# LINEAR ESTIMATION

Linear estimators are an important class of estimators

- simplicity
- dependence on first and second order moments only
- ease of implementation

This last point is especially important for filtering problems where we must process data real-time.

Our discussion of filtering will rely on Bayesian linear estimators, but for completeness we begin with classical linear estimators.

We make the following distinctions:

constant: $\qquad \hat{\underline{\theta}}(\underline{x}) = \underline{b}$

linear: $\qquad \hat{\underline{\theta}}(\underline{x}) = H\underline{x}$

affine: $\qquad \hat{\underline{\theta}}(\underline{x}) = H\underline{x} + \underline{b}$

where $H \in R^{P \times N}$, $\underline{b} \in R^{P}$. We'll study all three cases.

## Best Linear Unbiased Estimation

The MVUE is often not computable.

- CRLB or Rao-Blackwell not applicable
- intractable mathematical model
- randomness is only known up to first and second order moments.

In such cases, we must be content with a suboptimal estimator. One approach is to compute

__Definition__ | The <u>best linear unbiased estimator</u> (BLUE) is the linear estimator

$$\hat{\underline{\theta}}(\underline{x}) = H\underline{x} \qquad , \qquad H \in \mathbb{R}^{P \times N}$$

with smallest variance among all linear, unbiased estimators.

Note that for $\hat{\underline{\theta}}$ to be unbiased we must have

$$\underline{\theta} = E\{\hat{\underline{\theta}}\} = E\{H\underline{x}\} = H E\{\underline{x}\}$$

Therefore the mean of the data must obey a linear relationship with the true parameter.

This relationship will not always be true, so even the suboptimal BLUE isn't always feasible.

However, there is an important class of problems where it does hold.

## Linear Model

Suppose $\underline{X}$ and $\underline{\theta}$ are related through

$$\underline{X} = A\underline{\theta} + \underline{W}$$

> $\underline{\theta}$ fixed but unknown

where

$A$ is $N \times p$ and known, full rank

$\underline{W}$ is random

$E[\underline{W}] = \underline{0}$

$E[\underline{W}\,\underline{W}^T] =: R$   is known, positive definite

Note: $\underline{W}$ is not necessarily Gaussian

Question: Can you give an $H$ such that

$$H\,E[\underline{X}] = \underline{\theta} \text{ ?}$$

Taking $H = (A^T A)^{-1} A^T$ we find

$$H E[\underline{x}] = \underbrace{(A^T A)^{-1} A^T A}_{I} \underline{\theta} + (A^T A)^{-1} A^T \underbrace{E[\underline{w}]}_{\underline{0}}$$

$$= \underline{\theta}.$$

So the pseudo inverse is <u>an</u> unbiased estimator.
But is its variance minimal?

<u>Theorem</u> (Gauss-Markov Theorem)

In the linear model described above, the BLUE is

$$\hat{\underline{\theta}}(\underline{x}) = (A^T R^{-1} A)^{-1} A^T R^{-1} \underline{x}$$

and its covariance matrix is

$$R_{\hat{\theta}} = E[(\hat{\underline{\theta}} - \underline{\theta})(\hat{\underline{\theta}} - \underline{\theta})^T] = (A^T R^{-1} A)^{-1}$$

**Proof 1** $\hat{\theta}$ is unbiased $\Longleftrightarrow$ $E\{H\underline{x}\} = \underline{\theta}$ $\forall\underline{\theta}$.

Now $E\{H\underline{x}\} = HE\{\underline{x}\} = HE\{A\underline{\theta} + \underline{w}\} = HA\underline{\theta}$.

Thus $\hat{\underline{\theta}}$ is unbiased $\Longleftrightarrow$ $HA\underline{\theta} = \underline{\theta}$ $\forall\underline{\theta}$ $\Longleftrightarrow$ $HA = I_{p\times p}$.

Now assume $\hat{\underline{\theta}}$ is unbiased and let's compute

$$Var(\hat{\underline{\theta}}) = E\left\{(\hat{\underline{\theta}} - \underline{\theta})^T(\hat{\underline{\theta}} - \underline{\theta})\right\}$$

$$= E\left\{(H\underline{x} - \underline{\theta})^T(H\underline{x} - \underline{\theta})\right\}$$

$$= E\left\{(HA\underline{\theta} + H\underline{w} - \underline{\theta})^T(HA\underline{\theta} + H\underline{w} - \underline{\theta})\right\}$$

$$= E\left\{(H\underline{w})^T(H\underline{w})\right\}.$$

Denote

$$H = \begin{bmatrix} \underline{h}_1 & \cdots & \underline{h}_p \end{bmatrix}^T = \begin{bmatrix} \underline{h}_1^T \\ \vdots \\ \underline{h}_p^T \end{bmatrix} \quad (p \times N)$$

Then

$$Var(\hat{\underline{\theta}}) = E\left\{\sum_{i=1}^{P}(\underline{h}_i^T\underline{w})^2\right\}$$

$$= \sum_{i=1}^{P} E\left\{(\underline{h}_i^T\underline{w})^2\right\}$$

$$= \sum_{i=1}^{P} E\left\{(\underline{h}_i^T\underline{w})(\underline{w}^T\underline{h}_i)\right\}$$

$$= \sum_{i=1}^{P} \underline{h}_i^T R \underline{h}_i$$

Thus, we need to solve the constrained optimization problem

$$\min_{H} \quad \sum_{i=1}^{P} \underline{h}_i^T R \, \underline{h}_i$$

$$st \quad H \cdot A = I$$

By the theory of Lagrange multipliers, it suffices to solve the unconstrained problem

$$\min_{H, \underline{\lambda}} \quad L(H, \underline{\lambda})$$

where $L$ is the Lagrangian and $\underline{\lambda}$ is a vector of real numbers called Lagrange multipliers, one for each equality constraint.

The Lagrangian is

$$L = \sum_{i=1}^{P} \underline{h}_i^T R \underline{h}_i + \sum_{i=1}^{P} \sum_{j=1}^{P} \lambda_j^{(i)} (\underline{h}_i^T \underline{a}_j - \delta_{ij})$$

Taking derivatives

$$\frac{\partial L}{\partial \underline{h}_i} = 2R\underline{h}_i + \sum_{j=1}^{P} \lambda_j^{(i)} \underline{a}_j$$

$$= 2R\underline{h}_i + A\underline{\lambda}^{(i)}$$

where $\underline{\lambda}^{(i)} = [\lambda_1^{(i)} \cdots \lambda_P^{(i)}]^T$.

$$\Rightarrow \hat{\underline{h}}_i = -\frac{1}{2} R^{-1} A \underline{\lambda}^{(i)}$$

From the constraint we know

ith position

$$A^T \hat{\underline{h}}_i = \underline{e}_i = [0 \cdots 1 \cdots 0]$$

$$\Rightarrow \underline{e}_i = -\frac{1}{2} A^T R^{-1} A \underline{\lambda}^{(i)}$$

$$\Rightarrow \underline{\lambda}^{(i)} = -2 (A^T R^{-1} A)^{-1} \underline{e}_i$$

$$\Rightarrow \hat{\underline{h}}_i = R^{-1} A (A^T R^{-1} A) \underline{e}_i$$

Assembling these scalar results, we have

$$\hat{H} = [\hat{\underline{h}}_1 \cdots \hat{\underline{h}}_p]^T$$

$$= \left( R^{-1}A(A^TR^{-1}A)^{-1} \cdot I_{p \times p} \right)^T$$

$$= (A^TR^{-1}A)^{-1}A^TR$$

The covariance matrix of $\hat{\underline{\theta}}$ is

$$\text{Cov}(\hat{\underline{\theta}}) =$$

**Proof 2]** From proof 1, we know that the BLUE depends only on the first and second order moments of $\underline{x}$ (or equivalently, of $\underline{w}$).

Therefore, we may assume

$$\underline{w} \sim N(\underline{0}, R).$$

We have previously seen that

$$\hat{\theta}(\underline{x}) = (A^T R^{-1} A)^{-1} A^T R^{-1} \underline{x}$$

is MVUE. Since this estimator is already linear, it is also the BLUE ▨

**Remark]** As the above discussion notes, when $\underline{w}$ is Gaussian, the BLUE is the MVUE for the linear model, i.e. the BLUE is optimal.

## Linear Minimum Mean Squared Error Estimation

Let us now turn to a Bayesian setting. There are some important similarities/differences w.r.t. classical linear estimation:

- the optimal linear estimator depends only on first and second order moments, but now for $\underline{X}$ and $\underline{\theta}$

- the optimal linear estimator always exists in the Bayesian setting

- Bayesian linear estimation has important geometric interpretations in terms of orthogonal projections in Hilbert space

<u>Definition</u> $\hat{\underline{\theta}}(\underline{x}) = \hat{H}\underline{x}$ is the LMMSE estimator if $\hat{H}$ minimizes

$$BMSE(H) := E_{X,\theta}\left[(\underline{\theta} - H\underline{X})^T(\underline{\theta} - H\underline{X})\right]$$

Introduce the notation

$$R_{\theta\theta} = E\left\{ (\underline{\theta} - E\underline{\theta})(\underline{\theta} - E\underline{\theta})^T \right\}$$
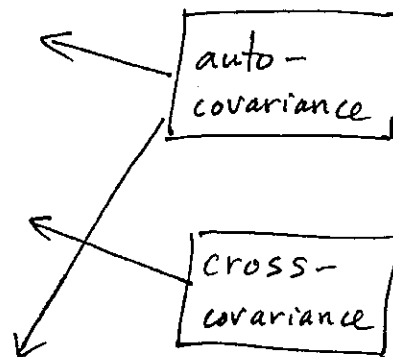
$$(p \times p)$$

$$R_{\theta x} = E\left\{ (\underline{\theta} - E\underline{\theta})(\underline{x} - E\underline{x})^T \right\}$$

$$(p \times N)$$

$$R_{xx} = E\left\{ (\underline{x} - E\underline{x})(\underline{x} - E\underline{x})^T \right\}$$

$$(N \times N)$$

auto-covariance

cross-covariance

__Theorem__ | If $E\underline{\theta} = \underline{0}$ and $E\underline{x} = \underline{0}$

then the LMMSE is

$$\hat{\underline{\theta}}(\underline{x}) = R_{\theta x} R_{xx}^{-1} \underline{x}$$

provided $R_{xx}$ is positive definite.

**Proof** :

$$BMSE(H) = E[(H\underline{x} - \underline{\theta})^T (H\underline{x} - \underline{\theta})]$$

$$= E\left[ \sum_{i=1}^{P} (\underline{h}_i^T \underline{x} - \theta_i)^2 \right]$$

$$= \sum_{i=1}^{P} E\left[ (\underline{h}_i^T \underline{x} - \theta_i)^2 \right]$$

To minimize this term w.r.t $H = [\underline{h}_1 \cdots \underline{h}_P]^T$, it suffices to optimize each $\underline{h}_i$ independently. That is, the vector linear estimation problem reduces to the scalar linear estimation problem.

Thus, let's drop the subscript $i$ and focus on minimizing $E[(\underline{h}^T \underline{x} - \theta)]$ w.r.t. $\underline{h}$.

Now

$$R_{\theta\theta} = Var(\theta) \quad \text{is} \quad 1 \times 1$$

$$R_{\theta x} \quad\quad\quad \text{is} \quad 1 \times N$$

Now

$$BMSE(\underline{h}) = E_{\underline{x},\theta}\left[(\theta - \underline{h}^T\underline{x})^2\right]$$

$$= E[\theta^2] - 2E[\underline{h}^T\underline{x}\cdot\theta] + E[(\underline{h}^T\underline{x})^2]$$

$$= E[\theta^2] - 2\underline{h}^TE[\underline{x}\cdot\theta] + E[\underline{h}^T\underline{x}\cdot\underline{x}^T\underline{h}]$$

$$= R_{\theta\theta} - 2\underline{h}^T R_{x\theta} + \underline{h}^T R_{xx}\underline{h} \qquad \circledast$$

By completing the square, we have

$$BMSE(\underline{h}) = \left(\underline{h} - R_{xx}^{-1} R_{x\theta}\right)^T R_{xx} \left(\underline{h} - R_{xx}^{-1} R_{x\theta}\right)$$

$$- \underbrace{R_{x\theta}^T R_{xx}^{-1} R_{x\theta} + R_{\theta\theta}}_{\text{independent of } \underline{h}}.$$

Since $R_{xx}$ is positive definite, the unique minimizer is

$$\hat{\underline{h}} = R_{xx}^{-1} R_{x\theta}$$

$$\Updownarrow$$

$$\hat{\underline{h}}^T = R_{\theta x} R_{xx}^{-1}.$$

Alternatively, you could apply vector calculus to minimize $\circledast$.

Now the vector LMMSE estimator is obtained by stacking these scalar estimators. 🔲

## Nonzero Means and Affine Estimators

When $E\underline{\theta} \neq \underline{0}$ or $E\underline{x} \neq \underline{0}$, the LMMSE is generalized by the affine MMSE estimator (AMMSE), which is defined to be the estimator $\hat{\underline{\theta}}(\underline{x}) = \hat{H}\underline{x} + \hat{\underline{b}}$

minimizing

$$BMSE(H, \underline{b}) = E\left[(\underline{\theta} - H\underline{x} - \underline{b})^T (\underline{\theta} - H\underline{x} - \underline{b})\right].$$

__Theorem__ The AMMSE is

$$\hat{\underline{\theta}}(\underline{x}) = E\underline{\theta} + R_{\theta x} \cdot R_{xx}^{-1}(\underline{x} - E\underline{x})$$

i.e.

$$\hat{H} = R_{\theta x} R_{xx}^{-1}$$

$$\hat{\underline{b}} = E\underline{\theta} - \hat{H} E\underline{x}$$

__Proof__ Introduce the variables

Skip in lecture

$$\underline{\theta}' = \underline{\theta} - E\underline{\theta}$$

$$\underline{x}' = \underline{x} - E\underline{x}$$

Then $BMSE(H, \underline{b})$

$$= E\left[(\underline{\theta}' - H\underline{x}' - (\underline{b} - E\underline{\theta} + HE\underline{x}))^T \cdot \right.$$

$$\left. (\underline{\theta}' - H\underline{x}' - (\underline{b} - E\underline{\theta} + HE\underline{x}))\right]$$

$$= E\left[(\underline{\theta}' - H\underline{x}')^T(\underline{\theta}' - H\underline{x}')\right]$$

$$-2E\left[(\underline{\theta}' - H\underline{x}')^T(\underline{b} - E\underline{\theta} + HE\underline{x})\right]$$

$$+ E\left[(\underline{b} - E\underline{\theta} + HE\underline{x})^T(\underline{b} - E\underline{\theta} + HE\underline{x})\right]$$

The second term is zero since $\underline{b} - E\underline{\theta} + HE\underline{x}$ is constant and $\underline{\theta}' - H\underline{x}'$ is zero mean.

For fixed $H$, the optimal $\underline{b}$ must minimize the third term, which can be made zero by taking

$$\underline{b} = E\underline{\theta} - HE\underline{x}. \qquad \circledast$$

Therefore, $H$ is obtained by minimizing the first term. Since $\underline{\theta}', \underline{x}'$ are zero mean, $\hat{H}$ is the LMMSE estimator for $\underline{\theta}', \underline{x}'$, i.e.

$$\hat{H} = R_{\theta'x'} R_{x'x'}^{-1} = R_{\theta x} R_{xx}^{-1} \qquad \boxed{\text{2}}$$

Remark 1  From the above argument and $\circledast$, we see that the optimal constant estimator $(H = 0)$ is

$$\hat{\underline{b}} = E\underline{\theta},$$

the prior mean.

# Connection to the Jointly Gaussian Case

The LMMSE and AMMSE estimators look exactly like the MMSE estimators for the case where $\theta, \underline{X}$ are jointly Gaussian. In fact, we can use the Gaussian case to give an alternate derivation of the LMMSE/AMMSE.

Let's consider the AMMSE. From the definition of the Bayesian MSE

$$BMSE(H, \underline{b}) = E\left[(\underline{\theta} - H\underline{X} - \underline{b})^T (\underline{\theta} - H\underline{X} - \underline{b})\right]$$

is can be seen that this criterion depends only on the first and second order moments of $\underline{\theta}$ and $\underline{X}$. Therefore, we can assume the higher order moments are whatever we want.

So let's assume $\underline{\theta}$ & $\underline{X}$ are jointly Gaussian with the given means and covariances:

$$\begin{bmatrix} \underline{X} \\ \underline{\theta} \end{bmatrix} \sim N\left( \begin{bmatrix} E\underline{X} \\ E\underline{\theta} \end{bmatrix}, \begin{bmatrix} R_{XX} & R_{X\theta} \\ R_{\theta X} & R_{\theta\theta} \end{bmatrix} \right)$$

We know the MMSE estimator is the posterior mean, which is

$$E\underline{\theta} + R_{\theta x} R_{xx}^{-1} (\underline{x} - E\underline{x})$$

by the Gaussian conditioning principle.

Since this estimator is affine, it is the AMMSE ∎

**Theorem** | Bayesian Gauss-Markov Theorem

Assume

$$\underline{X} = A\underline{\theta} + \underline{W}$$

where

$A$ is full rank, $N \times p$, known

$E\underline{\theta}$, $R_{\theta\theta}$ known

$E\underline{W} = \underline{0}$, $R_{ww}$ known

$R_{\theta w} = 0_{p \times N}$ ($\underline{\theta}$ & $\underline{w}$ uncorrelated)

Then the affine MMSE estimator is

$$\hat{\underline{\theta}}(\underline{x}) = E\underline{\theta} + R_{\theta\theta}A^T\left(AR_{\theta\theta}A^T + R_{ww}\right)^{-1}(\underline{x} - AE\underline{\theta})$$

$$= E\underline{\theta} + \left(R_{\theta\theta}^{-1} + A^TR_{ww}^{-1}A\right)^{-1}A^TR_{ww}^{-1}(\underline{x} - AE\underline{\theta})$$

How would you prove this?

**Proofs**

1. Show $R_{\theta x} =$

(a)

   $R_{xx} =$

   and apply the AMMSE result.

2. Argue that the BMSE depends only on the first and second order moments of $\underline{\theta}$ and $\underline{W}$. Then assume $\underline{\theta}, \underline{W}$ are jointly Gaussian. The MMSE estimator has the stated form, and since it is affine it is the AMMSE.

**Exercise**    Suppose

$$X_i = A + W_i, \qquad i = 1, \ldots, N$$

where

$$A \sim \text{unif}(-A_0, A_0)$$
$$W_i \stackrel{iid}{\sim} N(0, \sigma_W^2)$$

} independent

1. Find the posterior $f(A/\underline{x})$

2. Can you compute $E[A/\underline{x}]$?

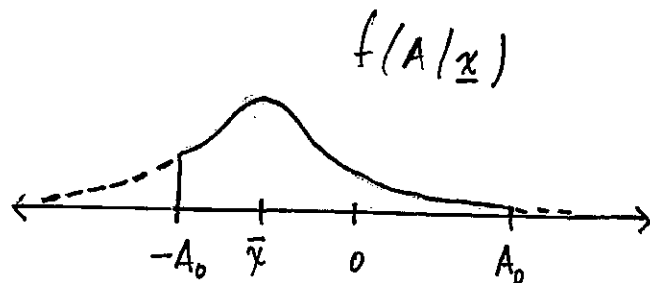3. Determine the LMMSE estimator of $A$.

<u>Solution</u>

1.  $f(A|\underline{x}) \propto f(\underline{x}|A) \cdot f(A)$

$$= \prod_{i=1}^{N} \phi(x_i | A, \sigma^2) \cdot \frac{1}{2A_0} I_{[-A_0, A_0]}(A)$$

$$\propto \exp\left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^{N} (x_i - A)^2 \right\}$$

$$= \exp\left\{ -\frac{1}{2\sigma^2} \left[ \sum (x_i - \bar{x})^2 + N(\bar{x} - A)^2 \right] \right\}$$

$$\Rightarrow f(A|\underline{x}) \propto \exp\left\{ \frac{(A - \bar{x})^2}{2(\sigma^2/N)} \right\} \cdot I_{[-A_0, A_0]}(A)$$

$f(A|\underline{x})$



Truncated normal.

2. $\quad E[A | \underline{x}] = \int_{-\infty}^{\infty} \overset{\circ}{A} f(A / \underline{x}) dA$

$$= \frac{\int_{-A_0}^{A_0} A \phi(A | \bar{x}, \frac{\sigma^2}{N}) dA}{\int_{-A_0}^{A_0} \phi(A | \bar{x}, \frac{\sigma^2}{N}) dA} \quad \longleftarrow \boxed{\begin{array}{l} Normalization \\ Constant \end{array}}$$

Numerator : closed form solution

Denominator : no closed form solution

Note : $E[A | \underline{x}]$ is non linear.

3.

$$\underline{X} = \underline{1} \cdot A + \underline{W}, \quad \underline{W} \sim N(\underline{0}, \sigma_w^2 I)$$

$$\sigma_A^2 = \int_{-A_0}^{A_0} \frac{x^2}{2 A_0} dx = \frac{1}{6 A_0} x^3 \Big|_{-A_0}^{A_0} = \frac{A_0^2}{3}$$

$$\Longrightarrow \hat{A} = \left( (\sigma_A^2)^{-1} + \left( \underline{1}^T \sigma_w^2 I \, \underline{1} \right)^{-1} \right)^{-1} (\sigma_w^2 I)^{-1} \underline{1}^T \underline{x}$$

$$= \frac{\sigma_A^2}{\sigma_A^2 + \frac{\sigma_w^2}{N}} \bar{x}$$

$$= \frac{A_0^2/3}{A_0^2/3 + \sigma^2/N} \bar{x}$$

<u>Remarks</u>

a) To compute the LMMSE estimator of $A$, we didn't need to know it was uniform, just its mean and variance.

b) Estimator is the same as the MMSE estimator for the prior

$$A \sim N\left(0, \frac{A_o^2}{3}\right)$$

<u>The Orthogonality Principle</u>

The LMMSE estimator satisfies

$$R_{\theta x} = \underline{\hat{h}}^T R_{xx}$$

or equivalently

$$\underline{0} = R_{\theta x} - \underline{\hat{h}}^T R_{xx}$$

$$= E\left[\theta \underline{x}^T - \underline{\hat{h}}^T \underline{x} \, \underline{x}^T\right]$$

$$= E\left[(\theta - \underline{\hat{h}}^T \underline{x}) \underline{x}^T\right]$$

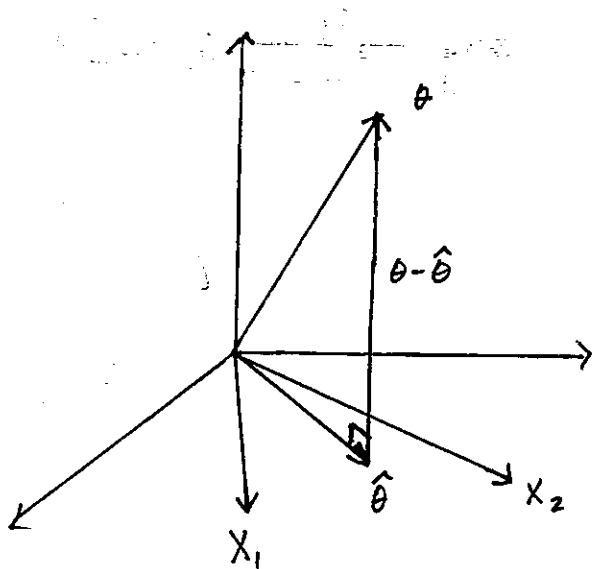$$= E\left[(\theta - \hat{\theta}) \underline{x}^T\right]$$

In words, this says that the prediction error is orthogonal to every measurement variable,

$$\theta - \hat{\theta} \perp X_i \qquad \forall i$$

where orthogonality is with respect to the inner product given by expectation.

As a consequence,

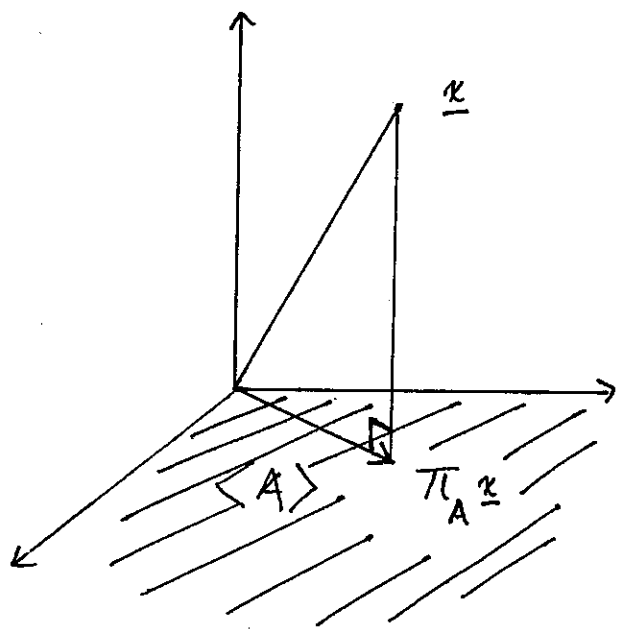$$\theta - \hat{\theta} \perp \underline{h}^T \underline{x} \qquad \forall \underline{h}.$$



Therefore, the prediction error is orthogonal to the entire <u>subspace of linear estimators.</u>

This is one manifestation of the $\boxed{\text{orthogonality principle.}}$

The orthogonality principle is also manifested in nonstatistical projections

In particular, if $\underline{x} \in \mathbb{R}^N$ and $\langle A \rangle \subseteq \mathbb{R}^N$, then

$$\underline{x} - \Pi_A \underline{x} \quad \perp \quad \underline{u} \qquad \forall \underline{u} \in \langle A \rangle$$



Here the inner product is the standard dot product.

To prove this fact, suppose $\underline{u} \in \langle A \rangle$.
Then $\underline{u} = A\underline{\phi}$ for some $\underline{\phi}$,

$$\langle \underline{x} - \Pi_A \underline{x}, \underline{u} \rangle = \underline{u}^T (\underline{x} - \Pi_A \underline{x})$$

$$= \underline{\phi}^T \underline{A}^T \left( \underline{x} - A(A^TA)^{-1} A^T \underline{x} \right)$$

$$= \underline{\phi}^T A^T \underline{x} - \underline{\phi}^T A^T \underline{x} = 0$$

In full generality...

**Definition]** Let $\mathcal{H}$ be an inner product space, and let $S \subseteq \mathcal{H}$ be a linear subspace. The _orthogonal complement_ of $S$ is

$$S^{\perp} = \{\, \underline{v} \in \mathcal{H} \mid \underline{v} \perp \underline{u} \;\; \forall \underline{u} \in S \,\}$$

**Theorem]** Let $\mathcal{H}$ be a Hilbert space and $S \subseteq \mathcal{H}$ a closed linear subspace.

1. _Projection Theorem_ : $\mathcal{H} = S \oplus S^{\perp}$, i.e. $\forall \underline{w} \in \mathcal{H}$, $\exists$ unique $\underline{u} \in S$, $\underline{v} \in S^{\perp}$ s.t. $\underline{w} = \underline{u} + \underline{v}$. Thus we can define the (orthogonal) projection onto $S$,

   $$\Pi_S(\underline{w}) = \underline{u}.$$

2. _Closest point property_ : $\Pi_S \underline{w}$ is the unique solution of $\displaystyle\min_{\underline{u} \in S} \| \underline{w} - \underline{u} \|$ where $\| \cdot \|$ is the norm induced by $\langle \cdot, \cdot \rangle$.

3. _Orthogonality principle_ : For all $\underline{w} \in \mathcal{H}$,

   $$\underline{w} - \Pi_S \underline{w} \;\perp\; \underline{u} \qquad \forall \underline{u} \in S$$

   i.e., $\underline{w} - \Pi_S \underline{w} \in S^{\perp}$.

**Note]** 2 & 3 follow easily from 1. See _Moon and Stirling_ for details.

# Application to LMMSE Estimation

Consider the space $\mathcal{H}$ of all scalar random variables with zero mean and finite variance. It can be shown that $\mathcal{H}$ is a Hilbert space with inner product

$$\langle V_1, V_2 \rangle = E\{V_1 \cdot V_2\}$$

Let $X_1, \ldots, X_N$ be random measurements and define the subspace

$$S = \left\{ \underline{h}^T \underline{X} \mid \underline{h} \in \mathbb{R}^N \right\}$$

where $\underline{X} = [X_1 \cdots X_N]^T$.

If $\theta \in \mathcal{H}$ is a scalar parameter of interest, the LMMSE estimator is the __projection__

$$\hat{\theta} = \Pi_S \, \theta$$

$$= \underset{\hat{\theta} \in S}{\arg\min} \; \| \theta - \hat{\theta} \|^2$$

$$= \underset{\underline{h} \in \mathbb{R}^N}{\arg\min} \; E\left[ (\theta - \underline{h}^T \underline{X})^2 \right]$$

By the orthogonality principle,

$$\theta - \hat{\theta} \perp u \qquad \forall u \in S.$$

Writing $\hat{\theta} = \underline{\hat{h}}^T \underline{x}$ and taking $u = X_i$, $i = 1, \dots, N$

we have

$$\theta - \hat{\theta} \perp X_i \iff \theta - \underline{\hat{h}}^T \underline{x} \perp X_i$$

$$\iff E\left[(\theta - \underline{\hat{h}}^T \underline{x}) X_i\right] = 0$$

Applying this for $i = 1, \dots, N$ we have

$$E\left[(\theta - \underline{\hat{h}}^T \underline{x}) \underline{x}^T\right] = [0 \ \cdots \ 0]$$

$$\Updownarrow$$

$$R_{\theta x} = \underline{\hat{h}}^T R_{xx}$$

$$\Updownarrow$$

$$\underline{\hat{h}}^T = R_{\theta x} R_{xx}^{-1}$$

**Conclusion**: The orthogonality principle gives us another proof of the form of the LMMSE estimator.

We will apply the orthogonality principle often in our study of filtering. Furthermore, the orthogonality principle applies when there are an <u>infinite</u> number of equations.

**Terminology**

The equations

$$R_{xx} \hat{\underline{h}} = R_{x\theta}$$

are called the <u>Wiener-Hopf</u> or <u>normal</u> equations.

The optimal $\hat{\underline{h}}$ is also called a <u>Wiener</u> estimator, after Norbert Wiener, especially in the context of estimating a signal in additive noise.

## Summary

- Classical and Bayesian linear estimation

  - only depends on first and second order moments
  - optimal in Gaussian case
  - proof: show that solution determined by 1st & 2nd order moments, assume Gaussianity, and invoke Gaussian results

- classical BLUE

  - solution of constrained quadratic minimization, doesn't always exist

- Bayesian LMMSE

  - solution of unconstrained quadratic minimization, always exists
  - projection in Hilbert space, obeys orthogonality principle

<u>Key</u>

a.
$$R_{\theta x} = E\left[(\underline{\theta} - E\underline{\theta})(\underline{x} - E\underline{x})^T\right]$$

$$= E\left[(\underline{\theta} - E\underline{\theta})(A\underline{\theta} + \underline{w} - AE\underline{\theta})^T\right]$$

$$= E\left[(\underline{\theta} - E\underline{\theta})\underline{w}^T\right] + E\left[(\underline{\theta} - E\underline{\theta})(\underline{\theta} - E\underline{\theta})^T A^T\right]$$

$$= 0 + R_{\theta\theta} A^T = R_{\theta\theta} A^T$$

$$R_{xx} = E\left[(\underline{X} - E\underline{X})(\underline{X} - E\underline{X})^T\right]$$

$$= E\left[(A\underline{\theta} + \underline{w} + AE\underline{\theta})(A\underline{\theta} + \underline{w} + AE\underline{\theta})^T\right]$$

$$= A R_{\theta\theta} A^T + R_{ww}$$

where we use the assumption that $\underline{\theta}, \underline{w}$ are uncorrelated, i.e. $E\left[(\underline{\theta} - E\underline{\theta})\underline{w}^T\right] = 0.$