



MMSE ESTIMATION

EE 541 – UNIT 3A

ESTIMATION, REGRESSION, CLASSIFICATION

statistical models

MMSE Estimation

Linear/Affine MMSE Est.

FIR Wiener filtering

Bayesian decision theory

Hard decisions

soft decisions (APP)

ML/MAP parameter
estimation

Karhunen-Loeve expansion

sufficient statistics

data driven

general regression

linear LS regression

stochastic gradient and

GD, SGD, LMS

Classification from data

linear classifier

logistical regression
(perceptron)

regularization

PCA

feature design

neural networks

for regression and
classification

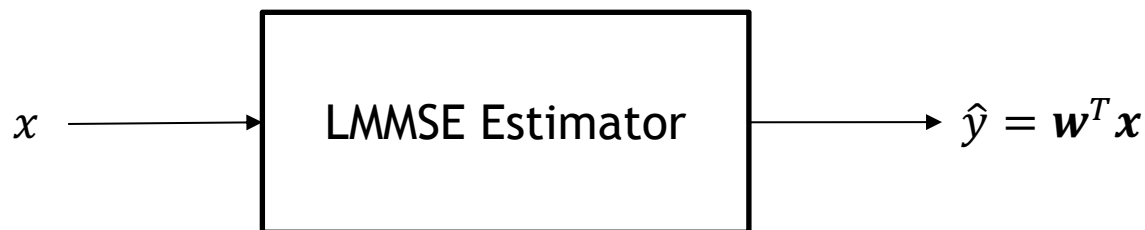
learning with SGD

working with data



PROBLEM: ESTIMATE $Y(t)$ FROM $X(t) = x$

- **Problem:** Given a vector observation $X(t) = x$, we would like to estimate $y(t)$ via linear filter $\hat{y} = \mathbf{w}^T \mathbf{x}$ with minimized mean squared error (MSE).



- The **objective** (cost function) to be minimized is $\text{MSE} = \mathbb{E}\{[y(t) - \mathbf{w}^T \mathbf{x}]^2\}$. The filter **design variables** are \mathbf{w} .



ESTIMATION

- What is Estimation?

- In **machine learning/signal processing/controls**, often need to make **predictions** based on **real world observations**.
- Process known as **inference** or **estimation**.

- What is LMMSE?

- **Estimates** are given as **linear combination** of observations!



KEY IDEAS FOR RANDOM VECTORS

- $N \times 1$ random vectors – generalization of 2×1
- Complete statistical description vs Second moment description
 - Directional preference (KL expansion)
- Gaussian processes and linear processing
- Linearity of the expectation operator

$$\mathbb{E}\{L(\mathbf{x}(t))\} = L(\mathbb{E}\{\mathbf{x}(t)\})$$

expectation commutes with any linear operation



RANDOM VECTORS

random vector

$$\mathbf{X}(t) = \begin{bmatrix} X_1(t) \\ X_2(t) \\ \vdots \\ X_N(t) \end{bmatrix} \quad (N \times 1)$$

Complete statistical
description

$$p_{\mathbf{X}(t)}(\mathbf{x}) = p_{X_1(t), X_2(t), \dots, X_N(t)}(x_1(t), x_2(t), \dots, x_N(t))$$

(pdf or cdf or pmf)

$$\mathbf{m}_X = \mathbb{E}\{\mathbf{X}(t)\} \quad \text{mean vector}$$

$$\mathbf{R}_X = \mathbb{E}\{\mathbf{X}(t)\mathbf{X}^T(t)\} \quad \text{correlation matrix}$$

Second Moment
Description

$$[\mathbf{R}_X]_{i,j} = \mathbb{E}\{X_i(t)X_j(t)\}$$

$$\begin{aligned} \mathbf{K}_X &= \mathbb{E}\{(\mathbf{X}(t) - \mathbf{m}_X)(\mathbf{X}(t) - \mathbf{m}_X)^T\} \quad \text{covariance matrix} \\ &= \mathbf{R}_X - \mathbf{m}_X \mathbf{m}_X^T \end{aligned}$$

$$[\mathbf{K}_X]_{i,j} = \text{Cov}[X_i(t), X_j(t)]$$



KARHUNEN–LOÈVE (KL) EXPANSION

Can always find orthonormal set of e-vectors of K

These are an alternate coordinate systems (**rotations, reflections**)
in this eigen-coordinate system, the components are uncorrelated

“principal components”

The eigen-values are the variance (energy) in each principal directions

(reduce dimensions by "throwing out" low-energy components)



KL-EXPANSION

$$\mathbf{K}_X \mathbf{e}_k = \lambda_k \mathbf{e}_k \quad k = 1, 2, \dots, N$$

Eigen equation

$$\mathbf{e}_k^T \mathbf{e}_l = \delta[k - l], \quad \lambda_k \geq 0$$

orthonormal eigen vectors

$$\mathbf{X}(t) = \sum_{k=1}^N X_k(t) \mathbf{e}_k$$

change of coordinates

$$X_k(t) = \mathbf{e}_k^T \mathbf{X}(t)$$

$$\mathbb{E}\{X_k(t)X_l(t)\} = \mathbf{e}_k^T \mathbf{K}_X \mathbf{e}_l = \lambda_k \delta[k - l]$$

uncorrelated components

$$\mathbf{K}_X = \sum_{k=1}^N \lambda_k \mathbf{e}_k \mathbf{e}_k^T = \mathbf{E} \mathbf{\Lambda} \mathbf{E}^T$$

Mercer's Theorem

$$\mathbb{E}\{\|\mathbf{X}(t)\|^2\} = \text{tr}(\mathbf{K}_x) = \sum_{k=1}^N \lambda_k$$

Total Energy

Always exists because \mathbf{K}_X is symmetric and positive semi- definite (PSD)

KL-EXPANSION EXAMPLE

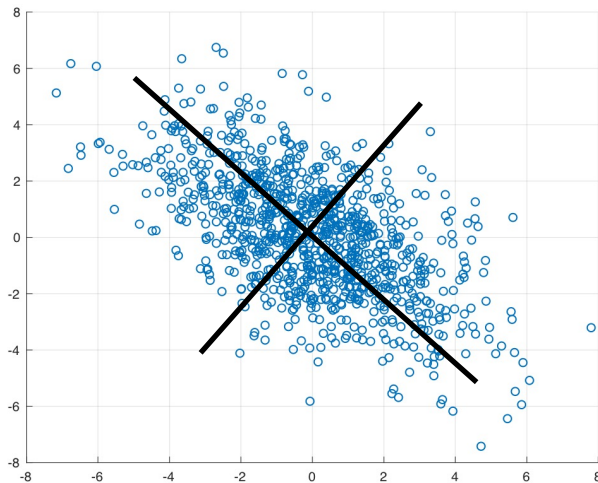
$$\mathbf{X}(t) = \mathbf{H} \mathbf{W}(t)$$

$$\mathbf{H} = \begin{bmatrix} 1 & 2 \\ 1 & -2 \end{bmatrix}$$

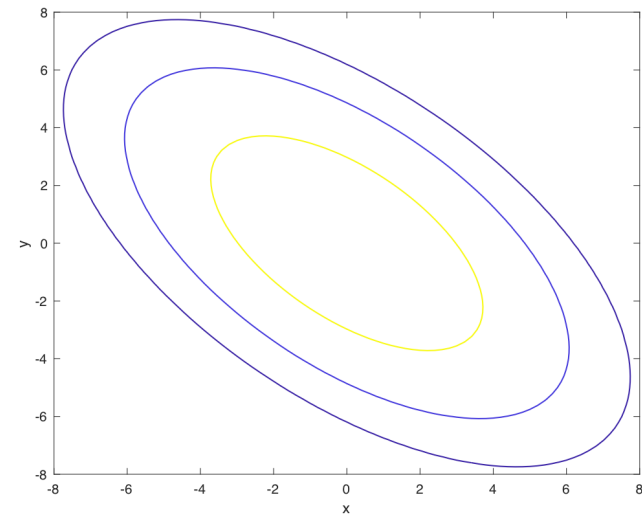
$$\mathbf{K} = \mathbf{H} \mathbf{K}_w \mathbf{H}^T = \mathbf{H} \mathbf{H}^T = \begin{bmatrix} 5 & -3 \\ -3 & 5 \end{bmatrix}$$

$$\mathbf{E} = \frac{1}{\sqrt{2}} \begin{bmatrix} +1 & +1 \\ +1 & -1 \end{bmatrix}$$

$$\mathbf{\Lambda} = \begin{bmatrix} 2 & 0 \\ 0 & 8 \end{bmatrix}$$



generated with $\mathbf{W}(t)$ Gaussian



Gaussian pdf contours



LINEAR/AFFINE MMSE (LMMSE AND AMMSE)

estimate $\mathbf{Y}(t)$ from $\mathbf{X}(t) = x$

$$\min_{f(x)} E \left\{ \|\mathbf{y}(t) - f(\mathbf{x}(t))\|^2 \right\}$$

with no constraint, this is the conditional expectation function

$$\begin{aligned} f_{\text{opt}}(x) &= m_{\mathbf{y}|\mathbf{X}}(x) \\ &= \mathbb{E}\{\mathbf{y}(t) | \mathbf{X}(t) = x(t)\} \\ &= \int \mathbf{y} p_{\mathbf{y}|\mathbf{X}}(\mathbf{y}|x) d\mathbf{y} \end{aligned}$$

LMMSE ASSUMPTIONS

Assume that we know the second order statistics of the joint distribution $p_{X(t),Y(t)}(\mathbf{x}, y)$

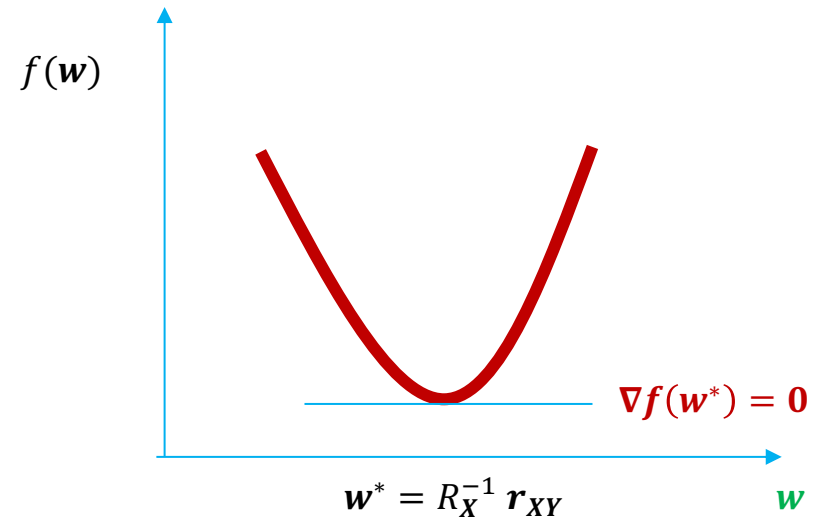
i.e., $\mathbf{m}_X, m_Y, R_X, r_Y, \mathbf{r}_{XY}$.

The objective is a quadratic function of \mathbf{w} ,

$$\begin{aligned} f(\mathbf{w}) &= \mathbb{E}\{[Y(t) - \mathbf{w}^T \mathbf{x}]^2\} \\ &= r_Y - 2\mathbf{r}_{YX}^T \mathbf{w} + \mathbf{w}^T R_X \mathbf{w} \end{aligned}$$

The global optimal occurs at \mathbf{w}^* .

i.e., $\nabla f(\mathbf{w})|_{\mathbf{w}=\mathbf{w}^*} = \mathbf{0}$





MINIMUM MEAN-SQUARE ERROR ESTIMATION (MMSE)

Estimate $y(u)$ from $x(u) = x$

cross covariance matrix

$$\mathbf{K}_{YX} = \mathbb{E}\{(\mathbf{y}(t) - \mathbf{m}_Y)(\mathbf{x}(t) - \mathbf{m}_X)^T\} = [\mathbf{K}_{YX}]^T$$

Affine MMSE

$$\min_{\mathbf{f}(\mathbf{x})=\mathbf{F}\mathbf{x}+\mathbf{b}} \mathbb{E}\{\|\mathbf{y}(t) - \mathbf{f}(\mathbf{x}(t))\|^2\}$$

$$\mathbf{F}_{\text{AMMSE}} = \mathbf{K}_{YX}\mathbf{K}_X^{-1}$$

$$\mathbf{b}_{\text{AMMSE}} = \mathbf{m}_Y - \mathbf{F}_{\text{AMMSE}}\mathbf{m}_X$$

$$\min_{\mathbf{F}, \mathbf{b}} \mathbb{E}\{\|\mathbf{y}(t) - [\mathbf{F}\mathbf{x}(t) + \mathbf{b}]\|^2\}$$

$$\hat{\mathbf{y}} = \mathbf{K}_{YX}\mathbf{K}_X^{-1}(\mathbf{x} - \mathbf{m}_X) + \mathbf{m}_Y$$

$$\text{AMMSE} = \text{Tr}(\mathbf{K}_Y - \mathbf{K}_{YX}\mathbf{K}_X^{-1}\mathbf{K}_{XY})$$

Linear MMSE

$$\min_{\mathbf{f}(\mathbf{x})=\mathbf{F}\mathbf{x}} \mathbb{E}\{\|\mathbf{y}(t) - \mathbf{f}(\mathbf{x}(t))\|^2\}$$

$$\mathbf{F}_{\text{AMMSE}} = \mathbf{R}_{YX}\mathbf{R}_X^{-1}$$

$$\min_{\mathbf{F}} \mathbb{E}\{\|\mathbf{y}(t) - \mathbf{F}\mathbf{x}(t)\|^2\}$$

$$\hat{\mathbf{y}} = \mathbf{R}_{YX}\mathbf{R}_X^{-1}\mathbf{x}$$

$$\text{LMMSSE} = \text{Tr}(\mathbf{R}_Y - \mathbf{R}_{YX}\mathbf{R}_X^{-1}\mathbf{R}_{XY})$$

- affine often called linear... same when means are zero
- Conditional Expectation better than affine, affine better than Linear

PROOF FOR LMMSE

$$\min_F \mathbb{E}\{\|y(t) - FX(t)\|^2\}$$

Proof:

$$\text{MSE}(F) = \mathbb{E}\{\|y(t) - FX(t)\|^2\}$$

$$= \mathbb{E}\left\{\left\|\left(y(t) - F_{\text{opt}}x(t)\right) + (F_{\text{opt}} - F)x(t)\right\|^2\right\}$$

$$= \mathbb{E}\{\|y(t) - F_{\text{opt}}x(t)\|^2\} + \text{Tr}\left((F_{\text{opt}} - F)R_X(F_{\text{opt}} - F)^T\right)$$

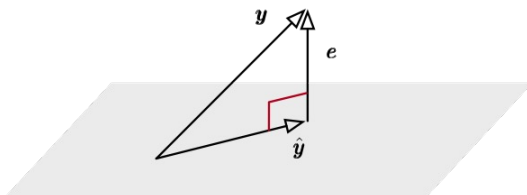
$$+ 2 \text{Tr}\left((R_{yx} - F_{\text{opt}}R_X)(F_{\text{opt}} - F)^T\right)$$

$$v^T w = \text{Tr}(wv^T)$$

if: $F_{\text{opt}}R_X = R_{XY}$ then: $\text{MSE}(F) = \mathbb{E}\{\|y(t) - F_{\text{opt}}x(t)\|^2\} + \underbrace{\text{Tr}\left((F_{\text{opt}} - F)R_X(F_{\text{opt}} - F)^T\right)}_{\geq 0 \forall F, \text{ since } R_X \text{ is psd}}$

↑
Wiener-Hopf equations, aka Orthogonality Principle

because of orthogonality principle (error and signal uncorrelated)

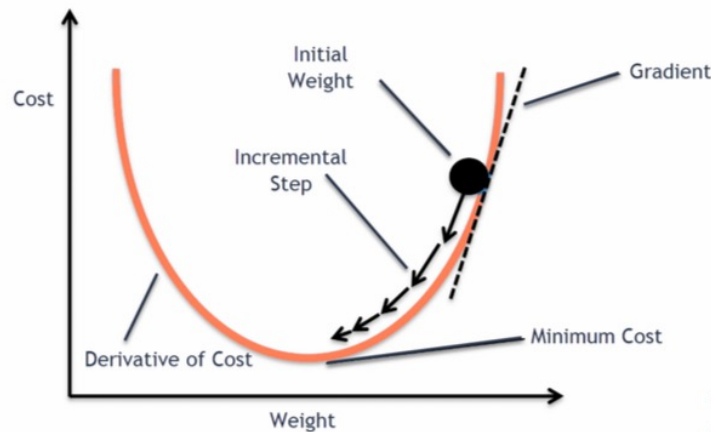


space of all estimates/approximations

$$\begin{aligned}\mathbb{E}\{\|y(t) - \hat{y}(t)\|^2\} &= \mathbb{E}\{\|y(t)\|^2\} + \mathbb{E}\{\|\hat{y}(t)\|^2\} \\ &= \text{Tr}(R_y - R_{yx}R_X^{-1}R_{xy})\end{aligned}$$

GRADIENT DESCENT

- Instead of obtaining optimal w^* directly, we can also find it iteratively via
 1. Initialize w_0
 2. $w_{n+1} = w_n - \eta \nabla f(w_n)$



#MLmuse
CLAIRVOYANT

source: <https://blog.clairvoyantsoft.com/>

GAUSSIAN RANDOM VECTORS

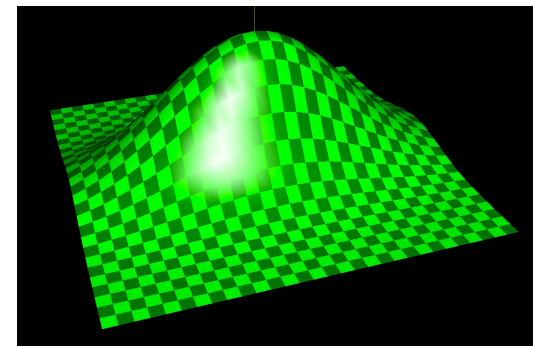
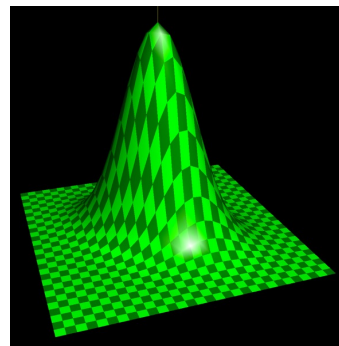
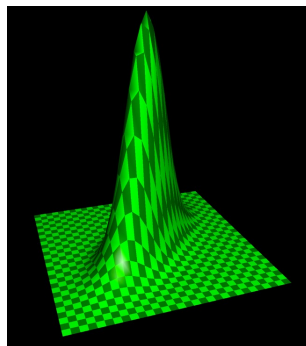
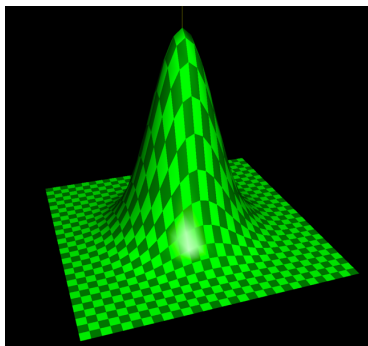
$$p_{x(t)}(x) = \mathcal{N}_N(x; m_x, K_x)$$

$$= \frac{1}{(2\pi)^{N/2} \sqrt{|K_x|}} \exp\left(-\frac{1}{2}(x - m_x)K_x^{-1}(x - m_x)\right)$$

$$p_{x(t)}(x) = \mathcal{N}_N(x; 0, I)$$

$$= \frac{1}{(2\pi)^{N/2}} \exp\left(-\frac{1}{2}\|x - m_x\|^2\right)$$

$$= \prod_{k=1}^N \mathcal{N}_1(x_k; 0, 1)$$



$$K = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

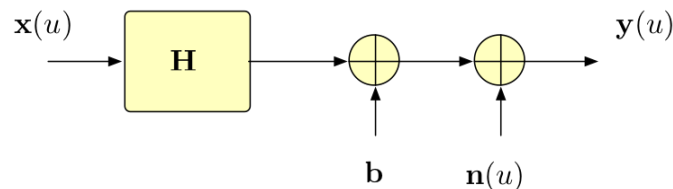
$$K = \begin{bmatrix} 1 & 0.7 \\ 0.7 & 1 \end{bmatrix}$$

$$K = \begin{bmatrix} 1 & -0.4 \\ -0.4 & 1 \end{bmatrix}$$

$$K = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$

$$K = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$$

GAUSSIAN RANDOM VECTORS



any linear processing of Gaussians yields Gaussians

if $\begin{bmatrix} x(t) \\ n(t) \end{bmatrix}$ is Gaussian (i.e., $x(u)$ and $n(u)$ jointly-Gaussian), then:

$y(t)$ is Gaussian

$\begin{bmatrix} x(t) \\ y(t) \end{bmatrix}$ is Gaussian

$\begin{bmatrix} x(t) \\ n(t) \\ y(t) \end{bmatrix}$ is Gaussian

Jointly-Gaussian
common in EE!

any subset of these random variables is also Gaussian

MMSE ESTIMATION: SPECIAL CASE JOINTLY-GAUSSIAN

Estimate $y(u)$ from $x(u) = x$

Conditional Gaussian pdf:

$$p_{y(t)|x(t)}(y|x) = \frac{p_{X(t),Y(t)}(x, y)}{p_{X(t)}(x)}$$

$$= \frac{\mathcal{N}_{M+N} \left(\begin{bmatrix} x \\ y \end{bmatrix}; \begin{bmatrix} m_X \\ m_Y \end{bmatrix}, \begin{bmatrix} K_Y & K_{XY} \\ K_{YX} & K_Y \end{bmatrix} \right)}{\mathcal{N}_N(x; m_X, K_X)}$$

$$= \mathcal{N}_M \left(y; \underbrace{m_Y + K_{YX}K_X^{-1}(x - m_X)}_{\text{Affine MMSE estimator}}, \underbrace{K_Y - K_{YX}K_X^{-1}K_{XY}}_{\text{error covariance}} \right)$$

JG... no ML or deep learning needed!

For JG observation and target: $\mathbb{E}[Y|x]$ is the Affine MMSE estimator



REMARKS

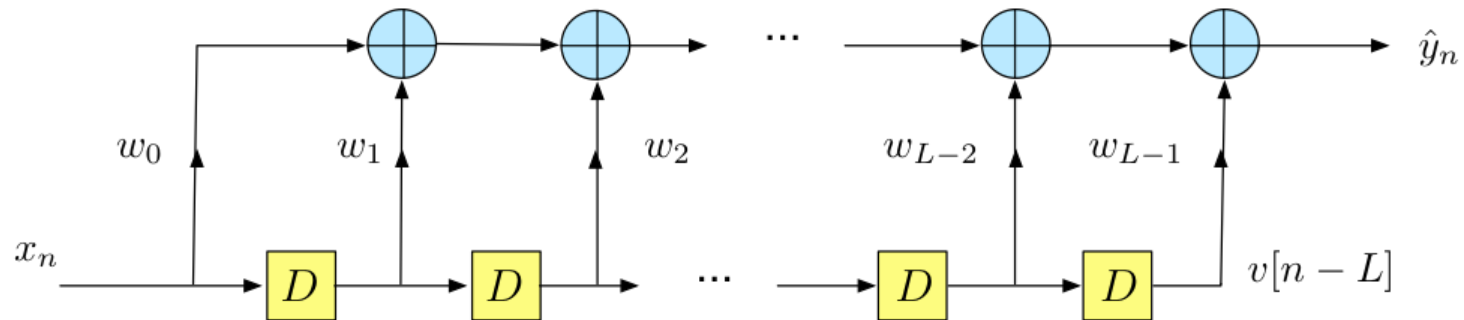
1. **Closed-form global optimality** can be derived if we have a convex and differentiable cost function.
2. **Gradient descent** works in any differentiable cost function.
3. What if we don't have second order statistics but lots of samples?
Use LMS!

Single point gradient descent or small batch gradient descent!

R_x computation not tractable for large number of samples

LMS ALGORITHM

- Suppose time index n and you have samples $\{(x_n, y_n)\}$ instead of the second order statistics then we use an online learning algorithm called least mean square **adaptive filtering**
 - Introduced by Widrow and Hopf (1959).



$$y_n = \sum_{l=1}^N w_l x_{n-l} = w^T v_n \quad w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_L \end{bmatrix} \quad v_n = x_{n-(L-1)} = \begin{bmatrix} x_n \\ x_{n-1} \\ \vdots \\ x_{n-L+1} \end{bmatrix}$$



LMS ALGORITHM

- Assumptions:
 - Observe $y(n) = s(n) + \text{Noise}$ and no access to the data signal $s(n)$.
 - Have a reference noise signal $x(n)$ with strong correlation to actual noise signal
- Why we need this?:
 - Slowly drifting interfering sinusoid so notch filter is insufficient
- Benefits:
 - Behaves as an adaptive notch filter
 - The notch can be very sharp, depending upon step size

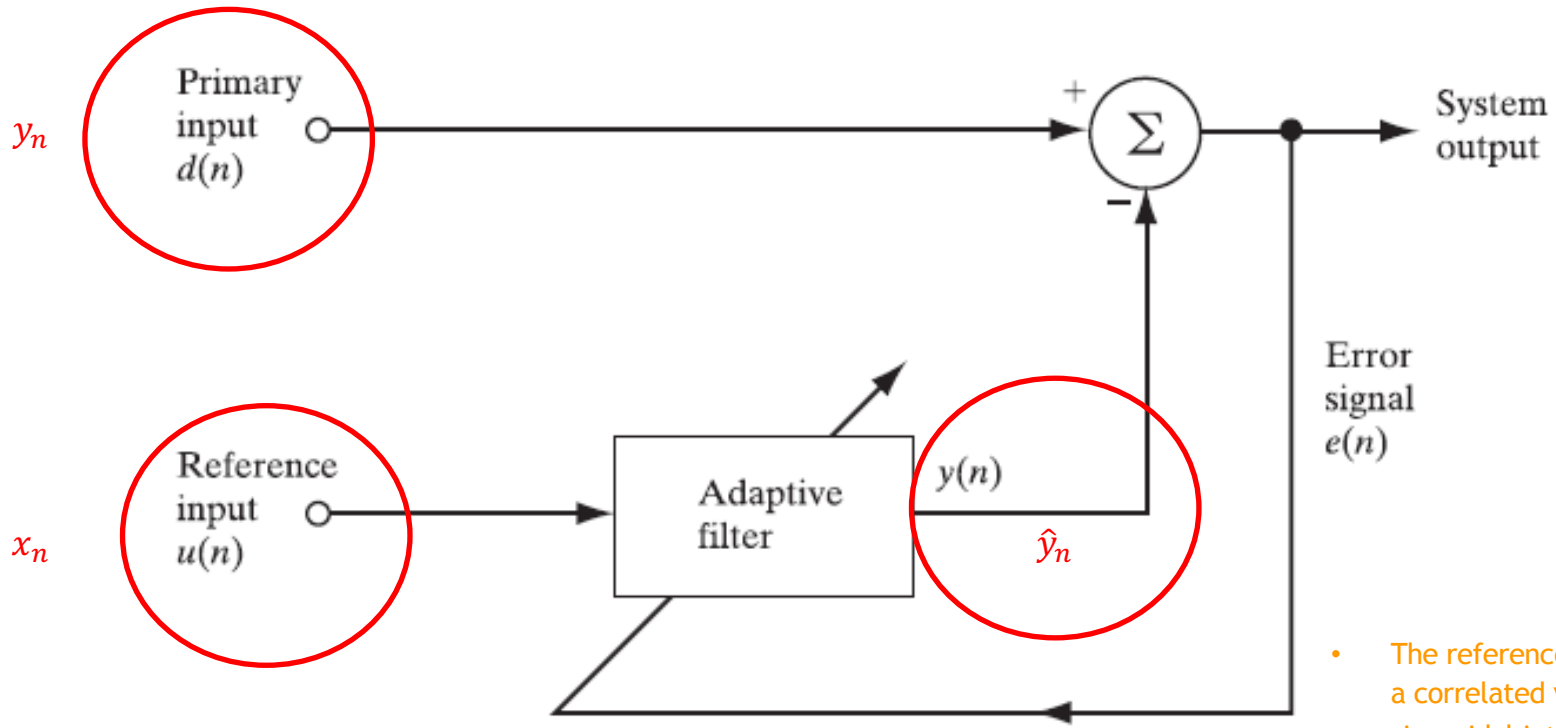
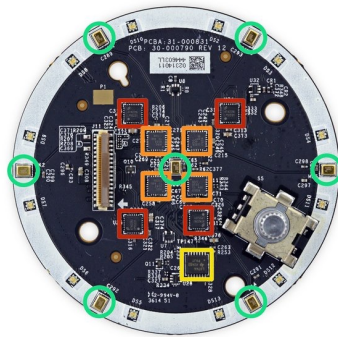


FIGURE 6.6 Block diagram of adaptive noise canceller.

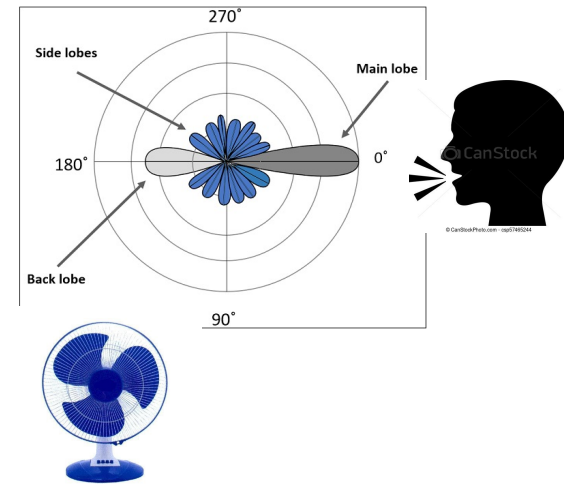
- The reference input supplies a correlated version of the sinusoidal interference.
- For the adaptive filter, we may use an FIR filter whose tap weights are adapted by means of the LMS algorithm

LMS HISTORY/EXAMPLE

Widrow and Hopf, Adaptive Linear Element (ADALINE)
(developed LMS for adaptive antenna array processing)



[ifixit.com](https://www.ifixit.com)



point a beam at the desired speaker and *learn*
to cancel noise energy in other directions using LMS



STEEPEST DESCENT AND LMS

Estimate $y(\mathbf{u})$ from $\mathbf{x}(\mathbf{u}) = \mathbf{x}$

$$E(w) = \mathbb{E}\{[y(t) - w^T x(t)]^2\}$$

$$\begin{aligned}\nabla_w E &= 2\mathbb{E}\{wx(t)x^T(t) - y(t)x(t)\} \\ &= 2(w^T R_x - r_{xy})\end{aligned}$$

Steepest descent using (ensemble average) gradient:

$$\begin{aligned}\hat{w}_{n+1} &= \hat{w}_n - \left(\frac{\eta}{2}\right) \nabla_w E \\ &= \hat{w}_n + \eta(r_{xy} - R_x \hat{w}_n)\end{aligned}$$

Single Point Stochastic Gradient Descent:

$$\begin{aligned}-\frac{1}{2}\nabla_w E &= r_{xy} - R_x w \\ &= \mathbb{E}\{y(t)X(t) - X(t)X^T(t)w\} \\ &\approx y_n x_n - x_n x_n^T w \\ &= (y_n - x_n^T \hat{w}_n)x_n\end{aligned}$$

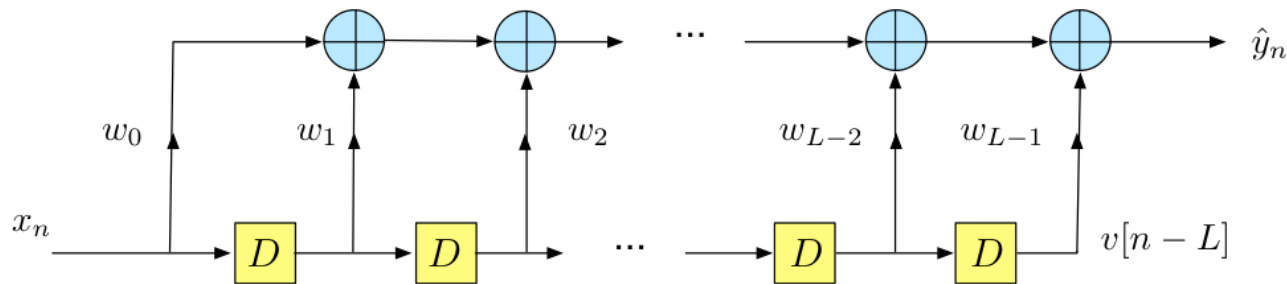
$$\hat{w}_{n+1} = \hat{w}_n + \eta(y_n - \hat{w}_n^T x_n)x_n$$

this is called “on-line learning”

when $n \sim$ time, this is the **Adaptive Least Mean Square (LMS)** filter

when n does not represent time can average the gradient over more data points to improve approximation (batching)

LMS ALGORITHM IS ADAPTIVE FIR FILTER



$$y_n = \sum_{l=1}^N w_l x_{n-l} = w^T v_n$$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_N \end{bmatrix}$$

$$v_n = x_{n-(L-1)} = \begin{bmatrix} x_n \\ x_{n-1} \\ \vdots \\ x_{n-L+1} \end{bmatrix}$$

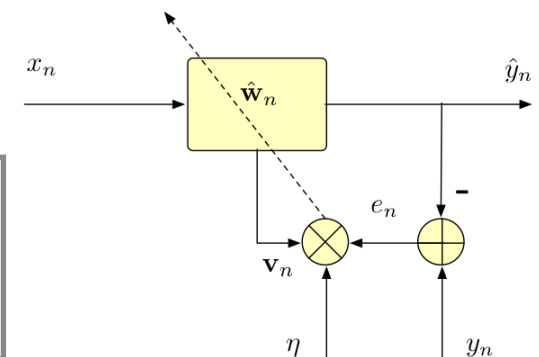
Single Point Stochastic Gradient Descent:

$$\begin{aligned} -\frac{1}{2} \nabla_w E &= r_{v_n y_n} - R_{v_n} w \\ &= \mathbb{E}\{y(t)v_n(t) - v_n(t)v_n^T(t)w\} \\ &\approx (y_n - w_n^T v_n)v_n \end{aligned}$$

LMS Algorithm

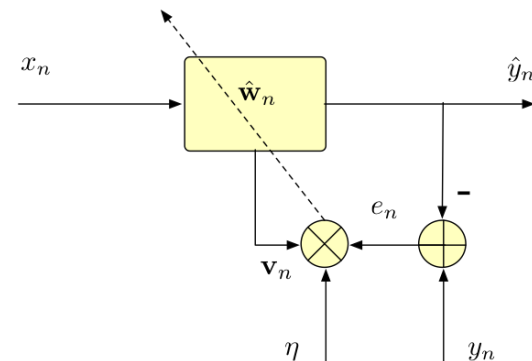
$$\begin{aligned} \hat{w}_{n+1} &= \hat{w}_n + \eta(y_n - w_n^T v_n)v_n \\ &= \hat{w}_n + \eta(y_n - \hat{y}_n)v_n \\ &= \hat{w}_n + \eta e_n v_n \end{aligned}$$

(an online linear regressor)



LMS ALGORITHM AS ADAPTIVE FIR FILTER

LMS algorithm: $\hat{w}_{n+1} = \hat{w}_n + \eta(y_n - \hat{w}_n^T v_n)v_n$



If R_{v_n} and $r_{v_n y} = \mathbb{E}\{v_n(t)y(t)\}$ do not change with n ,

$$\hat{w}_n \rightarrow \approx w_{LMMSE} = R_v^{-1} r_{vy}$$

If these correlations vary with time, the LMS filter will **adaptively track** them



LMS ALGORITHM: FOR SINUSOIDAL NOISE CANCELLATION

Primary input:

$$y(n) = \boxed{s(n)} + \boxed{A_0 \cos(\omega_0 n + \phi_0)} \quad (1)$$

Information bearing part Sinusoidal noise

Reference input:

$$x(n) = \boxed{A \cos(\omega_0 n + \phi)} \quad (2)$$

Same nature as Sinusoidal noise

$s(n)$ = information bearing signal / data signal.

LMS: get current estimate and error:

$$\begin{aligned} \hat{y}(n) &= \sum_{i=0}^{L-1} \hat{w}_i(n) x(n-i) \\ &= \hat{\mathbf{w}}_n^T \mathbf{v}_n \end{aligned} \quad (3)$$

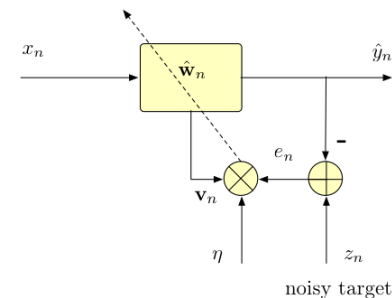
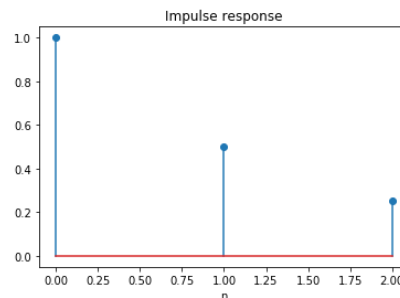
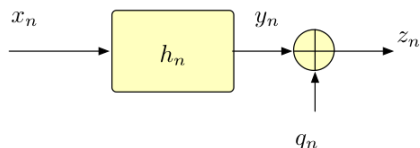
$$e(n) = y(n) - \hat{y}(n) = y(n) - \hat{\mathbf{w}}_n^T \mathbf{v}_n \quad (4)$$

Update filter tap weights:

$$\begin{aligned} \hat{\mathbf{w}}_{n+1} &= \hat{\mathbf{w}}_n + \eta e(n) \mathbf{v}_n \\ &= \hat{\mathbf{w}}_n + \eta (y(n) - \hat{\mathbf{w}}_n^T \mathbf{v}_n) \mathbf{v}_n \end{aligned} \quad (5)$$

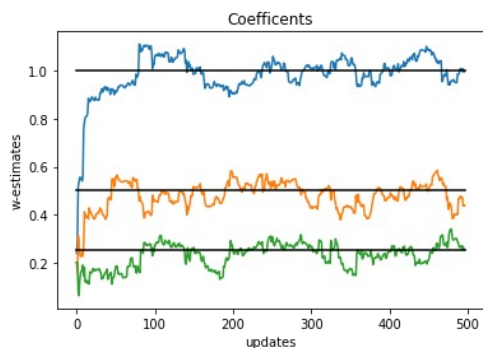
EXPERIMENT: LMS EXAMPLE

Generate data:

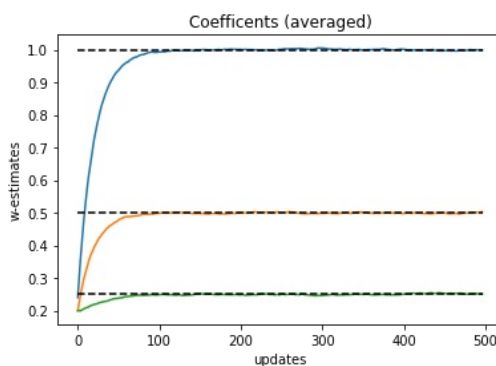


this is the ideal case where model
and observed data are matched

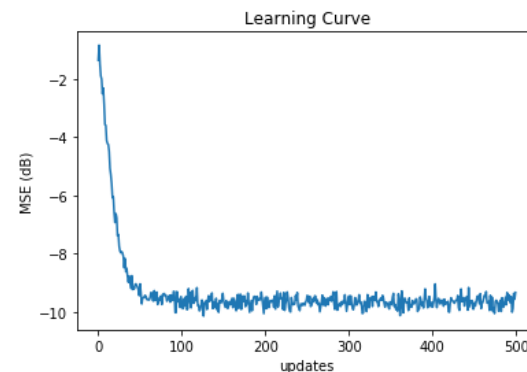
$$\eta = 0.05, \text{SNR} = 10 \text{ dB}$$



single run



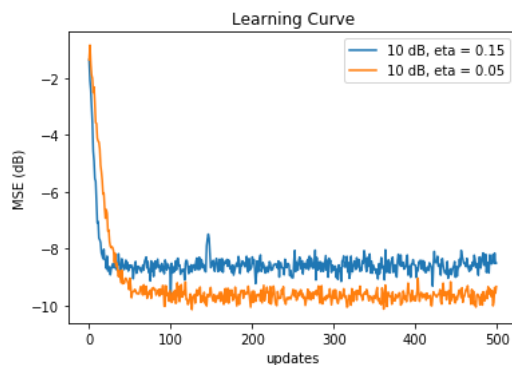
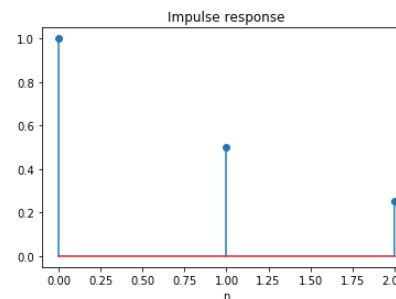
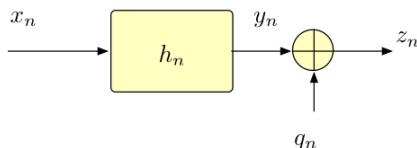
averaged over 500 runs



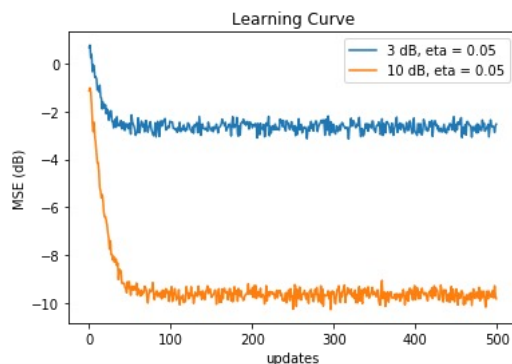
averaged over 500 runs

EXPERIMENT: LMS EXAMPLE

Generate data:



larger learning rate means faster convergence
but more misalignment (gradient noise)



even the optimal Wiener (LMMSE) filter will
have higher MMSE when the SNR is lower

EVOLUTION OF THE LMS ALGORITHM

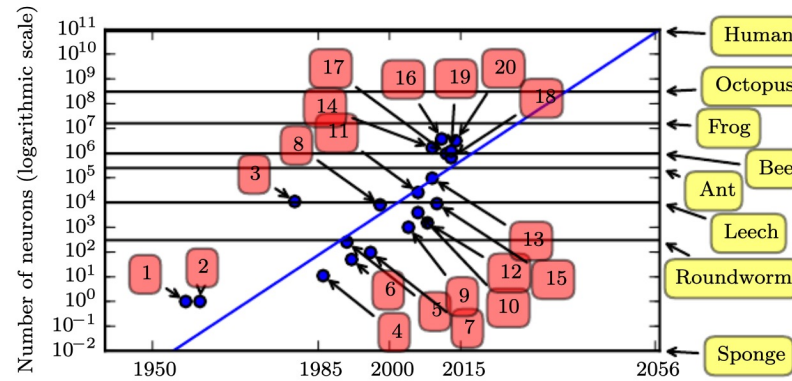


Figure 1.11: Increasing neural network size over time. Since the introduction of hidden units, artificial neural networks have doubled in size roughly every 2.4 years. Biological neural network sizes from [Wikipedia \(2015\)](#).

1. Perceptron ([Rosenblatt, 1958, 1962](#))
2. Adaptive linear element ([Widrow and Hoff, 1960](#))
3. Neocognitron ([Fukushima, 1980](#))
4. Early back-propagation network ([Rumelhart et al., 1986b](#))
5. Recurrent neural network for speech recognition ([Robinson and Fallside, 1991](#))
6. Multilayer perceptron for speech recognition ([Bengio et al., 1991](#))
7. Mean field sigmoid belief network ([Saul et al., 1996](#))
8. LeNet-5 ([LeCun et al., 1998b](#))
9. Echo state network ([Jaeger and Haas, 2004](#))
10. Deep belief network ([Hinton et al., 2006](#))
11. GPU-accelerated convolutional network ([Chellapilla et al., 2006](#))
12. Deep Boltzmann machine ([Salakhutdinov and Hinton, 2009a](#))
13. GPU-accelerated deep belief network ([Raina et al., 2009](#))
14. Unsupervised convolutional network ([Jarrett et al., 2009](#))
15. GPU-accelerated multilayer perceptron ([Ciresan et al., 2010](#))
16. OMP-1 network ([Coates and Ng, 2011](#))
17. Distributed autoencoder ([Le et al., 2012](#))
18. Multi-GPU convolutional network ([Krizhevsky et al., 2012](#))
19. COTS HPC unsupervised convolutional network ([Coates et al., 2013](#))
20. GoogLeNet ([Szegedy et al., 2014a](#))



MMSE SUMMARY

1. Estimation using statistical models
2. Best MMSE estimator (unconstrained) is conditional expectation
 - Requires complete statistical description of observed and desired – i.e., $p(\mathbf{y}|\mathbf{x})$
3. Linear/affine MMSE estimator have closed form equations
 - Require only the second moment description of observed and desired – i.e., means, correlations
4. For jointly Gaussian observed and desired - 2 & 3 are the same!
5. The LMS algorithm is an algorithm that approximating the LMMSE cost function gradient with a single realization.