


Module 1 - Elements of machine learning

① Regression error function : $E(w) = \frac{1}{2} \sum_{n=1}^N [y(x_n, w) - t_n]^2$

Find w that can minimize $E(w)$

② Underfitting (high train and test error)

Overttting (Low train error but high test error)

larger dataset can afford a more complex model

③ Regularization = penalty term to training objectives

$$E(w) = \frac{1}{2} \sum_{n=1}^N [y(x_n, w) - t_n]^2 + \frac{1}{2} \|w\|^2$$

④ Model selection :

1. Validation set

2. K-fold - computational expensive

3. leave one out - computational expensive

Probabilistic theory concepts

① max likelihood = max log likelihood = min negative log likelihood

② Bootstrap = sampling with replacement, compute maximum likelihood of model parameters
Higher variance = more uncertain

③ Bayes theorem : $P(w|D) = \frac{P(D|w) P(w)}{P(D)}$

KNN

Module 2 - Linear model for regression

$$y(x, w) = w \cdot \phi(x)$$

$w = [w_0 \ w_1 \ w_2]$ $\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$

$$\text{logistic sigmoid} = \frac{1}{1 + e^{-x}}$$

$$\tanh h = \frac{1 - e^{-2h}}{1 + e^{-2h}}$$

$$E(w) = \frac{1}{2} \sum_{n=1}^N \{ t_n - w \phi(x_n) \}^2$$

Gradient descent

Batch

$$w_t = w_{t-1} - \eta \times \nabla E(w_{t-1})$$

↓

$$- \sum_{n=1}^N \{ t_n - w \phi(x_n) \} \phi(x_n)$$

Stochastic

$$w^t = w^{t-1} - \eta^t \nabla E(w^{t-1})$$

$$= w^{t-1} - \eta^t \{ t_n - w \phi(x_n) \} \phi(x_n)$$

L2 regression

$$\frac{1}{2} \sum_{n=1}^N (t_n - w \phi(x))^2 + \frac{\lambda}{2} \|w\|^2$$

L1 regression

$$\frac{1}{2} \sum_{n=1}^N (t_n - w \phi(x))^2 + \frac{\lambda}{2} \|w\|$$

If λ large : some weights $\rightarrow 0$ \therefore fast prediction
more interpretable

Average prediction over all datasets differs from the desired regression function



Bias: consistently learn wrong thing

Variance: learn random things irrespective to real signals



Solutions for individual datasets vary around their average.

Complex model: high variance, low bias

Simple low high

Stochastic gradient descent

\mathcal{L}_0	x_1	x_2	x_3	x_4	/	w
1	0.46	0.33	0.56	0.49	0.3	0.2

$$w_{\text{new}} = \underline{0.2} + 0.01 \times (\text{True} - \text{pred}) \times \underline{0.46}$$

Module 3 - Models for classification

Discriminative models \Rightarrow perceptrons

Probabilistic disc model \Rightarrow logistic regression

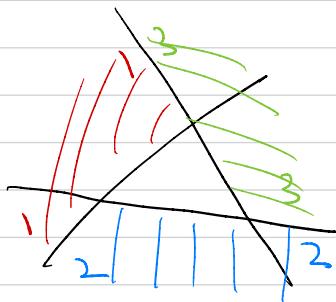
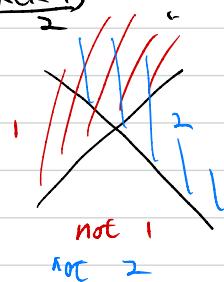
generative model \Rightarrow more parameters, higher complexity \rightarrow overfitting

$$y(x) = w_0 + w^T x \quad \begin{array}{l} \text{Assign to class 1 if } y(x) \geq 0 \\ \text{else class 2} \end{array}$$

One versus rest
One vs One

(K-1) classifier
 $\frac{K(K-1)}{2}$

Ambiguous region



Perception:

pred true

miscalculated when $w \phi(x_n) \times t_n < 0$

Training objective: minimize $- \sum_{n \in \text{em}} w \phi(x_n) t_n$

\uparrow
negative for misclassified
data point

$$w^{t+1} = w^t - \eta \nabla E(w) = w^t + \eta \underline{\phi(x_n) t_n} \leftarrow 1 / -1$$

Prob Gen Model

$$P(C_1|x) > P(C_2|x) = \frac{P(x|C_1)P(C_1)}{P(x)} > \frac{P(x|C_2)P(C_2)}{P(x)}$$

$$P(\alpha|C_1)P(C_1) > P(\alpha|C_2)P(C_2)$$

$$\alpha = \ln \frac{P(x|C_1)P(C_1)}{P(x|C_2)P(C_2)} \quad \begin{cases} \alpha > 0 \rightarrow C_1 \\ \text{else} \rightarrow C_2 \end{cases}$$

$$w_X \quad w_0$$

Max Likelihood

$$\text{Likelihood} = \prod [e^{N(x_n | \mu_1, \Sigma)}]^{t_n} [(1-e)^N(x_n | \mu_2, \Sigma)]^{(1-t_n)}$$
$$\text{Log Likelihood} = \sum t_n \ln [e^{N(x_n | \mu_1, \Sigma)}] + (1-t_n) \ln [(1-e)^N(x_n | \mu_2, \Sigma)]$$

$$\frac{dL}{dp} = e = \frac{N_1}{N_1+N_2}$$

$$\mu_1 = \frac{1}{N_1} \sum_{n=1}^N t_n x_n$$
$$\mu_2 = \frac{1}{N_2} \sum_{n=1}^N (1-t_n) x_n$$

$$\Sigma = \frac{N_1}{N} S_1 + \frac{N_2}{N} S_2 \quad S_1 = \frac{1}{N_1} \sum_{n \in C_1} (x_n - \mu_1)(x_n - \mu_1)^T$$

Logistic regression

$$P(C_1 | x) = \frac{1}{1 + e^{-w \cdot x}}$$

$$L(w) = \log \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n}$$

t_n = True class
 $y_n = P(C_1 | x)$

$$w^{t+1} = w^t - \eta^t (y_n - t_n) x_n \quad (1 \quad x_1 \quad x_2)$$

Perceptron

Phi	w	T
1	0.5 0.4	0.2 0.3 0.5

$$w_i = w \times \eta \times \text{Phi} \times T$$

weighted avg of
class corr

Bayes Classification

- 1: Calculate prior prob $p(C_k)$ for each class
- 2: calculate μ_k , S_k class cov, \sum share cov
- 3: use PDR calculate $p(x|C_k)$ for each pf.
- 4: calculate $p(C_k) p(x|C_k) \leftarrow$ find max

Module 4 - Latent Variable and EM

$$\gamma(z_{nk}) = p(z_n=k|x_n) = \frac{p_k N(x_n | \mu_k, \Sigma_k)}{\sum_i p_i N(x_n | \mu_i, \Sigma_i)}$$

$$\text{Log-likelihood} = \sum_{n=1}^N \ln p(x_n) = \sum_{n=1}^N \ln \frac{\sum_{k=1}^K p_k N(x_n | \mu_k, \Sigma_k)}{p(x_n)}$$

$$① \mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) x_n$$

↑
sum of γ

$$② \sum_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (x_n - \mu_k) (x_n - \mu_k)^T$$

$$③ p_k = \frac{N_k}{N}$$

Document Clustering

$$p(k|d) = p_k \prod_{w \in d} \frac{c(w,d)}{\mu_{kw}}$$

$\downarrow \mu_{kw}$ $\leftarrow p$

$$\gamma = p(z|d) = \frac{p(d|z) p(z)}{p(d)}$$

$$p_k = \frac{N_k}{N} \quad N_k = \sum \gamma(z_{n,k})$$

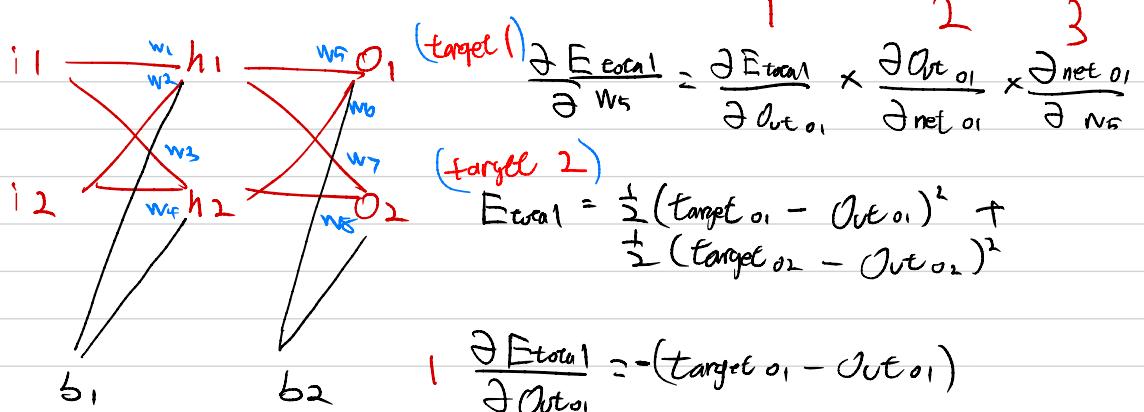
$$\mu_{kw} = \frac{\sum \gamma(z_{n,k}) c(w,d)}{\sum_{w' \in A} \sum_{n=1}^N \gamma(z_{n,k}) c(w',d_n)}$$

Module 5

Neural Network

$$\begin{aligned}\frac{\partial E}{\partial w_{jk}} &= \frac{\partial E}{\partial o} \times \frac{\partial o}{\partial w_{jk}} \\&= \frac{\partial (t_k - o_k)^2}{\partial o_k} \times \frac{\partial o_k}{\partial w_{jk}} \\&= -2(t_k - o_k) \times \frac{\partial o_k}{\partial w_{jk}} \\&= -2(t_k - o_k) \times \frac{\partial \sigma(\sum_i w_{ik} \cdot o_i)}{\partial w_{jk}} \\&= -2(t_k - o_k) \times \left[\sigma(\sum_i w_{ik} \cdot o_i)(1 - \sigma(\sum_i w_{ik} \cdot o_i)) \right. \\&\quad \left. \times \frac{\partial}{\partial w_{jk}} (\sum_i w_{ik} \cdot o_i) \right] \\&= -2(t_k - o_k) \times \sigma(\sum_i w_{ik} \cdot o_i)(1 - \sigma(\sum_i w_{ik} \cdot o_i)) \\&\quad \times o_j\end{aligned}$$

$$d_3 = \rightarrow (True - pred) * \text{derivative of sigmoid } (z_3)$$



$$\begin{aligned} \frac{\partial \text{Out}_{o_1}}{\partial \text{net}_{o_1}} &= \frac{\partial \sigma(\text{net}_{o_1})}{\partial (\text{net}_{o_1})} \\ &= \sigma(\text{net}_{o_1}) \times (1 - \sigma(\text{net}_{o_1})) \\ &= \text{Out}_{o_1} \times (1 - \text{Out}_{o_1}) \end{aligned}$$

$$\frac{\partial \text{net}_{o_1}}{\partial w_5} = \frac{\partial (w_5 \times \text{out}_{h_1} + w_6 \times \text{out}_{h_2} + b_2 \times 1)}{\partial w_5}$$

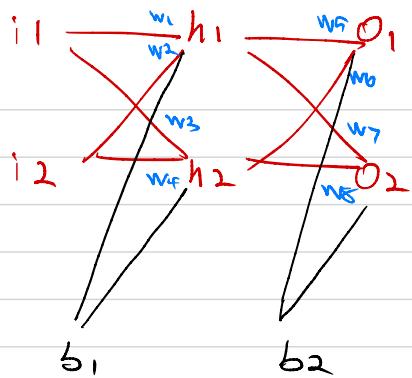
$$d_{o_1} = \text{Out}_{h_1}$$

$$\frac{\partial E_{\text{total}}}{\partial w_5} = -(\text{target}_{o_1} - \text{Out}_{o_1}) \times (\text{Out}_{o_1} \times (1 - \text{Out}_{o_1})) \times \text{Out}_{h_1}$$

$$\frac{\partial E_{\text{total}}}{\partial w_5} = d_{o_1} \times \text{Out}_{h_1}$$

SAME FOR
 w_6
 w_7
 w_8

$$w_5 = w_5 - \eta \times (d_{o_1} \times \text{Out}_{h_1})$$



$$\frac{\partial E_{\text{total}}}{\partial w_1} = \left[\frac{\partial E_{\text{total}}}{\partial \text{out}_{h1}} \right] \times \left[\frac{\partial \text{out}_{h1}}{\partial \text{net}_{h1}} \right] \times \left[\frac{\partial \text{net}_{h1}}{\partial w_1} \right]$$

$$\frac{\partial E_{\text{total}}}{\partial \text{out}_{h1}} = \frac{\partial E_{o1}}{\partial \text{out}_{h1}} + \frac{\partial E_{o2}}{\partial \text{out}_{h1}}$$

$$\frac{\partial E_{o1}}{\partial \text{out}_{h1}} = \left[\frac{\partial E_{o1}}{\partial \text{net}_{o1}} \right] \times \left[\frac{\partial \text{net}_{o1}}{\partial \text{out}_{h1}} \right]$$

$$\left(\frac{\partial E_{o1}}{\partial \text{out}_{o1}} \times \frac{\partial \text{out}_{o1}}{\partial \text{net}_{o1}} \right) \text{ from previous step}$$

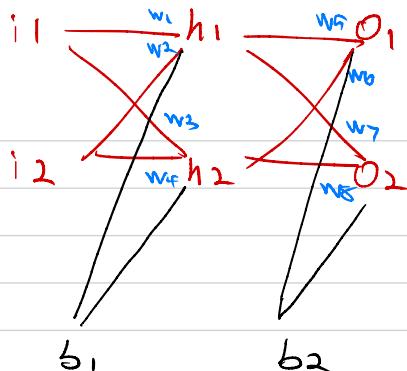
$$\frac{\partial \text{out}_{o1}}{\partial \text{net}_{o1}} \times (1 - \text{out}_{o1})$$

$$\text{out}_{h1} = \sigma(\text{net}_{h1})$$

$$\frac{\partial \text{out}_{h1}}{\partial \text{net}_{h1}} = \text{out}_{h1} \times (1 - \text{out}_{h1})$$

$$\text{net}_{h1} = i_1 \times w_1 + i_2 \times w_2 + b_1$$

$$\frac{\partial \text{net}_{h1}}{\partial w_1} = i_1$$



$$\frac{\partial E_{\text{TOTAL}}}{\partial w_5} = \frac{\partial E_{\text{TOTAL}}}{\partial \text{Out}_{o1}} \times \frac{\partial \text{Out}_{o1}}{\partial \text{Net}_{o1}} \times \frac{\partial \text{Net}_{o1}}{\partial w_5}$$

$$\frac{\partial E_{\text{TOTAL}}}{\partial \text{Out}_{o1}} = \frac{\partial E_{o1}}{\partial \text{Out}_{o1}} = \frac{\partial (t_1 - \text{Out}_{o1})^2}{\partial \text{Out}_{o1}}$$

$$= -2(t_1 - \text{Out}_{o1})$$

b_1

b_2

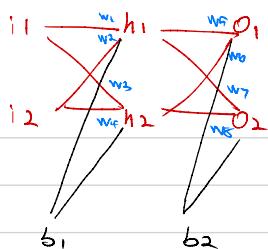
$$\frac{\partial \text{Out}_{o1}}{\partial \text{Net}_{o1}} = \text{Out}_{o1} \times (1 - \text{Out}_{o1})$$

$$\frac{\partial \text{Net}_{o1}}{\partial w_5} = \frac{\partial (b_2 + w_5 \text{Out}_{h1} + w_6 \text{Out}_{h2})}{\partial w_5}$$

$$= \text{Out}_{h1}$$

$$\frac{\partial E_{\text{TOTAL}}}{\partial w_5} = -2(t_1 - \text{Out}_{o1}) \times \text{Out}_{o1} \times (1 - \text{Out}_{o1}) \times \text{Out}_{h1}$$

$$\frac{\partial E_{\text{TOTAL}}}{\partial b_2} = -2(t_1 - \text{Out}_{o1}) \times \text{Out}_{o1} \times (1 - \text{Out}_{o1}) \times 1$$



$$\frac{\partial E_{\text{TOTAL}}}{\partial w_1} = \frac{\partial E_{\text{TOTAL}}}{\partial \text{Out}_{\text{o1}}} \times \left[\frac{\partial \text{Out}_{\text{o1}}}{\partial \text{Out}_{\text{h1}}} \right] \times \left[\frac{\partial \text{net}_{\text{h1}}}{\partial w_1} \right]$$

$$\frac{\partial E_{\text{TOTAL}}}{\partial \text{Out}_{\text{h1}}} = \boxed{\frac{\partial E_{\text{o1}}}{\partial \text{Out}_{\text{h1}}}} + \boxed{\frac{\partial E_{\text{o2}}}{\partial \text{Out}_{\text{h1}}}}$$

$$\boxed{\frac{\partial E_{\text{o1}}}{\partial \text{Out}_{\text{h1}}}} = \frac{\partial E_{\text{o1}}}{\partial \text{Out}_{\text{o1}}} \times \frac{\partial \text{Out}_{\text{o1}}}{\partial \text{net}_{\text{o1}}} \times \frac{\partial \text{net}_{\text{o1}}}{\partial \text{Out}_{\text{h1}}}$$

$$= -2(t_1 - \text{Out}_{\text{o1}}) \times \text{Out}_{\text{o1}} \times (1 - \text{Out}_{\text{o1}}) \times w_5$$

$$\boxed{\frac{\partial E_{\text{o2}}}{\partial \text{Out}_{\text{h1}}}} = \frac{\partial E_{\text{o2}}}{\partial \text{Out}_{\text{o2}}} \times \frac{\partial \text{Out}_{\text{o2}}}{\partial \text{net}_{\text{o2}}} \times \frac{\partial \text{net}_{\text{o2}}}{\partial \text{Out}_{\text{h1}}}$$

$$= -2(t_2 - \text{Out}_{\text{o2}}) \times \text{Out}_{\text{o2}} \times (1 - \text{Out}_{\text{o2}}) \times w_7$$

$$\boxed{\frac{\partial \text{Out}_{\text{h1}}}{\partial \text{net}_{\text{h1}}}} = \text{Out}_{\text{h1}} \times (1 - \text{Out}_{\text{h1}}) \quad \text{h.d. } (\Sigma_2)$$

$$\boxed{\frac{\partial \text{net}_{\text{h1}}}{\partial w_1}}, \frac{\partial (\text{w}_1 \times i_1 + \text{w}_2 \times i_2 + b_1)}{\partial w_1} = i_1 \quad \text{d}_1$$

w_L

d_2

Linear Regression Walk through

w

$$0.25 \quad 0.82 \quad 0.25 \quad 0.11 \quad 0.41$$

Observations

	x_0	x_1	x_2	x_3	x_4	T
①	1	2	3	4	5	1
②	2	3	4	5	6	1

Stochastic

gradient, if batch, then sum all

$$\begin{aligned} w_0 &= w_0 + \eta \left[(T_i - \hat{y}_i) * \Phi[i, j] \right] \\ &= 0.25 + \eta (1 - \hat{y}_1) * 1 \end{aligned}$$

$$\text{Error func} = \frac{1}{2} (T - \hat{y})^2$$

$$\text{Gradient} = -(T - \hat{y}) * \Phi[i, j]$$

Regularisation - Ridge

$$\text{Error func} = \frac{1}{2} (T - \hat{y})^2 + \lambda \sum_{i=0}^n b_i^2$$

Perceptron Walkthrough

$$\text{Error Func} = - \sum \hat{y} \times t_n \\ = - \sum w \cdot \phi(x) \cdot t_n$$

For a single pt : $-w \cdot \phi(x) \cdot t_n$

$$- (w_0 x_0 + w_1 x_1 + w_2 x_2) \cdot t_n$$

$$\text{Gradient} = -t_n x_i \quad \text{for } w_i$$

$$\text{Update} = w_i = w_i - \eta \nabla \\ = w_i + \eta t_n x_i$$

W

$$w_0 \quad w_1 \quad w_2 \\ 0.13 \quad 0.49 \quad 0.18$$

Observation	T
1 -0.28 3.46	-1

$$\text{pred} = 0.13 - 0.28 \times 0.49 + 0.18 \times 3.46 \\ = 0.61 \Rightarrow \text{classified as } +1$$

Wrong classification = need update

Bayes Classifier

1. Calculate $P(C)$
2. Calculate $M_C \Sigma_C$
3. Use result in step 2 to calculate $P(x|C)$
4. Calculate $P(C|x)$ using the result from above

Logistic regression

$$\text{Log likelihood} = \log \sum y_n t_n (1 - y_n)^{1-t_n}$$

$$\text{Gradient} = (y_n - t_n) x_n$$

$$\text{update} = -\eta (y_n - t_n) x_n$$

	A	B	C	D	
Doc 1	1	3	7	8	← count
Doc 2	2	5	1	2	

$$\begin{matrix} r_1 & 0.2 & 0.3 & 0.4 & 0.1 \\ r_2 & 0.3 & 0.5 & 0.1 & 0.1 \end{matrix}$$

2 cluster $\leftarrow 1, 2$

$$M_{1,A} = 0.2 \times 1 + 0.3 \times 2$$

$$0.2 \times 1 + 0.2 \times 3 + 0.2 \times 7 + 0.2 \times 8 + 0.3 \times 2 + 0.3 \times 5 + 0.3 \times 1 + 0.3 \times 2$$

$$\ell_k = \frac{N_k}{N} \quad N_k = \sum_i^n r_i$$

$$\begin{aligned} & -W \Phi(x) t_n \\ & + \eta \Phi(x) t_n \end{aligned}$$

For l to N :

For l to K :

$$Y_{nlk} = \frac{p_k N(x_n | M_k, \Sigma_k)}{\sum_{k=1}^K p_k N(x_n | M_k, \Sigma_k)}$$

$$N_{lk} = N_{lk} + Y_{nlk}$$

$$M_{lk} = M_{lk} + Y_{nlk} x_n$$

$$\Sigma_{lk} = \Sigma_{lk} + Y_{nlk} x_n x_n^\top$$

For l to K :

$$p_k = N_{lk} / N$$

$$M_{lk} = M_{lk} / N_{lk}$$

$$\Sigma_{lk} = \Sigma_{lk} / N_{lk} - M_{lk} M_{lk}^\top$$

$$\begin{aligned} & \sum_{m=1}^M n_{km} / \sum_{k=1}^K \sum_{m=1}^M n_{km} \\ & \sum n_{km} / \sum_{m=1}^M n_{km} \\ & \sum_{m=1}^M \sum_{k=1}^K n_{km} / \sum_{m=1}^M n_{km} \end{aligned}$$

For l to N

$$k^* = \arg \min_l \text{dist}(x_n, M_k)$$

$$\tilde{M}_{lk} = M_k + x_n$$

$$n_{lk} = n_{lk} + 1$$

$$\sum_{m=1}^M \tilde{M}_{km} / \sum_{m=1}^M n_{km}$$

For i to N^2

For i to K :

$$N_{ik} = N_k + \gamma_{ik} N_k$$

$$M_{ik} = M_k + \gamma_{ik} x_n$$

$$\Sigma_{ik} = \Sigma_k + \gamma_{ik} x_n x_n^T$$

For i to K :

$$N_{ik} = N_k / N$$

$$M_{ik} = M_k / N_k$$

$$\Sigma_{ik} = \Sigma_k / N_k - M_k M_k^T$$

$$N_{ik} = \sum_{m=1}^m n_{k,m} / \sum_{k=1}^K \sum_{m=1}^m n_{k,m}$$

$$M_{ik} = \sum_{m=1}^m M_{k,m} / \sum_{m=1}^m n_{k,m}$$

$$\Sigma_{ik} = \sum_{m=1}^m \Sigma_{k,m} / \sum_{m=1}^m n_{k,m} - M_k M_k^T$$

$$\gamma_{ik} = \frac{p_{ik}}{\sum_{k=1}^K p_{ik}} \frac{N(x_i | M_k, \Sigma_k)}{N(x_i | M_k, \Sigma_k)}$$

$$N_{ik} = \sum_{n=1}^N \gamma_{n,k}$$

$$M_{ik} = \frac{1}{N_{ik}} \left(\sum_{n=1}^N \gamma_{n,k} x_n \right)$$

$$p_{ik} = \frac{N_{ik}}{N}$$

$$\Sigma_{ik} = \frac{1}{N_{ik}} \sum_{n=1}^N \gamma_{n,k} (x_n - M_{ik}) (x_n - M_{ik})^T$$

