

Topic 1 Discriminant Function

So far, we are not strange to classification. In "Probability" chapter, we have know MLE and MAP. They can be used to do classification. We will introduce another method which uses linear model.

1. Definition of Discriminant Function.

Discriminant is a function which turns input \underline{x} to one of k classes, denoted C_k . Let's limit ourselves to linear discriminants at present.

Our model is

$$y(\underline{x}) = f(\underline{w}^T \underline{x} + b)$$

active function
can be non-linear

linear discriminant
function.

It is true that, because of " $f(\cdot)$ ", our model is non-linear. But, if we see the \underline{x} that

$f(\underline{w}^T \underline{x} + b) = \text{const}$, all \underline{x} satisfying the condition will form up a linear hyper plane.

monotonous
function

① Two classes.

Aka, $k=2$.

Let's start from the simplest model. $\underline{x} \in \mathbb{R}^D$

$$y(\underline{x}) = \underline{w}^T \underline{x} + b$$

Our classification strategy is

$$\begin{cases} y(\underline{x}) \leq 0 \rightarrow C_1 \\ y(\underline{x}) > 0 \rightarrow C_2 \end{cases}$$

$\{\underline{x} : y(\underline{x}) = 0\}$ will form up a decision hyperplane of $D-1$

dimension

Suppose $\underline{x}_A, \underline{x}_B$ are points $y(\underline{x}_A) = y(\underline{x}_B) = 0$

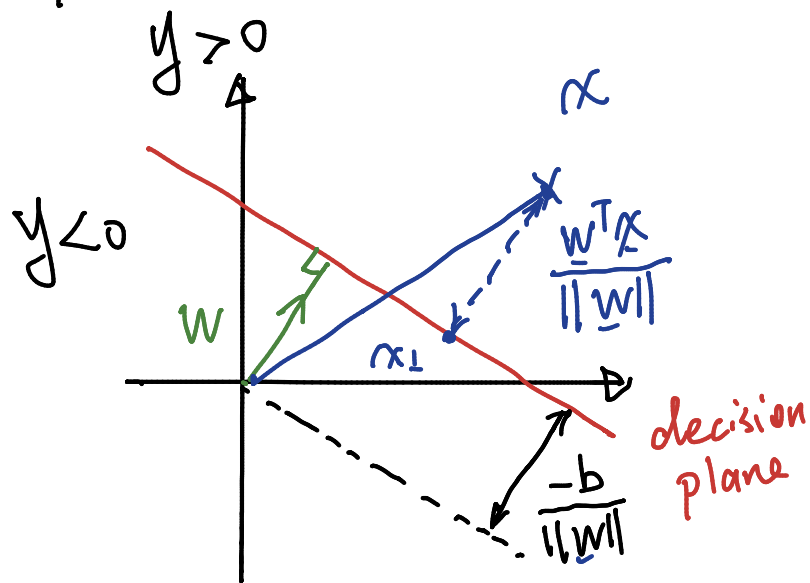
$$\underline{w}^T (\underline{x}_A - \underline{x}_B) = 0 \Leftrightarrow \underline{w} \perp \underline{x}_A - \underline{x}_B$$

\vec{w} is orthogonal to any vector within decision plane.

$\Leftrightarrow \vec{w}$ is the normal vector of decision plane.

Also, if \underline{x} is on decision plane

$$\frac{\underline{w}^T \underline{x}}{\|\underline{w}\|} = \frac{-b}{\|\underline{w}\|}$$



If the projection of \underline{x} onto decision plane is \underline{x}_\perp , then

$$\left. \begin{aligned} \underline{x} &= \underline{x}_\perp + \underline{r} \cdot \frac{\underline{w}}{\|\underline{w}\|} \\ y(\underline{x}) &= \underline{w}^T \underline{x} + b \\ y(\underline{x}_\perp) &= \underline{w}^T \underline{x}_\perp + b \end{aligned} \right\}$$

$$\underline{r} = \frac{y(\underline{x})}{\|\underline{w}\|}$$

The distance between \underline{x} and decision plane.

Let's use homogenous coordinate the linear function

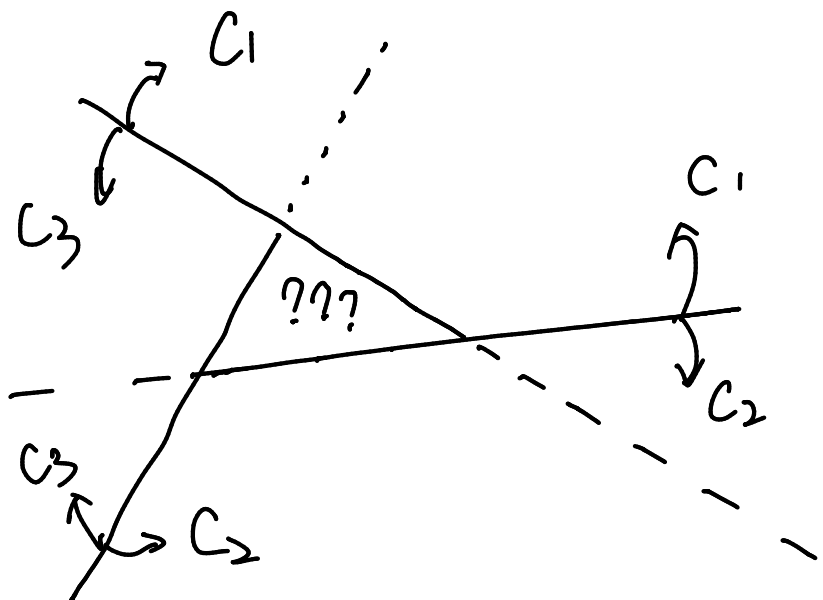
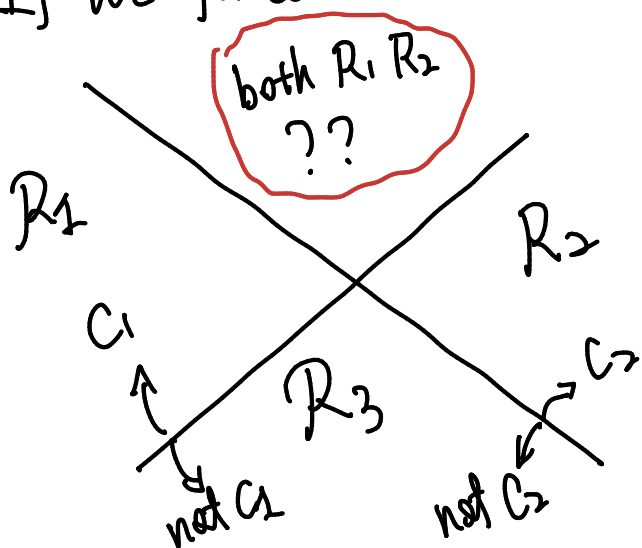
$$y(\underline{x}) = [\underline{w}, b] \begin{bmatrix} \underline{x} \\ 1 \end{bmatrix}$$

2. Multiple Class.

When $k > 2$, things turn to be complicated

If we follow our method in section 1, we may have

We have a region can be both C_1 & C_2 .



If we create decision boundary for every pair of Class, problems still happen, but in another mode.

We need a compound discriminant function!

We have a series of discriminant function

$$y_k(\underline{x}) = \underline{w}_k^T \underline{x} + b_k \quad k = 1, \dots, K$$

we assign \underline{x} to class C_k iff $y_k(\underline{x}) > y_j(\underline{x}) \quad j \neq k$.

Let's investigate the decision boundary.

for C_k & C_j

$$\begin{aligned} & (\underline{w}_k - \underline{w}_j)^T \underline{x} + (b_k - b_j) = 0 \\ \hookrightarrow & \underline{w}'^T \underline{x} + b' = 0 \quad \text{same form as 2-class case.} \end{aligned}$$

According to knowledge from Convex set, the decision region is a union of several half plane, so it is convex.

We can prove this by definition

Proof:

Suppose we have $\underline{x}_A, \underline{x}_B$ in R_k . According to definition,

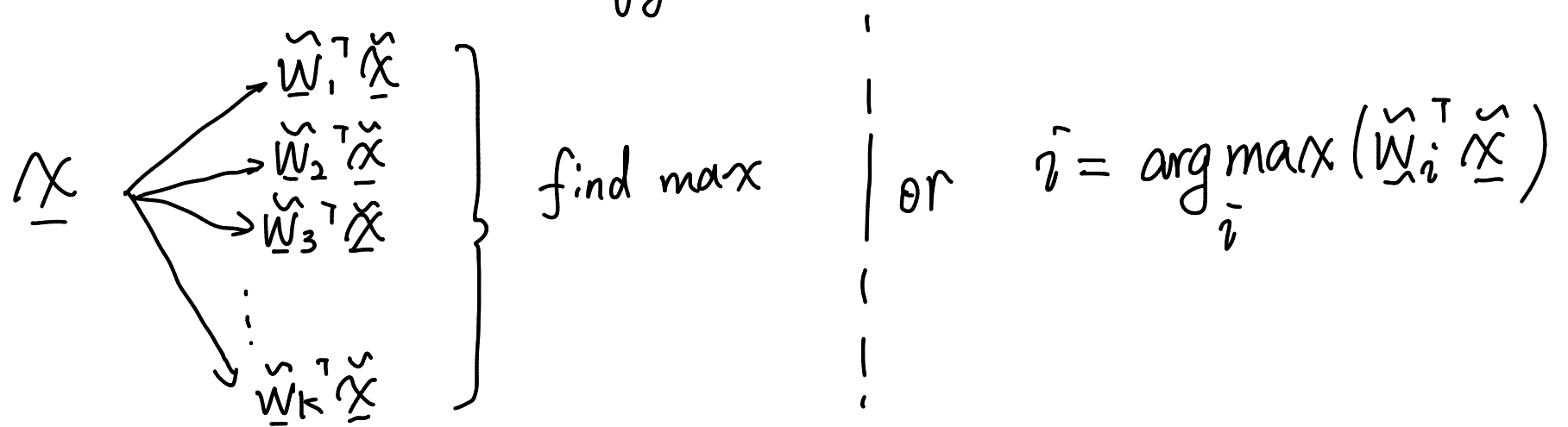
$$\hat{\underline{x}} = \lambda \hat{\underline{x}}_A + (1-\lambda) \hat{\underline{x}}_B \quad \lambda \in [0, 1]$$

From linearity

$$y(\hat{\underline{x}}) = \lambda y(\hat{\underline{x}}_A) + (1-\lambda) y(\hat{\underline{x}}_B) \Rightarrow y_k(\hat{\underline{x}}) < y_j(\hat{\underline{x}}) \Rightarrow \hat{\underline{x}} \text{ still in } R_k$$

3. Least Square for classification

Let's claim our strategy.



Suppose we have a dataset $\{\underline{x}_n, t_n\} \quad n = 1, \dots, N$.

$$\tilde{X} = \begin{bmatrix} \tilde{x}_1^T \\ \vdots \\ \tilde{x}_N^T \end{bmatrix} \quad T = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix} \quad \tilde{W} = \begin{bmatrix} w_1 & w_2 & \dots & w_k \\ b_1 & b_2 & \dots & b_k \end{bmatrix}$$

Let's choose the following sum-of-squares error function

$$\begin{aligned} E_D(\tilde{W}) &= \frac{1}{2} \text{Tr} \{ (\tilde{X}\tilde{W} - T)^T (\tilde{X}\tilde{W} - T) \} \\ &= \frac{1}{2} \left\{ \text{Tr} \{ \tilde{W}^T \tilde{X}^T \tilde{X} \tilde{W} \} - \text{Tr} \{ \tilde{W}^T \tilde{X}^T T \} - \text{Tr} \{ T^T \tilde{X} \tilde{W} \} + \text{Tr} \{ T^T T \} \right\} \end{aligned}$$

↔ Same
↔ No \tilde{W}

$$\begin{aligned} \frac{\partial}{\partial \tilde{W}} E_D(\tilde{W}) &= \frac{1}{2} \frac{\partial}{\partial \tilde{W}} \text{Tr} \{ \tilde{W}^T \tilde{X}^T \tilde{X} \tilde{W} \} - \frac{\partial}{\partial \tilde{W}} \text{Tr} \{ T^T \tilde{X} \tilde{W} \} \\ &= \tilde{X}^T \tilde{X} \tilde{W} - \tilde{X}^T T = 0 \Rightarrow \tilde{W}_{LS} = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T T \\ &= \tilde{X}^+ T \end{aligned}$$

$$\therefore y(\underline{x}) = T^T (\tilde{X}^+)^T \tilde{x}$$

An Interesting property is that.

if we can find \underline{a}, b that

$$\underline{a}^T \underline{t}_j + b = 0 \quad j = 1, \dots, N \text{ then}$$

$$\underline{a}^T y(\underline{x}) + b = \underline{a}^T T^T (\tilde{X}^T)^T \tilde{\underline{x}} + b$$

$$a^T y(\underline{x}) + b = 0$$

Recall that

$$T^T = [t_1, t_2, \dots, t_N], \text{ so } \underline{a}^T T^T = -b \mathbb{1}_{1 \times N}$$

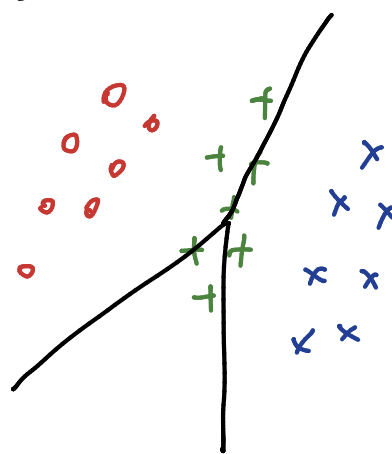
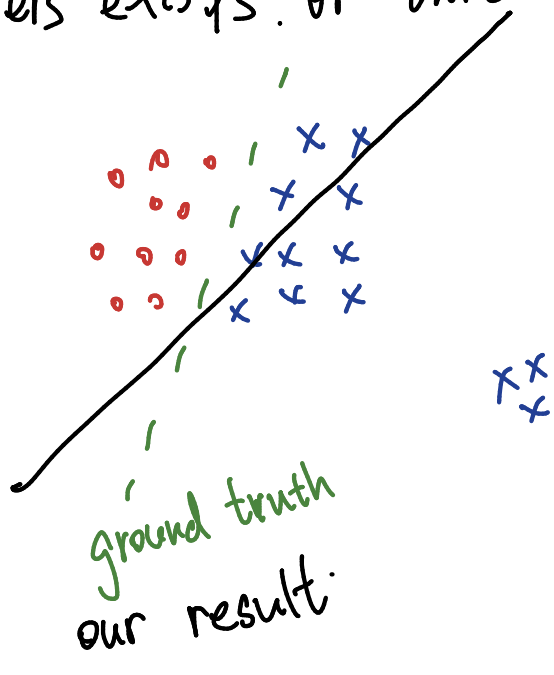
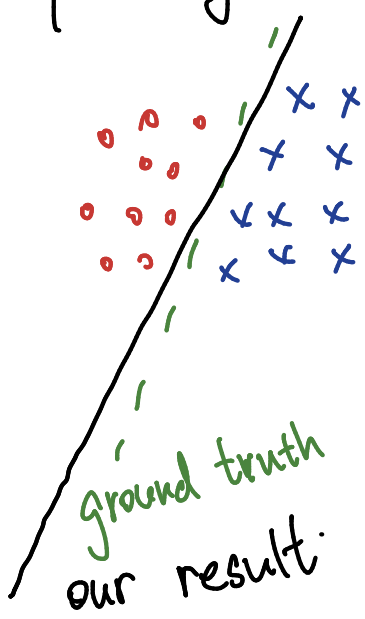
$$(\tilde{X}^T)^T \tilde{\underline{x}} = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix}$$

So, if we use 1-bit hotkey, $\underline{a} = \mathbb{1}$, then

$$\mathbb{1}^T \cdot \underline{t}_n = 1.$$

So, $y(\underline{x})$'s summation of all elements is 1. But it may not be interpreted as probability, because the element may not guaranteed to be in $[0, 1]$

We should be clear that least-square method lacks robustness, especially when outliers exists. or three classes show in 2-dimension



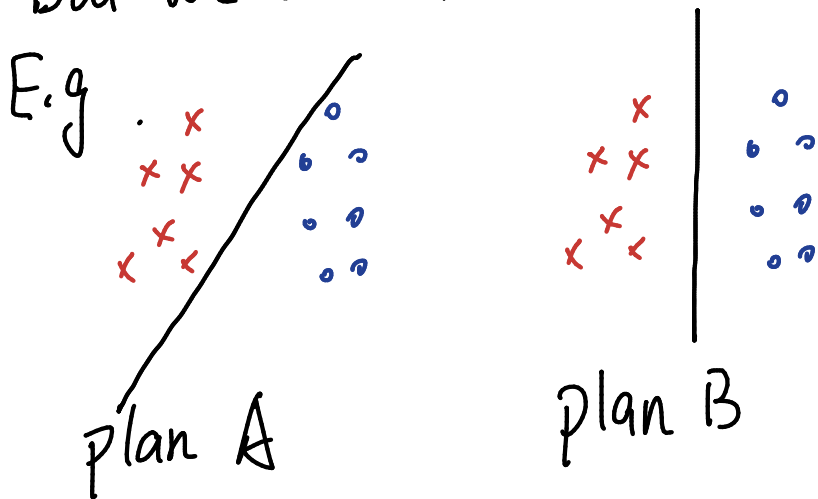
4. Fisher Linear Discriminant

In previous section, our method is project \underline{x} onto one dimension

by $y = \underline{w}^T \underline{x}$

and find out whether y is > 0 or < 0 .

But we didn't make use of points fully



Obviously, plan B is better. because the gap between points & boundary is larger.

Our target in this section is to find a good direction of boundary by adjust w to maximize the gap.

Suppose we have N_1 points for C_1 , N_2 points for C_2 . The mean vector for two classes are

$$\underline{m}_1 = \frac{1}{N_1} \sum_{n \in C_1} \underline{x}_n \quad \underline{m}_2 = \frac{1}{N_2} \sum_{n \in C_2} \underline{x}_n$$

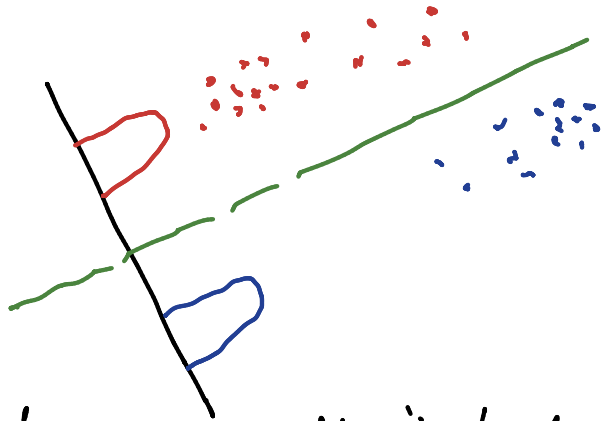
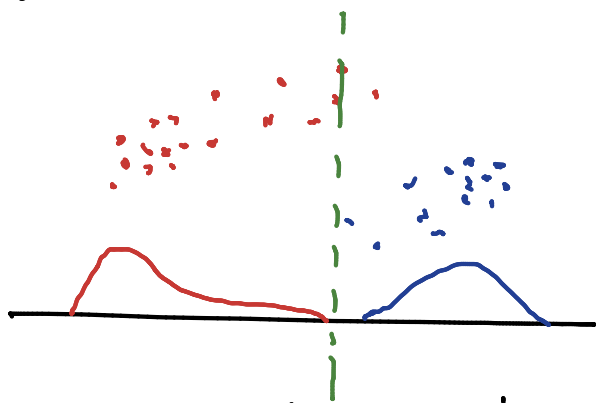
One simple strategy is to maximize the projection of $\underline{m}_1, \underline{m}_2$

$$\arg \max_{\underline{w}} \underline{w}^T (\underline{m}_2 - \underline{m}_1)$$

As you may feel, we need to constrain the magnitude of \underline{W} .
 The overall expression is

$$\operatorname{argmax}_{\underline{W}, \|\underline{W}\|=1} \underline{W}^T (\underline{m}_1 - \underline{m}_2)$$

Besides from separation of class's center, we want projected points have low variance.



The left decision has large "center separation" but high inner-class variance and not that good performance. The right decision is obviously better. So, let's give out a new objective function

$$J(\underline{W}) = \frac{[\underline{W}^T (\underline{m}_1 - \underline{m}_2)]^2}{\underbrace{\sum_{n \in C_1} [\underline{W}^T (\underline{m}_1 - \underline{x}_n)]^2 + \sum_{n \in C_2} [\underline{W}^T (\underline{m}_2 - \underline{x}_n)]^2}_{\text{projected points' var}}}$$

$$= \frac{\underline{W}^T S_B \underline{W}}{\underline{W}^T S_W \underline{W}}$$

$$S_B = (\underline{m}_2 - \underline{m}_1)(\underline{m}_2 - \underline{m}_1)^T$$

$$S_W = \sum_{n \in C_1} (\underline{m}_1 - \underline{x}_n)(\underline{m}_1 - \underline{x}_n)^T + \sum_{n \in C_2} (\underline{m}_2 - \underline{x}_n)(\underline{m}_2 - \underline{x}_n)^T$$

Our next target is to find out \underline{W}^*

$$\frac{\partial}{\partial \underline{W}} J(\underline{W}) = \frac{2 \underline{W}^T \underline{S}_W \underline{W} \underline{S}_B \underline{W} - 2 \underline{W}^T \underline{S}_B \underline{W} \underline{S}_W \underline{W}}{\underline{W}^T \underline{S}_W \underline{W} \underline{W}^T \underline{S}_W \underline{W}} = 0$$

$$\Rightarrow \underbrace{\underline{W}^T \underline{S}_B \underline{W}}_{\text{Scalar}} \underline{S}_W \underline{W} = \underbrace{\underline{W}^T \underline{S}_W \underline{W}}_{\text{Scalar}} \underline{S}_B \underline{W}$$

We don't care scalar, or the magnitude of \underline{W} , let's use

$$\begin{cases} \underline{W}^T \underline{S}_B \underline{W} = a \\ \underline{W}^T \underline{S}_W \underline{W} = b \end{cases} \quad a \underline{S}_W \underline{W} = b \underline{S}_B \underline{W}$$

$$\underline{S}_B \underline{W} = (\underline{m}_2 - \underline{m}_1) \underbrace{(\underline{m}_2 - \underline{m}_1)^T \underline{W}}_{\text{scalar} \leftarrow \text{use } C}$$

Then we have

$$a \underline{S}_W \underline{W} = b \cdot C (\underline{m}_2 - \underline{m}_1) \Rightarrow \underline{W} = \alpha \cdot \underline{S}_W^{-1} (\underline{m}_2 - \underline{m}_1)$$

Recall that $\|\underline{W}\| = 1$, the final answer is

$$\underline{W} = \frac{\underline{S}_W^{-1} (\underline{m}_2 - \underline{m}_1)}{\|\underline{S}_W^{-1} (\underline{m}_2 - \underline{m}_1)\|}$$

Fisher's Linear Discriminant

5. The perceptron algorithm

This is also a linear discriminant model but can work online. This algorithm can be viewed as the first generation of Gradient Descent learning.

The model is

$$y(x) = f(w^T \phi(x))$$

$\phi(x)$ is a fixed kernel function.

w is the parameter that we want to learn

$f(\cdot)$ is the active function given by a step function

$$f(x) = \begin{cases} +1 & x \geq 0 & \dots \text{Class } C_1 \\ -1 & x < 0 & \dots \text{Class } C_2 \end{cases}$$

We need to sequentially update w to make the classifier more and more accurate.

The update is made based on minimizing error
Error function, also known as perceptron criterion is

$$E_p(w) = \sum_{n \in M} w^T \phi(x) t_n$$

\nwarrow miss labeled data.

Note that, if the algorithm miss-labels sample i , then

$$w^T \phi(x_i) t_i < 0$$

So the definition above is reasonable

Then the update turns to be

$$W^+ = W - \eta \nabla_w E_p(w) = W + \eta \phi(x) t$$