

Εκτίμηση Μίσου Τετραγωνικού Σφάλματος (Mean Square Error Estimation)

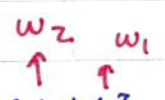
Ας επιμερευθούμε και πάλι σε ένα πρόβλημα στατιστικής.
 Στο προηγούμενο κεφάλαιο είδαμε πως η έξοδος του πεπτηρον αυήκε στο σύνολο $\{-1, +1\}$ ανάλογα με την κλάση πρόβλεψης του διαυήματος χαρακτηριστικών \underline{x} . Μάλιστα, υπό την προϋπόθεση της γραμμικής διαχωρισιμότητας των δύο κατηγοριών, η έξοδος του πεπτηρον ταυίζεται με την σωστή κατηγορία για κάθε διάυσμα χαρακτηριστικών του συνόλου εκπαίδευσης αφού έχει επιλθεί η σύχυλση του αλγορίθμου.

Στο ενόγω αυή, θα εστιάσουμε στην σχεδιασμέ ενή γραμμική ταξινόμηση, η επιθυμητή έξοδος του οποίου θα είναι κλι πάλι στατιστική $\{-1, +1\}$ εξαρτήμενη από την κλάση πρόβλεψης του διαυήματος εισόδου. Ωστόσο, στην συχευρήνη περίπτωση θα επιτρέψουμε την ύπαρξη ταξινόμητων σφαλμάτων. Περικώς, δηλαδή, όπου οι έξοδος του ταξινόμηση θα είναι διαφορετική από την επιθυμητή έξοδο.

Αν επιθυμητή έξοδος του ταξινόμηση συμβολισθεί με $\gamma(\underline{x}) = \gamma \in \{-1, +1\}$ όπου \underline{x} είναι το διάυσμα των χαρακτηριστικών που προφοδύζεται ως είσοδος στο ταξινόμηση, η έξοδος αυτού θα δίνεται από την σχέση:

$$\gamma(\underline{x}) = \underline{w}^T \cdot \underline{x} \quad [1]$$

⊛ Η ύπαρξη σταθερού όρου στην εξίσωση (1) αντιστοιχεί στην θεώρηση επιτετομήτων διαυημάτων βάρων και διαυημάτων χαρακτηριστικών.



Το διάνυσμα των βαρών θα υπολογισθεί έτσι ώστε να ελαχιστοποιηθεί το μέσο τετραγωνικό σφάλμα ανάμεσα στην επιθυμητή και την πραγματική έξοδο του ταξινόμητή. Επομένως, θα πρέπει να ορίσουμε το συνάρτησης κόστους ως εξής:

$$J(\underline{w}) = E [|y - \underline{x}^T \underline{w}|^2] \quad [2]$$

Με βάση του παραπάνω ορισμό του συνάρτησης κόστους, η βέλτιστη τιμή του διανύσματος των βαρών προκύπτει ως:

$$\underline{w}^* = \arg \min_{\underline{w}} J(\underline{w}) \quad [3]$$

Εύκολα μπορεί να συμπεράνει κανείς πως: $J(\underline{w}) = E [(y - \underline{w}^T \underline{x})^2] \Leftrightarrow$

$$J(\underline{w}) = \int_{-\infty}^{+\infty} p(\underline{x}) \cdot (y(\underline{x}) - \underline{w}^T \underline{x})^2 d\underline{x} \quad (4). \quad \text{Οστόσο, ισχύει ότι:}$$

$$(i): \quad p(\underline{x}) = \begin{cases} P(\omega_1), & \underline{x} \in \omega_1; \\ P(\omega_2), & \underline{x} \in \omega_2. \end{cases} \quad \text{και} \quad (ii): \quad y(\underline{x}) = \begin{cases} +1, & \underline{x} \in \omega_1; \\ -1, & \underline{x} \in \omega_2. \end{cases}$$

που οδηγεί στην παρακάτω σχέση:

$$J(\underline{w}) = P(\omega_1) \int_{\underline{x} \in \omega_1} (1 - \underline{w}^T \underline{x})^2 d\underline{x} + P(\omega_2) \int_{\underline{x} \in \omega_2} (-1 - \underline{w}^T \underline{x})^2 d\underline{x} \quad [5]$$

Η ελαχιστοποίηση της συνάρτησης κόστους που εμφανίζεται στην εξίσωση (1) πραγματοποιείται κατά τα γνωστά αίτια της συνθήκης πρώτης τάξης, ως εξής:

$$\frac{\partial J(\underline{w})}{\partial \underline{w}} = \underline{0} \quad [6]$$

Η περτελίρω σφάλμα της εδίσωμς (6) δίνη:

$$\frac{\partial J(\underline{w})}{\partial \underline{w}} = \underline{0} \Leftrightarrow \frac{\partial}{\partial \underline{w}} E[(y - \underline{x}^T \underline{w})^2] = \underline{0} \Leftrightarrow E\left[\frac{\partial}{\partial \underline{w}} (y - \underline{x}^T \underline{w})^2\right] = \underline{0} \Leftrightarrow$$

$$\left[\left[2(y - \underline{x}^T \underline{w}) \frac{\partial}{\partial \underline{w}} (y - \underline{x}^T \underline{w}) \right] \right] = \underline{0} \Leftrightarrow -2 E[\underline{x}(y - \underline{x}^T \underline{w})] = \underline{0} \Leftrightarrow$$

$$E[y \cdot \underline{x} - \underline{x} \underline{x}^T \underline{w}] = \underline{0} \Leftrightarrow E[y \cdot \underline{x}] - E[\underline{x} \underline{x}^T \underline{w}] = \underline{0} \Leftrightarrow$$

$$E[y \cdot \underline{x}] - E[\underline{x} \underline{x}^T] \underline{w} = \underline{0} \Leftrightarrow E[\underline{x} \underline{x}^T] \cdot \underline{w} = E[y \cdot \underline{x}] \Leftrightarrow$$

$$\underline{w}^* = \underline{R}_x^{-1} \cdot E[y \cdot \underline{x}] \quad (7), \text{ όπου } \underline{R}_x = E[\underline{x} \underline{x}^T] \quad (8)$$

$$\psi \in \underline{R}_x \in \mathbb{R}^{p \times p}$$

(*) Για του πίνακα \underline{R}_x έχουμε ότι:

$$\underline{R}_x \equiv E \left[\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix} \begin{bmatrix} x_1 & x_2 & \dots & x_p \end{bmatrix} \right] = E \begin{bmatrix} x_1 x_1 & x_1 x_2 & \dots & x_1 x_p \\ x_2 x_1 & x_2 x_2 & \dots & x_2 x_p \\ \vdots & \vdots & \ddots & \vdots \\ x_p x_1 & x_p x_2 & \dots & x_p x_p \end{bmatrix} \quad \Leftrightarrow$$

$$\underline{R}_x = \begin{bmatrix} E[x_1 x_1] & \dots & E[x_1 x_p] \\ \vdots & \ddots & \vdots \\ E[x_p x_1] & \dots & E[x_p x_p] \end{bmatrix} \quad (9)$$

(*) Ο πίνακς \underline{R}_x ονομάζεται πίνακς συσχίζιμς ή αυτοσυσχίζιμς [correlation or autocorrelation matrix] και τανίζτμ με του πίνακα συβδσκήμης συν πείζωμς κατὰ τήν οποία το μίσο δίαυνομα είνμ μηδενικό $\underline{\mu} = \underline{0}$.

(*) Αντίστοιχα, το δίαυνομα $E[y \cdot \underline{x}] = E \left[\begin{bmatrix} y x_1 \\ \vdots \\ y x_p \end{bmatrix} \right]$ ποσοτικοποιήμν συσχίζιμς (cross-correlation) (εξέρο συσχίζιμς) ανάμετα συν επιθυμής ίζοδο και σεα δίαυόσματα χαρακτηριστικόν που ζαζινομήμς λαμβάνμ συν είσοδο.

Κατά συνέπεια, το βέλτιστο διάνυσμα βάρων στο οποίο οδηγεί η ελαχιστοποίηση της συνάρτησης κόστους που χρησιμοποιείται από τον μέθοδο της εκτίμησης του μέσου τετραγωνικού σφάλματος αντιστοιχεί στην επίλυση ενός συνόλου γραμμικών εξισώσεων υπό την προϋπόθεση ότι ο πίνακας αυτοσυσχίτισης R_x είναι αντιστρέψιμος.

Είναι αξιοσημείωτο το γεγονός πως η βέλτιστη λύση για το διάνυσμα των βάρων που συνδέεται με την εκτίμηση του μέσου τετραγωνικού σφάλματος επιδράται γεωμετρικώς ερμηνείας. Συγκεκριμένα, γνωρίζουμε πως οι τυχαίες μεταβλητές μπορούν να θεωρηθούν ως σημεία σε ένα πολυδιάστατο διανυσματικό χώρο. Μάλιστα, γίνεται εύκολα αντιληπτό πως ο τελεστής της αναμενόμενης τιμής $E[\cdot]$ μετατρέπει δύο τυχαίων μεταβλητών ικανοποιεί τις ιδιότητες του εσωτερικού γινομένου. Συγκεκριμένα, έχουμε ότι:

$$E[x \cdot y] = \text{Cov}(x, y) + E[x]E[y] \quad (*)$$

$$\varphi(x, y) = E[x \cdot y] \quad (**)$$

INNER
PRODUCT

$$(*) \text{Cov}(x, y) = E[(x - E[x])(y - E[y])] \Leftrightarrow$$

$$\text{Cov}(x, y) = E[xy - E[y] \cdot x - E[x] \cdot y + E[x] \cdot E[y]] \Leftrightarrow$$

$$\text{Cov}(x, y) = E[xy] - E[y]E[x] - E[x]E[y] + E[E[x] \cdot E[y]] \Leftrightarrow$$

$$\text{Cov}(x, y) = E[xy] - E[x]E[y] - E[x]E[y] + E[x]E[y] \Leftrightarrow$$

$$\text{Cov}(x, y) = E[xy] - E[x]E[y] \Rightarrow$$

$$E[x \cdot y] = \text{Cov}(x, y) + E[x]E[y]$$

$$(i): \varphi(x, x) \geq 0, \text{ για } E[x^2] = \text{Cov}(x, x) + E[x]E[x] \Rightarrow$$

$$E[x^2] = \text{Var}(x) + E^2[x] \Rightarrow$$

$$E[x^2] = \sigma^2(x) + \mu_x^2 \quad (13)$$

$$(ii): \varphi(x, y) = \varphi(y, x), \text{ για } E[x \cdot y] = \text{Cov}(x, y) + E[x]E[y] \text{ και}$$

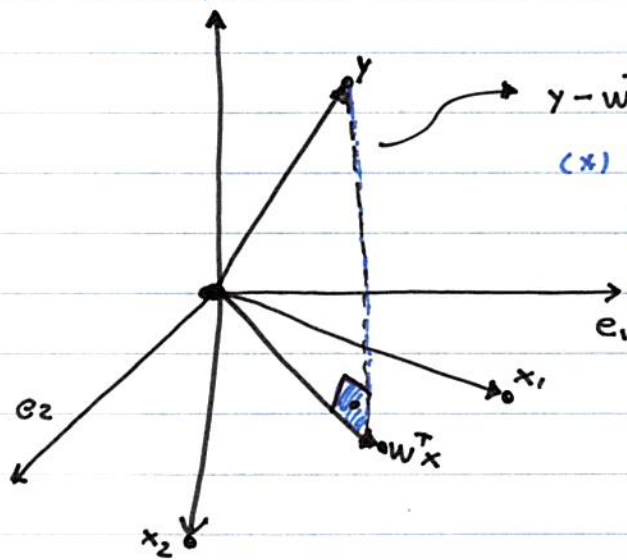
$$E[y \cdot x] = \text{Cov}(y, x) + E[x]E[y] \text{ και}$$

$$E[xy] = E[yx] \text{ καθώς } \text{Cov}(x, y) = \text{Cov}(y, x). \quad (14)$$

$$(iii): \varphi(x, \lambda y + \mu z) = \lambda \varphi(x, y) + \mu \varphi(x, z), \text{ για } \varphi(x, \lambda y + \mu z) = E[x \cdot (\lambda y + \mu z)] =$$

$$= E[\lambda xy + \mu xz] = \lambda E[xy] + \mu E[xz] = \lambda \varphi(x, y) + \mu \varphi(x, z). \quad (15)$$

Στον συγκεκριμένο διανυσματικό χώρο η ποσότητα $\underline{w}^T \underline{x} = w_1 x_1 + \dots + w_n x_n$ υποδηλώνει έναν γραμμικό συνδυασμό διανυσμάτων ο οποίος θα βρίσκεται εντός του υποχώρου που ορίζεται από τα x_i .



$y - \underline{w}^T \underline{x}$ (approximation error)
 (*) Το σφάλμα προσέγγισης του (y) από του υποχώρου που ορίζουν τα x_1 και x_2 γίνεται ελάχιστο όταν γίνεται ορθογώνιο στο εν λόγω υποχώρο.

(*) Η ποσότητα ανάλυση παρουσιάζεται γραφικά στο προηγούμενο σχήμα. Η προσέγγιση του (y) από του γραμμικού συνδυασμού $w_1 x_1 + \dots + w_n x_n$ επιφέρει ένα σφάλμα της μορφής $y - \underline{w}^T \underline{x}$.

(*) Η συνθήκη πρώτης τάξης, που εκφράζεται μέσω της σχέσης (6) $\{ E [\underline{x} \cdot (y - \underline{w}^T \underline{x})] = \underline{0} \}$, υποδηλώνει πως η βέλτιστη τιμή του διανύσματος των βερών προκύπτει όταν το διάνυσμα του σφάλματος είναι ορθογώνιο σε σχέση από τα x_i ($1 \leq i \leq n$). Με άλλα λόγια, η βέλτιστη εκτίμηση του (y) προκύπτει από μια ορθογώνια προβολή του στο διανυσματικό χώρο που ορίζεται από τα x_i .

Στοχαστική Προσέγγιση και ο Αλγόριθμος Μέσου Τετραγωνικού Σφάλματος

Η επίλυση της εξίσωσης (7) απαιτεί του υπολογισμό του πίνακα συσχέτισης (covariance matrix) \underline{R}_x και του διανύσματος ετεροσυσχέτισης $E [\underline{x} \cdot y]$ (cross-covariance vector). Ο υπολογιστής, ωστόσο, των εν λόγω διανυσμάτων προϋποθέτει γνώση των υποκείμενων κατανομών που στην γενική περίπτωση είναι άγνωστές. Εξόλητου, αν είχαμε γνώση των κατανομών αυτών για ποιο λόγο να μην χρησιμοποιούσαμε έναν Μαξίμιζαδό Ζαξισοτήτα!

Το πρόβλημα, ενόψει, του υπολογισμού ^{του} εκτιμητή των ελαχίστων τετραγώνων ανάμεσα στην επίλυση της εξίσωσης (7) χωρίς την γνώση της απαραίτητης στατιστικής πληροφορίας.

Robbins-Monro Algorithm

Το συγχευρισμένο πρόβλημα λύθηκε από τον αλγόριθμο που πρότειναν οι Herbert Robbins και Sutton Monro το 1951. Η μεθοδολογία που προτάθηκε από τους Robbins και Monro αποτελεί στην πραγματικότητα μια υπολογιστική διαδικασία προσδιορισμού των ριζών μιας συνάρτησης $f(\theta)$, η οποία όμως είναι στην ουσία η αναμενόμενη τιμή μιας συνάρτησης $F(\theta, \xi)$ η οποία εξαρτάται από την τυχαία μεταβλητή ξ . Ισχύει, δηλαδή, ότι:

$$f(\theta) = E_{\xi} [F(\theta, \xi)] \quad (16)$$

Οι αλγόριθμοι στοχαστικής προσέγγισης (Stochastic Approximation Algorithms) στοχεύουν στην ανεύρεση ιδιοτήτων της συνάρτησης $f(\cdot)$ παρουσιάζοντας την ανάγκη άμεσης αποτίμησής της. Παραδείγματα, των επιθυμητών ιδιοτήτων της $f(\cdot)$ είναι οι ρίζες της ή τα ακρότατά της.

Οι Robbins και Monro θεώρησαν το πρόβλημα της επίλυσης της εξίσωσης

$$f(\theta) = a \quad (17)$$

όπου a είναι μια σταθερά υπό την προϋπόθεση ότι η εξίσωση:

$$f(\theta) - a = 0 \quad (18)$$

έχει μια μοναδική ρίζα για $\theta = \theta^*$.

Η κεντρική ιδέα των Robbins και Monro βασίζεται στην υπόθεση πως ενώ η συνάρτηση $f(\theta)$ δεν είναι άμεσα παρατηρήσιμη, είναι δυνατή η απόκτηση δειγμάτων της τυχαίας μεταβλητής $N(\theta)$ για την οποία ισχύει ότι:

$$E[N(\theta)] = f(\theta) \quad (19)$$

Η βασική δομή του αλγορίθμου συνίσταται στην διαμεριστική μιας ακολουθίας εκτιμήσεων για την ζητούμενη τιμή του θ , η οποία έχει την ακόλουθη μορφή:

$$\theta_{n+1} = \theta_n + \alpha_n \cdot (N(\theta_n) - a) \quad (20)$$

$n \geq 0$

Η ακολουθία $\alpha_1, \alpha_2, \dots$; είναι μια ακολουθία θετικών βημάτων που χρησιμοποιούνται για την επαναζήτηση της τιμής του θ .

Οι Robbins και Monro απέδειξαν ότι η ακολουθία $(\theta_n)_{n \geq 0}$ συρτάνει στην ζητούμενη τιμή θ^* με βήμα της μετρικής L^2 .

$$\lim_{n \rightarrow \infty} \|\theta_n\|^2 = \theta^* \quad (21)$$

Η εν λόγω σύγκλιση συνεπάγεται και σύγκλιση της ακολουθίας $(\theta_n)_{n \geq 0}$ σε πιθανότητα:

$$\lim_{n \rightarrow \infty} \text{Prob} [\theta_n = \theta^*] = 1 \quad (22)$$

Τα δύο παραπάνω είδη σύγκλισης εξασφαλίζονται υπό την προϋπόθεση ότι:

$$\left. \begin{array}{l} \text{(i): } \sum_{n=0}^{\infty} \alpha_n = \infty \quad (23) \\ \text{(ii): } \sum_{n=0}^{\infty} \alpha_n^2 < \infty \quad (24) \end{array} \right\} \Rightarrow \lim_{n \rightarrow \infty} \alpha_n = 0 \quad (25)$$

★ Μια ειδική ακολουθία βημάτων που ικανοποιεί τις παραπάνω συνθήκες είναι η εξής:

$$\alpha_n = \frac{a}{n}, \quad n \geq 1 \quad \text{ή} \quad a > 0 \quad (26)$$

(*) Στην περίπτωση του ειδικού του μέσου τετραγωνικού σφάλματος η άρνηση συνάρτησης $F(\cdot, \cdot)$ έχει την ειδική μορφή $F(x, \underline{w}) = (\gamma - \underline{w}^T x)^2$, ενώ η συνάρτησης $f(\cdot)$ έχει την ειδική μορφή $f(\underline{w}) = E_x [F(x, \underline{w})]$.

(*) Από την στιγμή που αναζητούμε την ρίζα της εξίσωσης:

$$E[F(x, \underline{w})] = 0 \quad (27)$$

η εφαρμογή της μεθοδολογίας των Robbins και Μοντο οδηγεί στην υιοθέτηση του επαναληπτικού σχήματος ενημέρωσης που βασίζεται σε μία ακολουθία τυχαίων διανυσμάτων \underline{x}_k με $k = 1, 2, \dots$ τα οποία έχουν ληφθεί από την ίδια κατανομή με την συνάρτησης $F(\cdot, \cdot)$, όπου \underline{w} είναι το ζητούμενο διάνυσμα των άγνωστων παραμέτρων:

$$\hat{\underline{w}}(k+1) = \hat{\underline{w}}(k) + \alpha_{k+1} \circ F(\underline{x}_{k+1}, \hat{\underline{w}}(k)) \quad (28) \quad \text{με } \alpha_k = \frac{a}{k} \text{ ή } k^2.$$

(*) Η σύγκλιση της ακολουθίας $(\alpha_k)_{k \geq 1}$ στο μηδέν, $\alpha_k \rightarrow 0$, εξασφαλίζει πως η διαδικασία ενημέρωσης των βαρών "πορμύ" για μεγάλα τιμή του k , ισοδύναμο στο άπειρο.

(*) Οσάουτο, η πρώτη συνθήκη πάνω στην ακολουθία των βαρών $(\sum_{k=1}^{\infty} \alpha_k = \infty)$ εξασφαλίζει πως η σύγκλιση δεν επιτρέπεται πρόωρα προκειμένου να εξασφαλιστεί πως η επαναληπτική διαδικασία δεν "πορμύ" μακριά από την ιδανική λύση.

(*) Η δεύτερη συνθήκη $(\sum_{k=1}^{\infty} \alpha_k^2 < \infty)$ εξασφαλίζει πως η συσσώρευση του θορύβου που υπεισέρχεται λόγω της στοχαστικής φύσης των μεταβλητών είναι διοχερίσιμος.

#9

Ας θεωρήσουμε για παράδειγμα την εφαρμογή του αλγοριθμικού πλαισίου των Robbins και Μονρο στην επίλυση της παρακάτω απλής Εξίσωσης:

$$E[\underline{x} - \underline{w}] = 0 \quad (29) \quad \text{με} \quad F(\underline{x}, \underline{w}) = \underline{x} - \underline{w} \quad (30)$$

πάνω στην βάση ενός συνόλου τυχαίων δείγματων \underline{x}_k για $k=1, 2, \dots$.

Η εφαρμογή της αναδρομικής σχέσης για τον υπολογισμό του \underline{w}^* δίνει:

$$\hat{w}(n+1) = \hat{w}(n) + \alpha_{n+1} \cdot F(\underline{x}_{n+1}, \hat{w}(n)) \Leftrightarrow$$

$$\hat{w}(n+1) = \hat{w}(n) + \frac{1}{n+1} \cdot (\underline{x}_{n+1} - \hat{w}(n)) \Leftrightarrow$$

$$\hat{w}(n+1) = \hat{w}(n) - \frac{1}{n+1} \hat{w}(n) + \frac{1}{n+1} \underline{x}_{n+1} \Leftrightarrow$$

$$\hat{w}(n+1) = \frac{n}{n+1} \hat{w}(n) + \frac{1}{n+1} \underline{x}_{n+1} \quad (31) \quad \text{για } \underline{n} \geq 0$$

Η σχέση (31) γράφεται ως εξής:

$$\hat{w}(n+1) = \frac{n}{n+1} \left[\frac{n-1}{n} \hat{w}(n-1) + \frac{1}{n} \underline{x}_n \right] + \frac{1}{n+1} \underline{x}_{n+1} \Rightarrow$$

$$\hat{w}(n+1) = \frac{n-1}{n+1} \hat{w}(n-1) + \frac{\underline{x}_n + \underline{x}_{n+1}}{n+1} \quad (32)$$

Παρατηρούμε επίσης την γενική μορφή της αναδρομικής σχέσης:

$$\hat{w}(k) = \frac{k-1}{k} \hat{w}(k-1) + \frac{1}{k} \underline{x}_k$$

Η σχέση (32) γράφεται:

$$\hat{W}(n+1) = \frac{n-1}{n+1} \cdot \left[\frac{n-2}{n-1} \cdot \hat{W}(n-1) + \frac{1}{n-1} \cdot x_{n-1} \right] + \frac{x_n + x_{n+1}}{n+1} \Leftrightarrow$$

$$\hat{W}(n+1) = \frac{n-2}{n+1} \cdot \hat{W}(n-1) + \frac{x_{n-1} + x_n + x_{n+1}}{n+1} \quad (33)$$

Από την σχέση (33) μπορούμε να παρατηρήσουμε την παρουσία γενική μορφή ως σχέσης:

$$\hat{W}(n+1) = \frac{n-k}{n+1} \cdot \hat{W}(n+1-k) + \frac{x_{n+1-k} + \dots + x_{n+1}}{n+1} \quad \Leftrightarrow$$

$$\hat{W}(n+1) = \frac{n-k}{n+1} \cdot \hat{W}(n+1-k) + \frac{1}{n+1} \cdot \sum_{r=n+1-k}^{n+1} x_r \quad (34)$$

Η σχέση (34) ισχύει προφανώς για $0 \leq k \leq n$ (κ διαφορετικού των δειγμάτων)

Επομένως, για $k=n$, θα έχουμε ότι:

$$\hat{W}(n+1) = \frac{1}{n+1} \cdot \sum_{r=1}^{n+1} x_r \quad (35)$$

* Η σχέση (35) υποδηλώνει πως η βέλτιστη εκτίμηση για τον παράμετρο \underline{w} καθώς αυξάνεται το πλήθος των δειγμάτων σύμφωνα με την εκτίμηση του μέσου τετραγωνικού σφάλματος είναι ίση με τον δεσφραζικό μέσο των παρατηρήσεων. Η συγκεκριμένη εκτίμηση αποτελεί τον πλήρη φυσική εκτίμηση.

* Για την περίπτωση της εκτίμησης του μέσου τετραγωνικού σφάλματος με $F(x, w) = x \cdot (y - w^T x)$ θα έχουμε ότι:

$$\hat{W}(n+1) = \hat{W}(n) + \alpha_{n+1} \cdot F(x_{n+1}, \hat{W}(n)) \Rightarrow$$

$$\hat{W}(n+1) = \hat{W}(n) + \alpha_{n+1} \cdot x_{n+1} \cdot (y_{n+1} - \hat{W}(n)^T x_{n+1}) \quad (36)$$