

Machine Learning I: Supervised Methods

B. Keith Jenkins

Announcements

- Homework 3 is due Friday
- Slido poll questions today
 - Event code: 3183558
 - Go to slido.com
 - Join as a participant
 - When asked, fill in your name and use email address
 - Leave your browser window open during lecture

Today's lecture

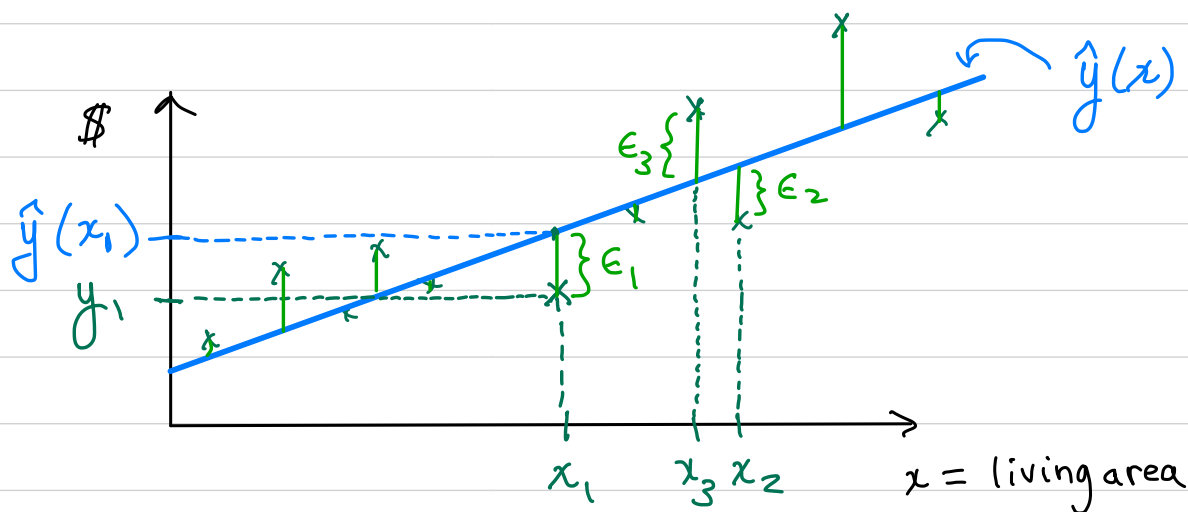
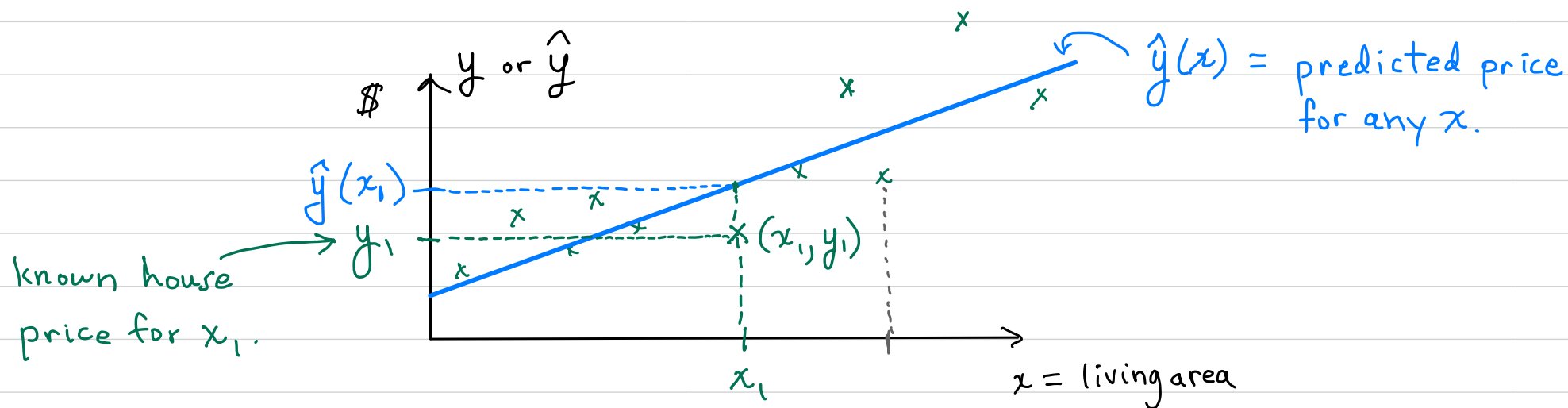
- Mean-squared-error (MSE) techniques for regression
 - Criterion function
 - Least squares
 - LMS
- MSE techniques for classification (1)
 - Criterion function and interpretation

⑤ shows locations of Slido questions.

Mean-squared-error techniques: Regression (augmented notation)

Ex: predict house prices: $\hat{y}(\underline{x})$ = price prediction
 y_i = actual price for data point \underline{x}_i .

Assume: linear model



slide notation: $\hat{y}(x_i)$ means $\hat{y}(x_i)$. p.3

Criterion function: $J(\underline{w}) = \frac{1}{N} \sum_{i=1}^N \epsilon_i^2$

$w^T x_i$ means $\underline{w}^T \underline{x}_i$
 w_0 means w_0 .

$$J(\underline{w}) = \frac{1}{N} \sum_{i=1}^N [\hat{y}(x_i) - y_i]^2 = \text{MSE} = \text{mean-squared error}$$

5

Linear model: $\hat{y}(x_i) = \underline{w}^T \underline{x}_i$

$$\Rightarrow J(\underline{w}) = \frac{1}{N} \sum_{i=1}^N [\underline{w}^T \underline{x}_i - y_i]^2$$

Let $N = N_{Tr}$

house price prediction known house price

Let $\underline{X}_{Tr} = \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_N^T \end{bmatrix}$ = data matrix, (training data)

$\underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$ = data-output (or target-output) vector

$$\Rightarrow \underline{X}_{Tr} \underline{w} = \begin{bmatrix} \underline{x}_1^T \underline{w} \\ \underline{x}_2^T \underline{w} \\ \vdots \\ \underline{x}_N^T \underline{w} \end{bmatrix} = \begin{bmatrix} \hat{y}(x_1) \\ \hat{y}(x_2) \\ \vdots \\ \hat{y}(x_N) \end{bmatrix} \triangleq \underline{\hat{y}} = \text{vector of output predictions}$$

We can write $J(\underline{w})$ as:

$$J(\underline{w}) = \frac{1}{N} \left\| \underline{X} \underline{w} - \underline{y} \right\|_2^2 = \frac{1}{N} (\underline{X} \underline{w} - \underline{y})^T (\underline{X} \underline{w} - \underline{y})$$

For learning (training), choose $\underline{X} = \underline{X}_{Tr}$ and $N = N_{Tr}$.

Minimization of $J(\underline{w})$:

Is $J(\underline{w})$ differentiable?

Is $J(\underline{w})$ convex?

Possible approaches:

1. $\nabla_{\underline{w}} J(\underline{w}) = \underline{0}$

Solve algebraically

\Rightarrow least squares

or ordinary least squares

or pseudoinverse

2. Gradient descent

Sequential or stochastic GD

\Rightarrow least mean squares (LMS)

Minimize $J(\underline{w})$ algebraically:

$$J(\underline{w}) = \frac{1}{N} (\underline{X}\underline{w} - \underline{y})^T (\underline{X}\underline{w} - \underline{y}) = \frac{1}{N} \left[\underline{w}^T \underline{X}^T \underline{X} \underline{w} - \underline{w}^T \underline{X}^T \underline{y} - \underline{y}^T \underline{X} \underline{w} + \underline{y}^T \underline{y} \right]$$

$$= \frac{1}{N} \left[\underline{w}^T \underline{X}^T \underline{X} \underline{w} - 2 \underline{w}^T \underline{X}^T \underline{y} + \underline{y}^T \underline{y} \right]$$

$$\nabla_{\underline{w}} J(\underline{w}) = \frac{1}{N} \left[2 \underline{X}^T \underline{X} \underline{w} - 2 \underline{X}^T \underline{y} \right]$$

$$\nabla_{\underline{w}} J(\hat{\underline{w}}) = \underline{0} = \frac{1}{N} \left[2 \underline{X}^T \underline{X} \hat{\underline{w}} - 2 \underline{X}^T \underline{y} \right]$$

$$\underline{X}^T \underline{X} \hat{\underline{w}} = \underline{X}^T \underline{y}$$

If $\underline{X}^T \underline{X}$ is nonsingular, then:

$$\hat{\underline{w}} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$$

Let $\underline{X}^- \triangleq (\underline{X}^T \underline{X})^{-1} \underline{X}^T$ = Moore-Penrose (left) pseudoinverse of \underline{X}

Then $\hat{\underline{w}} = \underline{X}^- \underline{y}$ ($\underline{X} = \underline{X}_{\text{tr}}$, augmented)

↗ Solution $\hat{\underline{w}}$ for least-squares regression

Comments on $\underline{\underline{X}}^{-}$

$$1. \underline{\underline{X}}^{-} \underline{\underline{X}} = \left[(\underline{\underline{X}}^T \underline{\underline{X}})^{-1} \underline{\underline{X}}^T \right] \underline{\underline{X}} = (\underline{\underline{X}}^T \underline{\underline{X}})^{-1} (\underline{\underline{X}}^T \underline{\underline{X}}) = \underline{\underline{I}}$$

$$2. \underline{\underline{X}} \underline{\underline{X}}^{-} \neq \underline{\underline{I}} \text{ in general.}$$

Comments on least-squares solution:

⑤

$$1. \text{ Predictions on dataset } \mathcal{D}: \underline{\underline{\hat{y}}} = \underline{\underline{X}}_{\mathcal{D}} \underline{\underline{\hat{w}}} = \underline{\underline{X}}_{\mathcal{D}} (\underline{\underline{X}}^{-} \underline{\underline{y}})$$

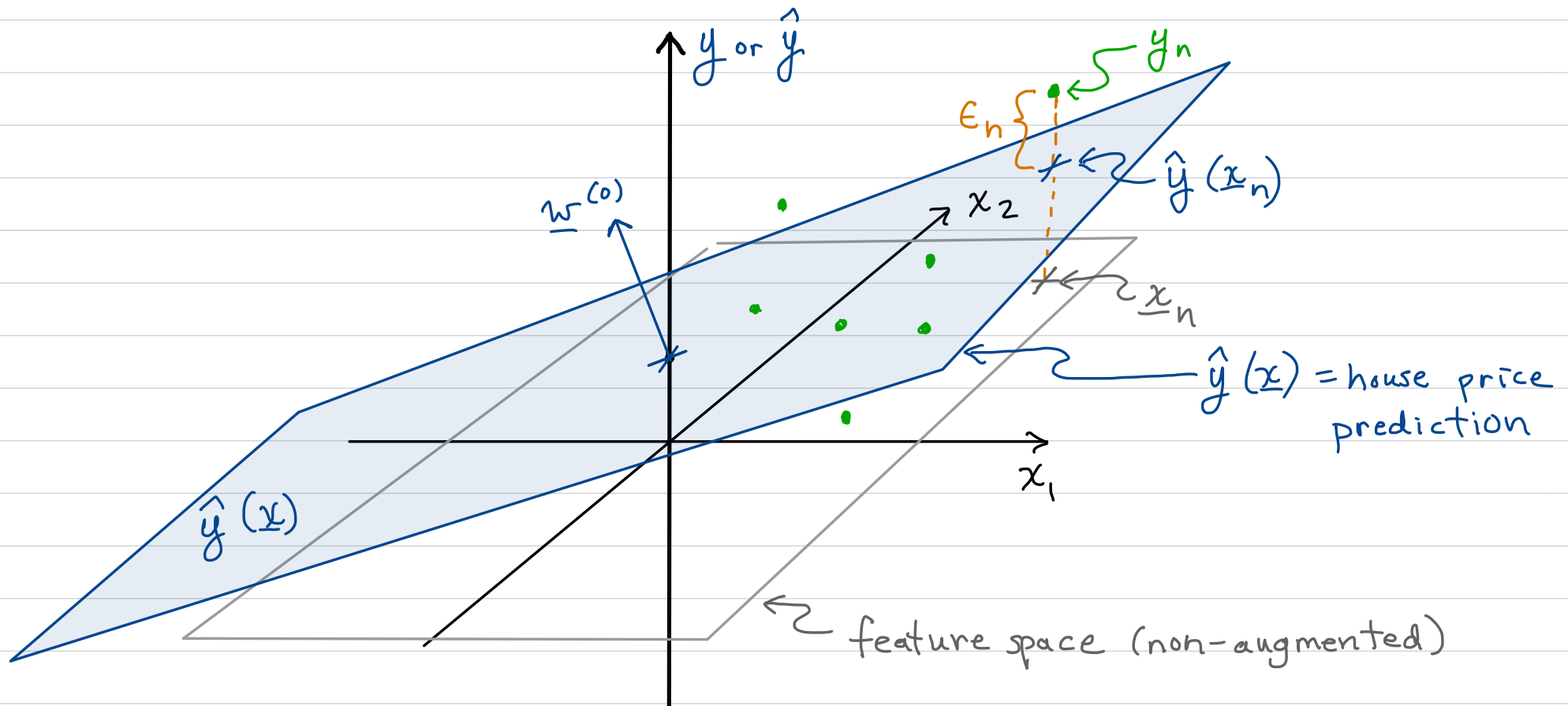
2. \therefore Mean-squared error on dataset \mathcal{D} :

$$J(\underline{\underline{\hat{w}}}) = \frac{1}{N} \left\| \underline{\underline{y}} - \underline{\underline{\hat{y}}} \right\|_2^2 = \frac{1}{N} \left\| \underline{\underline{y}} - \underline{\underline{X}}_{\mathcal{D}} \underline{\underline{X}}^{-} \underline{\underline{y}} \right\|_2^2 = \frac{1}{N} \left\| (\underline{\underline{I}} - \underline{\underline{X}}_{\mathcal{D}} \underline{\underline{X}}^{-}) \underline{\underline{y}} \right\|_2^2$$

MSE Regression: plot in feature space

Criterion function: $J(\underline{w}) = \frac{1}{N} \sum_{i=1}^N [\hat{y}(x_i) - y_i]^2 = \text{MSE}$

Linear model: $\hat{y}(x_i) = \underline{w}^T \underline{x}_i \Rightarrow J(\underline{w}) = \frac{1}{N} \sum_{i=1}^N [\underline{w}^T \underline{x}_i - y_i]^2$



$$J_n(\underline{w}) = \frac{1}{N} \epsilon_n^2 ,$$

$$J(\underline{w}) = \sum_{n=1}^N J_n(\underline{w})$$

MSE Regression: Least-Mean-Squares (LMS)

→ Sequential GD (or stochastic GD) optimization of MSE criterion.

$$J(\underline{w}) = \frac{1}{N} \sum_{n=1}^N [\hat{y}(\underline{x}_n) - y_n]^2 = \frac{1}{N} \sum_{n=1}^N [\underline{w}^T \underline{x}_n - y_n]^2 = \sum_{n=1}^N J_n(\underline{w})$$

$$J_n(\underline{w}) = \frac{1}{N} [\underline{w}^T \underline{x}_n - y_n]^2$$

$$\nabla_{\underline{w}} J_n(\underline{w}) = \frac{2}{N} [\underline{w}^T \underline{x}_n - y_n] \underline{x}_n$$

Sequential GD:

$$\underline{w}(i+1) = \underline{w}(i) - \eta'(i) \frac{2}{N} (\underline{w}(i)^T \underline{x}_n - y_n) \underline{x}_n, \quad n = (i \bmod N) + 1$$

(i: 0, 1, 2, ...; n: 1, 2, ..., N)

$$\text{Let } \eta(i) = \frac{2}{N} \eta'(i)$$

$$\underline{w}(i+1) = \underline{w}(i) - \eta(i) (\underline{w}(i)^T \underline{x}_n - y_n) \underline{x}_n, \quad n = (i \bmod N) + 1$$

(i: 0, 1, 2, ...; n: 1, 2, ..., N)

↪ LMS algorithm weight update

$(\underline{w}^T \underline{x}_n - y_n)$ is a signed error measure: $(\hat{y}(\underline{x}_n) - y_n)$; \underline{x}_n gives its direction.

Plot in weight space: $\epsilon_n = (\underline{w}^T \underline{x}_n - y_n) = 0 \Rightarrow 0 \text{ error for data point } \underline{x}_n.$

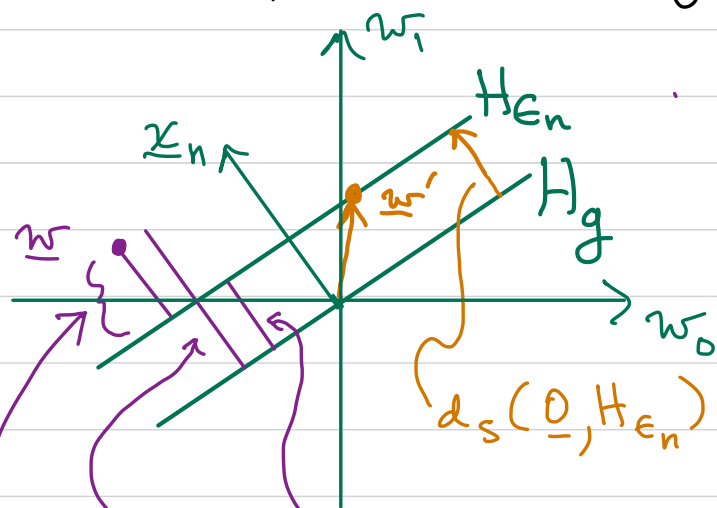
(S) $H_g: \underline{w}^T \underline{x}_n = 0$ is hyperplane:

$H_g \perp \underline{x}_n?$

H_g passes through origin?

What set of points is $\epsilon_n = 0$?

$H_{\epsilon_n}: \epsilon_n = \underline{w}^T \underline{x}_n - y_n = 0$ is hyperplane $\perp \underline{x}_n$, offset from origin by $\frac{y_n}{\|\underline{x}_n\|}$



Let \underline{w}' be any \underline{w} on H_{ϵ_n}
 Let $d_s(\underline{0}, H_{\epsilon_n}) \triangleq$ signed distance between $\underline{0}$ and H_{ϵ_n} .

$$= \underline{w}'^T \frac{\underline{x}_n}{\|\underline{x}_n\|} = \frac{y_n}{\|\underline{x}_n\|} = d_s(\underline{0}, H_{\epsilon_n})$$

and $d_s(\underline{0}, H_{\epsilon_n}) > 0$ iff $\underline{w}'^T \underline{x}_n > 0$

$$d_s(H_{\epsilon_n}, \underline{w}) = d_s(H_g, \underline{w}) - d_s(H_g, H_{\epsilon_n})$$

$$d_s(H_{\epsilon_n}, \underline{w}) = \underline{w}^T \frac{\underline{x}_n}{\|\underline{x}_n\|} - \frac{y_n}{\|\underline{x}_n\|} = \frac{\epsilon_n}{\|\underline{x}_n\|}$$

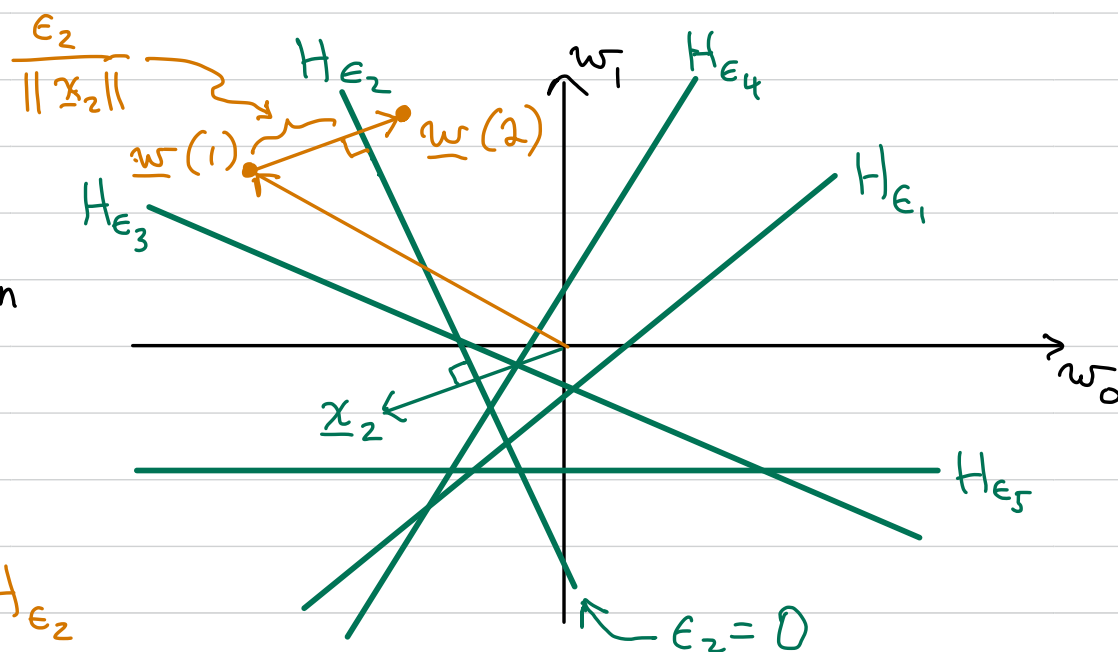
H_{ϵ_n} is target hyperplane for data point \underline{x}_n .

LMS example (1 weight update)

$N=5$

$$\underline{w}(i+1) = \underline{w}(i) - \eta(i) (\underline{w}(i)^T \underline{x}_n - y_n) \underline{x}_n$$

$$\underline{w}(2) = \underline{w}(1) - \eta(1) (\underbrace{\underline{w}(1)^T \underline{x}_2}_{\epsilon_2} - y_2) \underbrace{\underline{x}_2}_{\underline{x}_2 \perp H_{\epsilon_2}}$$



Note: if instead we had:

$$\underline{w}(2) - \underline{w}(1) = - \underbrace{\frac{(\underline{w}(1)^T \underline{x}_2 - y_2)}{\|\underline{x}_2\|}}_{\frac{\epsilon_2}{\|\underline{x}_2\|}} \frac{\underline{x}_2}{\|\underline{x}_2\|}$$

$$\frac{\epsilon_2}{\|\underline{x}_2\|} = d_s(H_{\epsilon_2}, \underline{w}(1))$$

then $\underline{w}(2) \Rightarrow 0$ error on \underline{x}_2

Mean-squared-error techniques: Classification (augmented notation) (2-class problems, reflected data points)

Can use above (MSE regression) techniques, except need a target output y_n .

Could choose: $y_n = z_n = \begin{cases} +1, & x_n \in S_1 \\ -1, & x_n \in S_2 \end{cases}$ ($\leftarrow x_n, y_n$ unreflected)

but can be more general: use b_n as a "target value", specified by the user.

$$(1) \therefore J(\underline{w}) = \frac{1}{N} \sum_{n=1}^N [g(z_n \underline{x}_n) - b_n]^2 = \frac{1}{N} \sum_{n=1}^N [\underline{w}^T z_n \underline{x}_n - b_n]^2, \quad b_n > 0 \quad \forall n.$$

Let H_B = decision boundary:

$$g(z\underline{x}) = \underline{w}^T z\underline{x} = 0$$

Let $b_n = b \quad \forall n$
for visualization

Let H_T = target
hyperplane
 $\underline{w}^T z_n \underline{x}_n - b = 0$

