations show that the above expression coincides with estimators of conditional MSEP proposed for special models, including Mack's Chain Ladder.

# Generalisation of conditional FPE

- We estimated  $\text{MSEP}_{\mathcal{F}_0}(X, \hat{X})$  by estimating

$$\text{FPE}_{\mathcal{F}_0}(X, \hat{X}) = \mathbb{E}[(X - h(\hat{\boldsymbol{\theta}}^\perp; \mathcal{F}_0))^2 \mid \mathcal{F}_0]$$

Why  $\hat{\boldsymbol{\theta}}^\perp$ ? Any  $\hat{\boldsymbol{\theta}}^*$  conditionally independent of  $X$  given  $\mathcal{F}_0$  makes most of the analysis stay the same when considering

$$\text{MSEP}_{\mathcal{F}_0}^*(X, \hat{X}) = \mathbb{E}[(X - h(\hat{\boldsymbol{\theta}}^*; \mathcal{F}_0))^2 \mid \mathcal{F}_0]$$

- $\text{MSEP}_{\mathcal{F}_0}^*(X, \hat{X})$  together with development-period time-series dynamics ( $\mathbf{C}_j$ ) with iid residuals ( $\mathbf{e}_j$ ) with iid components

$$\mathbf{C}_{j+1} = \mathbf{A}_j \boldsymbol{\beta}_j + \sigma_j \mathbf{D}_j \mathbf{e}_{j+1}$$

allows for estimation of  $\text{MSEP}_{\mathcal{F}_0}^*(X, \hat{X})$  by parametric bootstrap.  
 $\text{FPE}_{\mathcal{F}_0}(X, \hat{X})$  does not.

## Natural next step

Recalling the originally intended use of Akaike's FPE, a natural next step would be to...

...use  $\text{MSEP}_{\mathcal{F}_0}^*(X, \hat{X})$  together with machine learning ideas for automatic predictor selection including assessment of prediction accuracy.