

Chapter 3: The Linear model

A line is intuitively our first choice for a decision boundary.

Recap: perceptron learning algorithm:

We started by taking the hypothesis set that included all possible lines (actually hyperplanes). The algorithm then searched for a good line in the hypothesis set by iteratively correcting the errors made by the current candidate line, in an attempt to improve E_{in} .

P43 Example 2.1: linear model: small VC dimension \rightarrow generalize well from E_{in} to E_{out}

3 important problems: Classification, regression(回归, 复原), probability estimation

1. Linear model – linear classification:

a) Recap:

i. Generalize? ($E_{out} \approx E_{in}$):

Exercise 2.4

Exercise 2.4

a). d_{VC} : for $M_H = 2^N$, $d_{VC} = \text{the max } N$

$x \rightarrow d+1$ points

$$\begin{pmatrix} x_0 & \dots & x_d \\ \vdots & & \vdots \\ x_0 & \dots & x_d \end{pmatrix} \xrightarrow{\text{perceptron}} \begin{pmatrix} w_0 \\ \vdots \\ w_d \end{pmatrix} = \begin{pmatrix} y_0 \\ \vdots \\ y_d \end{pmatrix} \rightarrow \text{output}$$

$d+1$ dimension \mathbb{R} has $d+1$ 个

$d+1$ 个 ~~linearly~~ independent $\uparrow \mathbb{R}$ / full rank \times

$\therefore () \rightarrow$ non singular.

$$AX = B \quad X = A^{-1}B$$

\downarrow
non singular \therefore

\therefore for every B (every results y), we can find X

$$\begin{pmatrix} y_0 \\ \vdots \\ y_d \end{pmatrix} \text{ max } \geq d+1 \text{ 个结果}$$

$$X = \begin{pmatrix} w_0 \\ \vdots \\ w_d \end{pmatrix} \text{ the perceptron coefficient}$$

$$\therefore d_{VC} \geq d+1$$

$H^1(\mathbb{R})$, $M_H = 2^{d+1}$ 满足
shatter

For the

Exercise 2.4 (b).

$$\begin{pmatrix} x_0 \\ \vdots \\ x_d \end{pmatrix}_1 \dots \begin{pmatrix} x_0 \\ \vdots \\ x_d \end{pmatrix}_{d+2}$$

$$\therefore \begin{pmatrix} x_0 \\ \vdots \\ x_d \end{pmatrix}_{d+2} = \sum_{i=1}^{d+1} a_i \begin{pmatrix} x_0 \\ \vdots \\ x_d \end{pmatrix}_i$$

↳ linearly dependent

any $d+2$ vectors of $d+1$ ^{length} (dim) have to be linearly dependent

$$y_{d+2} = \sum_{i=1}^{d+1} a_i (x_0 \dots x_d)_i \begin{pmatrix} w_0 \\ \vdots \\ w_d \end{pmatrix}$$

For $i=1, \dots, d+1$.

$$\therefore y_1, \dots, y_{d+1} \rightarrow \text{determinant} \neq y_{d+2}.$$

\therefore Perceptron cannot shatter $d+2$.

$$d_{vc} < d+2. \quad d_{vc} \leq d+1.$$

$$\therefore d_{vc} = d+1$$

of perceptron.

dvc of linear model = $d+1$

ii. VC generalization bound

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \sqrt{\frac{8}{N} \ln \frac{4m_{\mathcal{H}}(2N)}{\delta}} \quad (2.12)$$

2.10 - bound on the growth function in terms of dvc

$$m_{\mathcal{H}}(N) \leq N^{d_{vc}} + 1. \quad (2.10)$$

we conclude that with high probability,

$$E_{\text{out}}(g) = E_{\text{in}}(g) + O\left(\sqrt{\frac{d}{N} \ln N}\right).$$

$$E_{out}(g) - E_{in}(g) \leq \sqrt{\frac{8}{N} \ln \frac{4m_H(2N)}{\delta}}$$

$$\leq \sqrt{\frac{8}{N} \ln \frac{4[(2N)^{d_{VC}} + 1]}{\delta}}$$

提出 Constant term

$$\sqrt{\frac{\ln N^{d_{VC}}}{N}} \quad d_{VC} = d+1$$

$$\sqrt{\frac{d \ln N}{N}}$$

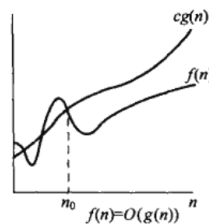
$O(\cdot)$ absolute value of this term is asymptotically smaller than a constant multiple of the argument

大O是我们在分析算法复杂度时最常用的一种表示法。

$f(x) = O(g(x))$ 表示的含义是 $f(x)$ 以 $g(x)$ 为上界

当函数的大小只有上界，没有明确下界的时候，则可以使用大O表示法，该渐进描述符一般用于描述算法的 **最坏复杂度**。

$f(x) = O(g(x))$ 正式的数学定义：存在正常数 c 、 n 、 n_0 ，当 $n > n_0$ 的时，任意的 $f(n)$ 符合 $0 \leq f(n) \leq c \cdot g(n)$ 。如下图所示



$\therefore N \uparrow$, $E_{out} \approx E_{in}$

b) E_{in} small? - Learn in sample a good hypothesis

Data is linearly separable: there is some hypothesis w^* with $E_{in}(w^*) = 0$.

PLA: increment based on one data point at a time; problem 1.3: finite step updating **any** misclassified data point $(\mathbf{x}(t), y(t))$, and update $\mathbf{w}(t)$ as follows:

$$\mathbf{w}(t+1) = \mathbf{w}(t) + y(t)\mathbf{x}(t).$$

PLA manages to search an infinite hypothesis set and output a linear separator in (provably) finite time.

Linearly separable data: from linearly separable target or not linearly separable target

But PLA converges & performance generalize well according to VC bound

2. Linear model – linear classification – non separable data

Cases:

1. Noisy or outlier;
2. Non-separable by a line

Both: PLA never terminates, jump from good to bad; No guarantee on E_{in}

Case1: tolerant small E_{in} ;

Case2: nonlinear transformation

a) Case 1:

Hypothesis with min E_{in}

solve the combinatorial optimization problem:

$$\min_{\mathbf{w} \in \mathbb{R}^{d+1}} \underbrace{\frac{1}{N} \sum_{n=1}^N [\text{sign}(\mathbf{w}^T \mathbf{x}_n) \neq y_n]}_{E_{in}(\mathbf{w})}.$$

NP-hard: no known efficient algorithm for it (NP 是指非确定性多项式 (non-deterministic polynomial, 缩写 NP)。所谓的非确定性是指，可用一定数量的运算去解决多项式时间内可解决的问题。)

∴ approximately minimising E_{in} -> pocket algorithm: slower than PLA; no guarantee for convergence speed

The pocket algorithm:

- 1: Set the pocket weight vector $\hat{\mathbf{w}}$ to $\mathbf{w}(0)$ of PLA.
- 2: **for** $t = 0, \dots, T - 1$ **do**
- 3: Run PLA for one update to obtain $\mathbf{w}(t + 1)$.
- 4: Evaluate $E_{in}(\mathbf{w}(t + 1))$.
- 5: If $\mathbf{w}(t + 1)$ is better than $\hat{\mathbf{w}}$ in terms of E_{in} , set $\hat{\mathbf{w}}$ to $\mathbf{w}(t + 1)$.
- 6: Return $\hat{\mathbf{w}}$.

b) Exercise 3.2 code explanation

[Learning-From-Data-A-Short-Course/Solutions to Chapter 3 The Linear Model.ipynb at master · niuers/Learning-From-Data-A-Short-Course \(github.com\)](#)

c) Example 3.1 (Handwritten digit recognition)

Raw input: preprocessed 16*16 pixel images

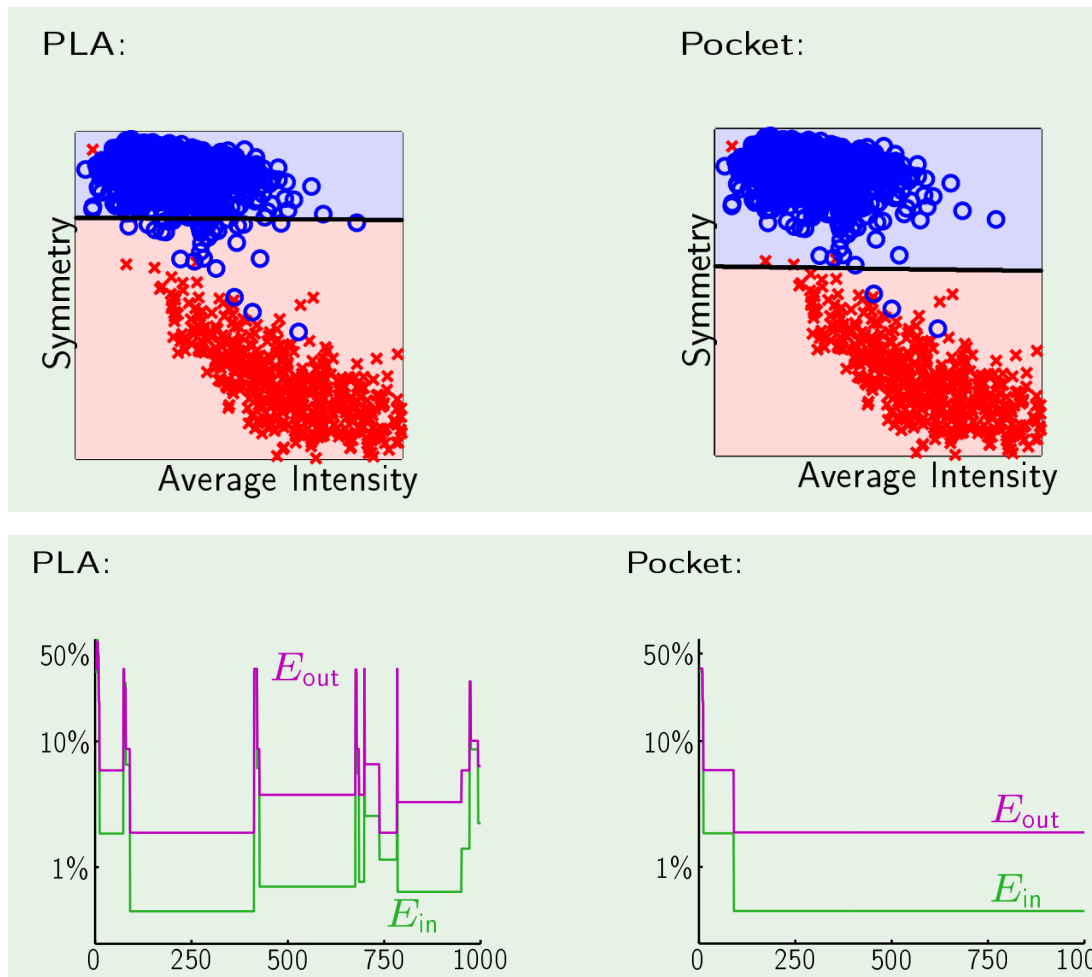
Decomposition: multiclass -> binary (classify {1,5} first)

Processed input: average intensity & symmetry (asymmetry: absolute difference between an image and its FLIPPED versions; symmetry is the negation)

Situation: roughly separable by a line; poorly written digits prevent a perfect linear separation.

PLA: not stop updating; behavior unstable

Pocket algorithm: E_{in} decreasing



3. Linear model – linear regression

- a) A useful model applied to real-valued target function

Example: Set a credit limit for approved customers (assume linear separable)

Input \mathbf{x}_n : customer info (eg. Historical records)

Output: Noisy target formalized as $P(y|\mathbf{x})$ – expert different views and variations within them.

$P(y|\mathbf{x}) \rightarrow (\mathbf{x}_n, y_n) \rightarrow g = \text{argmin}(\text{Error between } g(\mathbf{x}) \text{ and } y \text{ from } P(y|\mathbf{x}))$

- b) The Algorithm

Minimizing the square error

$$E_{\text{out}}(h) = \mathbb{E} \left[(h(\mathbf{x}) - y)^2 \right]$$

The Expectation uses joint probability distribution $P(\mathbf{x}, y)$ – unknown

Resort to E_{in}

$$E_{\text{in}}(h) = \frac{1}{N} \sum_{n=1}^N (h(\mathbf{x}_n) - y_n)^2$$

data matrix $X \in \mathbb{R}^{N \times (d+1)}$

$$\begin{array}{c}
 \begin{array}{c} N \\ \text{data} \\ \text{points} \end{array} \left\{ \begin{array}{c} (x_0 \dots x_d) \\ \vdots \\ (x_0 \dots x_d) \end{array} \right\} \xrightarrow{\text{1st data point}} \begin{pmatrix} w_0 \\ \vdots \\ w_d \end{pmatrix} = \begin{pmatrix} h_0 \\ \vdots \\ h_d \end{pmatrix} \\
 \underbrace{\hspace{10em}}_{\substack{d+1 \\ x \in \{x_i\} \times \mathbb{R}^d \\ x_0 = 1}} \hspace{10em} h = XW
 \end{array}$$

$$\begin{aligned}
 \therefore E_{in}(h) &= \frac{1}{N} \sum_{n=1}^N (h(x_n) - y_n)^2 \\
 &= \frac{1}{N} \sum_{n=1}^N (W^T x_n - y_n)^2 \\
 &= \frac{1}{N} \|XW - y\|^2 \quad \rightarrow \begin{pmatrix} h_0 - y_0 \\ \vdots \\ h_N - y_N \end{pmatrix} \text{ is Euclidean norm}^2 \\
 &\quad \downarrow \\
 &= \frac{1}{N} (W^T X^T X W - 2 W^T X^T y + y^T y)
 \end{aligned}$$

The linear regression algorithm is derived by minimizing $E_{in}(w)$ over all possible $w \in \mathbb{R}^{d+1}$

$$w_{lin} = \underset{w \in \mathbb{R}^{d+1}}{\operatorname{argmin}} E_{in}(w).$$

find the w that minimizes $E_{in}(w)$ by requiring that the **gradient of E_{in} with respect to w is the zero vector**

$$\nabla E_{in}(w) = 0$$

$$[\nabla E_{in}(w)]_i = \frac{\partial}{\partial w_i} E_{in}(w).$$

$$\begin{aligned}
 \nabla_w E_{in}(w) &= \begin{pmatrix} \frac{\partial}{\partial w_0} E_{in}(w) \\ \frac{\partial}{\partial w_1} E_{in}(w) \\ \vdots \\ \frac{\partial}{\partial w_d} E_{in}(w) \end{pmatrix} = \frac{1}{N} (X^T X + X^T X) W - 2 X^T y \\
 &= \frac{2}{N} (X^T X W - X^T y) \\
 &\quad (d+1) \times 1
 \end{aligned}$$

Rules: $\nabla_{\mathbf{w}}(\mathbf{w}^T \mathbf{A} \mathbf{w}) = (\mathbf{A} + \mathbf{A}^T) \mathbf{w}, \quad \nabla_{\mathbf{w}}(\mathbf{w}^T \mathbf{b}) = \mathbf{b}.$

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}.$$

If $\mathbf{X}^T \mathbf{X}$ is invertible, $\mathbf{w} = \mathbf{X}^\dagger \mathbf{y}$ where $\mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is the *pseudo-inverse*
 If $\mathbf{X}^T \mathbf{X}$ not invertible, pseudo-inverse can be defined, but solution not unique

Problem 3.15.

(a) $\text{rank}(\mathbf{X}) = p \quad \mathbf{X} \in \mathbb{R}^{N \times (d+1)}$

$\mathbf{X}^T \mathbf{X} \rightarrow (d+1) \times (d+1)$ no inverse
 $(d+1) \times N \quad N \times (d+1) \quad \therefore \text{rank}(\mathbf{X}^T \mathbf{X}) < d+1$

\downarrow
 $p \leq \min(d+1, N)$
 \downarrow
 $p \leq d+1$
 by the property of rank
 $\text{rank}(\mathbf{X}) = \text{rank}(\mathbf{X}^T \mathbf{X}) = p$
 $\therefore p < d+1$

$\mathbf{X} = \mathbf{U} \mathbf{\Gamma} \mathbf{V}^T$
 $(N \times p) \quad p \times (d+1)$

(b) $\mathbf{w}_{\text{lin}} = \mathbf{V} \mathbf{\Gamma}^{-1} \mathbf{U}^T \mathbf{y}$

$$\begin{aligned} \mathbf{X}^T \mathbf{X} \mathbf{w}_{\text{lin}} &= \mathbf{X}^T \mathbf{X} \mathbf{V} \mathbf{\Gamma}^{-1} \mathbf{U}^T \mathbf{y} \\ &= (\mathbf{U} \mathbf{\Gamma} \mathbf{V}^T)^T (\mathbf{U} \mathbf{\Gamma} \mathbf{V}^T) \mathbf{V} \mathbf{\Gamma}^{-1} \mathbf{U}^T \mathbf{y} \\ &= \mathbf{V} \mathbf{\Gamma}^T \mathbf{U}^T \mathbf{U} \mathbf{\Gamma} \mathbf{V}^T \mathbf{V} \mathbf{\Gamma}^{-1} \mathbf{U}^T \mathbf{y} \\ &= \mathbf{V} \mathbf{\Gamma}^2 \mathbf{\Gamma}^{-1} \mathbf{U}^T \mathbf{y} = \mathbf{V} \mathbf{\Gamma} \mathbf{U}^T \mathbf{y} = \mathbf{X}^T \mathbf{y} \end{aligned}$$

c.c)

$$X^T X W = X^T Y$$

$$X^T X W_{lin} = X^T Y$$

$$X^T (W - W_{lin}) = 0$$

$$V U^T U V^T (W - W_{lin}) = 0$$

$$V \Gamma^2 V^T (W - W_{lin}) = 0$$

$$\underbrace{V^{-2} V^T V \Gamma^2 V^T}_{\neq 0} (W - W_{lin}) = \underbrace{V^T (W - W_{lin})}_{=0} = 0$$

$\Gamma V^T \neq 0$ $\Gamma \neq 0$ $\Gamma^2 \neq 0$
 $\Gamma^T V^T \neq 0$ $\Gamma^2 V^T \neq 0$

let $W - W_{lin} = \delta$
 $W_{lin} = W - \delta$ $W_{lin} \perp \delta$
 orthogonal i.e. $\delta \perp W_{lin}$

$$\begin{aligned}
 W_{lin}^T \delta &= W_{lin}^T (W - W_{lin}) \\
 &= Y^T U^T V^T (W - W_{lin}) \\
 &= 0
 \end{aligned}$$

$$W_{lin} = V \Gamma^+ U^T Y$$

$\rightarrow X^T X$ $\rightarrow X^T Y$
 $A W_{lin} = Z$
 not full rank $p < d+1$
 $A W_{lin} + A A' = A (W_{lin} + A) = Z$
 for $A \perp A'$

\rightarrow not full rank
 \rightarrow not full rank

$= 0$
 $\|W\|^2 = \|W_{lin} + \delta\|^2 \geq 0$
 $= W_{lin}^T W_{lin} + 2 W_{lin}^T \delta + \delta^T \delta$
 $= \|W_{lin}\|^2 + \|\delta\|^2$
 $\geq \|W_{lin}\|^2$ For the Minimum norm

$X \in \mathbb{R}^{N \times (d+1)}$, N is much bigger than $d+1$, very likely there are $d+1$ linearly independent points

So $\text{rank}(X^T X) = d+1$, invertible

Linear regression algorithm:

Linear regression algorithm:

- 1: Construct the matrix X and the vector y from the data set $(x_1, y_1), \dots, (x_N, y_N)$, where each x includes the $x_0 = 1$ bias coordinate, as follows

$$\underbrace{X = \begin{bmatrix} -x_1^T \\ -x_2^T \\ \vdots \\ -x_N^T \end{bmatrix}}_{\text{input data matrix}}, \quad \underbrace{y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\text{target vector}}$$

- 2: Compute the pseudo-inverse X^\dagger of the matrix X . If $X^T X$ is invertible,

$$X^\dagger = (X^T X)^{-1} X^T$$

- 3: Return $w_{lin} = X^\dagger y$.

w_{lin} has a decent E_{out} \rightarrow learning has occurred.

an analytic formula for learning \rightarrow widely used.

Estimated output: $\hat{y} = X w_{lin}$ $\rightarrow \hat{y} = X(X^T X)^{-1} X^T y$.

linear transformation of the actual y through $H = X(X^T X)^{-1} X^T$. (H 'puts a hat' on y)

c) Exercise 3.3

$$\begin{aligned} (a) \quad H^T &= X \left[(X^T X)^{-1} \right]^T X^T \\ &= X \left[(X^T X)^{-1} \right] X^T = H \end{aligned}$$

$$\begin{aligned} (b) \quad H \circ H &= X (X^T X)^{-1} X^T X (X^T X)^{-1} X^T \\ &= X (X^T X)^{-1} X^T = H \end{aligned}$$

$$\therefore H^k = H$$

$$\begin{aligned} (c) \quad (I - H)(I - H) &= I - H I - I H + H^2 \\ &= I - H - H + H \\ &= I - H \\ (I - H)^k &= I - H \end{aligned}$$

$$\begin{aligned} (d) \quad \text{trace}(H) &= \text{trace} \left(X (X^T X)^{-1} X^T \right) = \text{trace} \left(X^T X (X^T X)^{-1} \right) \\ &\quad \begin{matrix} N \times N & N \times (d+1) & (d+1) \times N & (d+1) \times N & N \times (d+1) \\ & \downarrow & & \downarrow & \\ & N \times N & & (d+1)(d+1) & \end{matrix} \\ &= \text{trace}(I_{d+1}) = d+1 \end{aligned}$$

d) Generalization Issues

A regression version of the VC generalization bound

$$E_{out}(g) = E_{in}(g) + O\left(\frac{d}{N}\right)$$

Proof: Exercise 3.4

$$E_{\mathcal{D}}[E_{in}(w_{lin})] = \sigma^2 \left(1 - \frac{d+1}{N} \right)$$

$$E_{\mathcal{D}, \epsilon'}[E_{test}(w_{lin})] = \sigma^2 \left(1 + \frac{d+1}{N} \right)$$

Exercise 3.4.

$$\underline{y} = \underline{W}^* \underline{X} + \underline{\varepsilon} \quad \underline{y} = \underline{X} \underline{W}^* + \underline{\varepsilon}$$

(a)

Linear regression $\rightarrow \hat{\underline{y}} = H \underline{y}$

没有假设 \underline{y} 是什么样的.

$$= H \cancel{\underline{W}^* \underline{X}}^T (\underline{X} \underline{W}^* + \underline{\varepsilon})$$

$$\hat{\underline{y}} = \underline{X} \underline{W}_{lin} = H \underline{y}$$

$$= H \underline{X} \underline{W}^* + H \underline{\varepsilon}$$

$$= \underline{X} (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{X} \underline{W}^* + H \underline{\varepsilon}$$

$$= \underline{X} \underline{W}^* + H \underline{\varepsilon}$$

(b)

$$\begin{aligned} \hat{\underline{y}} - \underline{y} &= \underline{X} \underline{W}^* + H \underline{\varepsilon} - \underline{X} \underline{W}^* - \underline{\varepsilon} = H \underline{\varepsilon} - \underline{\varepsilon} \\ &= \underline{(H - I)} \underline{\varepsilon} \end{aligned}$$

(c)

$$\begin{aligned} E_{in}(\underline{W}_{lin}) &= \frac{1}{N} \sum_{n=1}^{N} (\underline{W}_{lin}^T \underline{x}_n - y_n)^2 \\ &= \frac{1}{N} \|\underline{X} \underline{W}_{lin} - \underline{y}\|^2 = \frac{1}{N} \|\hat{\underline{y}} - \underline{y}\|^2 \\ &= \frac{1}{N} \|\underline{(H - I)} \underline{\varepsilon}\|^2 = \frac{1}{N} \underline{\varepsilon}^T \underline{(H - I)}^T \underline{(H - I)} \underline{\varepsilon} \end{aligned}$$

Exercise 3.3

$$\begin{aligned} \text{(c)} &= \frac{1}{N} \underline{\varepsilon}^T (\underline{H}^T - \underline{I}) (\underline{H} - \underline{I}) \underline{\varepsilon} \\ &= \frac{1}{N} \underline{\varepsilon}^T (\underline{I} - \underline{H}) \underline{\varepsilon} \end{aligned}$$

For the Visnapier

$$H = X(X^T X)^{-1} X^T$$

$$(d) E_D[E_{in}(w_{lin})] = E_D\left[\frac{1}{N} \sum_{i=1}^N \epsilon_i^T (I - H) \epsilon_i\right]$$

$1 \times N$
 $(\epsilon_1 \dots \epsilon_N)$
 $(1 \times N)$
 (\dots)

$N \times N$
 $\begin{pmatrix} H_{11} & \dots & H_{1N} \\ \vdots & \ddots & \vdots \\ H_{M1} & \dots & H_{MN} \end{pmatrix}$
 $(N \times N)$

$$= \frac{1}{N} [E_D(\sum_{i=1}^N \epsilon_i^T \epsilon_i) - E_D(\sum_{i=1}^N \epsilon_i^T H \epsilon_i)]$$

$$= \frac{1}{N} [E_D(\sum_{n=1}^N \epsilon_n^2) - E_D(\sum_{i=1}^N \sum_{j=1}^N \epsilon_i h_{ij} \epsilon_j)]$$

$$= \frac{1}{N} [E_D(\sum_{n=1}^N \epsilon_n^2) - \sum_{i=1}^N \sum_{j=1}^N E_D(\epsilon_i h_{ij} \epsilon_j)]$$

$$= \frac{1}{N} [N\sigma^2 - \sum_{i=1}^N E_D(\epsilon_i^2 h_{ii})]$$

$$= \frac{1}{N} [N\sigma^2 - \sum_{i=1}^N E_D(\epsilon_i^2) E_D(h_{ii})]$$

$$= \frac{1}{N} [N\sigma^2 - \sigma^2 \sum_{i=1}^N E_D(h_{ii})]$$

$$= \frac{1}{N} [N\sigma^2 - \sigma^2 E_D(\text{trace}(H))]$$

$$= \frac{1}{N} [N\sigma^2 - \sigma^2 \text{trace}(H)]$$

$$E_D(\epsilon_i h_{ij} \epsilon_j) = E_D(\epsilon_i) E_D(h_{ij} \epsilon_j) = 0$$

\swarrow
independent
if $i \neq j$

\parallel
0

$$= \frac{1}{N} [N\sigma^2 - \sigma^2 (d+1)]$$

$$= \sigma^2 - \frac{\sigma^2 (d+1)}{N}$$

$$= \sigma^2 \left(1 - \frac{d+1}{N}\right)$$

(e)

$$E_{D, \epsilon'} [E_{\text{test}}(w_{\text{lin}})] = E_{D, \epsilon'} \left[\frac{1}{N} \|\hat{y} - y'\|^2 \right]$$

$$= E_{D, \epsilon'} \left[\frac{1}{N} \|\underbrace{H y - y'}_{\downarrow} \|^2 \right]$$

$$\downarrow XW^* + H\epsilon - (XW^* + H\epsilon')$$

$$\downarrow H\epsilon - \epsilon'$$

y' from D_{test}

y from D_{train}

Share same X_n

but different ϵ

$$= E_{D, \epsilon'} \left[\frac{1}{N} \|H\epsilon - \epsilon'\|^2 \right]$$

$$= E_{D, \epsilon'} \left[\frac{1}{N} (H\epsilon - \epsilon')^T (H\epsilon - \epsilon') \right]$$

$$= E_{D, \epsilon'} \left[\frac{1}{N} (\epsilon^T H^T - \epsilon'^T) (H\epsilon - \epsilon') \right]$$

$$= E_{D, \epsilon'} \left[\frac{1}{N} (\epsilon^T H H \epsilon - \epsilon'^T H \epsilon - \epsilon^T H \epsilon' + \epsilon'^T \epsilon') \right]$$

$$= E_{D, \epsilon'} \left[\frac{1}{N} (\underbrace{\epsilon^T H \epsilon}_{\text{wavy}} - \underbrace{\epsilon'^T H \epsilon}_{\text{wavy}} - \underbrace{\epsilon^T H \epsilon'}_{\text{wavy}} + \underbrace{\epsilon'^T \epsilon'}_{\text{wavy}}) \right]$$

$$= E_{D, \epsilon'} \left[\frac{1}{N} (\epsilon^T H \epsilon + \epsilon'^T \epsilon') \right]$$

$$= \frac{1}{N} \left\{ \sum_{i=1}^N \sum_{j=1}^N E_{D, \epsilon'} [\epsilon_i h_{ij} \epsilon_j] + \sum_{i=1}^N E_{D, \epsilon'} [\epsilon_i^2] \right\}$$

$$= \frac{1}{N} \left\{ \sigma^2 (d+1) + N \sigma^2 \right\} = \sigma^2 \left(1 + \frac{d+1}{N} \right)$$

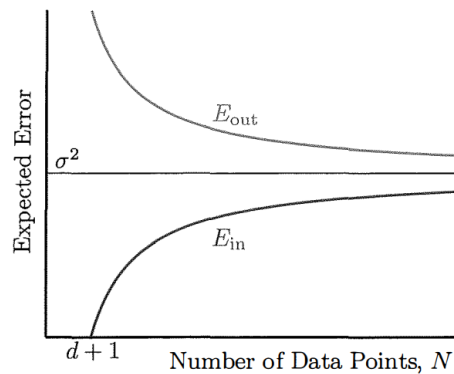


Figure 3.4: The learning curve for linear regression.

E_{in} : learned linear fit has degree of freedom $d+1$ ($\therefore E_{in}$ 可以 = 0 when $d+1$ points) \rightarrow eat into in-sample noise (\therefore cannot distinguish)
 E_{out} : best with σ^2 , additional error due to fitting in the in-sample noise

Recap:

Linear model - 'draw a line' between 2 categories

1) linear classification

2) linear regression

linear classification - perceptron learning algorithm (PLA)

as long as data linearly separable, $E_{in} = 0$

what is non linearly separable?

-two cases:

1. linearly separable after discarding a few examples

2. not linearly separable, but separable by a more sophisticated curve

case 1 - The pocket algorithm

- put the best in the pocket

- get the best of the random selected lines

An example - handwritten digit recognition

1. multiclass decomposition: recognise 1 / 5

2. 256 pixels \rightarrow 2 features (symmetry and intensity)

test: use pocket, E_{out} lower than human error

2) linear regression:

target function real-valued, not y generated from $p(y|x)$

what is the E_{in} and E_{out} ?

The hypothesis $h(x)$ is still a linear combination of the components of x

Now we want to minimise E_{in}

$E_{in}(w)$ differentiable \rightarrow differentiation wrt w ,

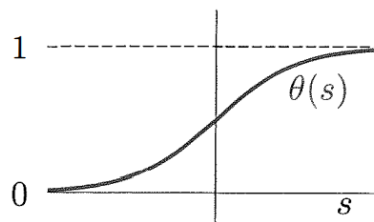
4. The linear model – Logistic Regression

Output a probability between 0 and 1: real (like linear regression) and bounded (like linear classification). Eg. Occurrence of heart attacks

a) The linear model – Logistic Regression (model) - Predicting a probability

$$h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x}),$$

where θ is the so-called *logistic* function $\theta(s) = \frac{e^s}{1+e^s}$ whose output is between 0 and 1.



Output: probability of a binary event – uncertain classification

i. Exercise 3.5

Let $s \rightarrow \infty$, we have $e^{-2s} \rightarrow 0$, thus $\tanh(s) \rightarrow 1$. Similarly, when $s \rightarrow -\infty$, $\tanh(s) \rightarrow -1$.

It's also easy to see that $\tanh(s) < 1$ and $\tanh(s) > -1$. So -1 and 1 are hard thresholds for $\tanh(s)$ when $|s|$ is large.

When $|s|$ is small, consider Taylor expansion of $e^s = 1 + s + \frac{s^2}{2} + O(s^3)$, then $e^{-s} = 1 - s + \frac{s^2}{2} - O(s^3)$, we have $\tanh(s) = \frac{e^s - e^{-s}}{e^s + e^{-s}} = \frac{2s + 2O(s^3)}{2 + s^2} \approx s$.

So $\tanh(s)$ is approximately linear when $|s|$ is small. There's no threshold in this case.

Exercise 3.5

$$(a) \quad \tanh(s) = \frac{e^s - e^{-s}}{e^s + e^{-s}} \quad \theta(s) = \frac{e^s}{1+e^s}$$

$$\tanh(s)+1 = \frac{e^s - e^{-s} + e^s + e^{-s}}{e^s + e^{-s}} = \frac{2e^s}{e^s + e^{-s}} = \frac{2e^{2s}}{e^{2s} + 1}$$

$$\tanh(s)+1 = 2\theta(2s)$$

$$\tanh(s) = 2\theta(2s) - 1$$

$$(b) \quad s \rightarrow \infty \quad \tanh(s) \rightarrow \frac{e^s - 0}{e^s + 0} = 1 \quad \Delta$$

$$s \rightarrow -\infty \quad \tanh(s) \rightarrow \frac{0 - e^{-s}}{0 + e^{-s}} = -1 \quad \Delta$$

$$\tanh(s) \rightarrow 1$$

$$e^s = 1 + s + \frac{s^2}{2} + O(s^3)$$

$$e^{-s} = 1 - s + \frac{s^2}{2} - O(s^3)$$

$$\tanh = \frac{2s + 2O(s^3)}{2 + s^2} \approx s \quad \text{linear}$$

$$P(y | \mathbf{x}) = \begin{cases} f(\mathbf{x}) & \text{for } y = +1; \\ 1 - f(\mathbf{x}) & \text{for } y = -1. \end{cases} \quad f(\mathbf{x}) = \mathbb{P}[y = +1 | \mathbf{x}].$$

The data is in fact generated by a noisy target $P(y | \mathbf{x})$,

Error measure: how close h to f in terms of the noisy ± 1 examples.

推导出 error measure:

$$P(y | \mathbf{x}) = \begin{cases} h(\mathbf{x}) & \text{for } y = +1; \\ 1 - h(\mathbf{x}) & \text{for } y = -1. \end{cases}$$

We substitute for $h(\mathbf{x})$ by its value $\theta(\mathbf{w}^T \mathbf{x})$, and use the fact that $1 - \theta(s) = \theta(-s)$ (easy to verify) to get

$$P(y | \mathbf{x}) = \theta(y \mathbf{w}^T \mathbf{x}). \quad (3.8)$$

the probability of :

given \mathbf{x} , get the actual y (generated by the noisy target $P(y|\mathbf{x})$)

data points in a set: independent: the probability of getting all the actual y from all the given \mathbf{x} :

$$\prod_{n=1}^N P(y_n | \mathbf{x}_n).$$

$$\prod_{n=1}^N P(y_n | \mathbf{x}_n).$$

Method of maximum likelihood: maximize

Equivalent to minimize:

$$-\frac{1}{N} \ln \left(\prod_{n=1}^N P(y_n | \mathbf{x}_n) \right) = \frac{1}{N} \sum_{n=1}^N \ln \left(\frac{1}{P(y_n | \mathbf{x}_n)} \right) = \frac{1}{N} \sum_{n=1}^N \ln \left(\frac{1}{\theta(y_n \mathbf{w}^T \mathbf{x}_n)} \right)$$

minimize the "error" with respect to the weight vector \mathbf{w}

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln \left(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n} \right).$$

Pointwise error measure

$$e(h(\mathbf{x}_n), y_n) = \ln(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n})$$

error is small when actual output (y_n) and the estimation ($\mathbf{w}^T \mathbf{x}_n$) are consistent (same sign)

error measure encourages \mathbf{w} to 'classify' each \mathbf{x}_n correctly

ii. Exercise 3.6 (prove error measure)

Exercise 3.6

1. (a) The probability to get y_n is $P(y_n|x_n)$, by maximum likelihood method, we need maximize the likelihood, $\prod_{n=1}^N P(y_n|x_n)$, this is equivalent to maximize the logarithm of it: $\sum_{n=1}^N \ln(P(y_n|x_n))$, or minimize the negative of it: $-\sum_{n=1}^N \ln(P(y_n|x_n))$

When $y_n = +1$, $P(y_n|x_n) = h(x_n)$, and when $y_n = -1$, $P(y_n|x_n) = 1 - h(x_n)$, separate the cases for $y_n = 1$ and $y_n = -1$, we have:

$$\begin{aligned} E_{in}(w) &= -\sum_{n=1}^N \ln(P(y_n|x_n)) \\ &= -\sum_{n=1}^N I(y_n = +1) \ln h(x_n) + I(y_n = -1) \ln(1 - h(x_n)) \\ &= \sum_{n=1}^N I(y_n = +1) \ln \frac{1}{h(x_n)} + I(y_n = -1) \ln \frac{1}{(1 - h(x_n))} \end{aligned}$$

1. (b) For $h(x) = \theta(w^T x) = \frac{e^{w^T x}}{1 + e^{w^T x}}$, we have $\ln \frac{1}{h(x_n)} = \ln(1 + e^{-w^T x_n})$ and $\ln \frac{1}{(1 - h(x_n))} = \ln(1 + e^{w^T x_n})$. Combine them together we have

$$\begin{aligned} E_{in}(w) &= \sum_{n=1}^N I(y_n = +1) \ln(1 + e^{-w^T x_n}) + I(y_n = -1) \ln(1 + e^{w^T x_n}) \\ &= \sum_{n=1}^N \ln(1 + e^{-y_n w^T x_n}) \end{aligned}$$

Which is equivalent of minimizing the one in equation (3.9).

这种 error measure 是 cross-entropy error measure

$$p \log \frac{1}{q} + (1 - p) \log \frac{1}{1 - q}.$$

iii. Exercise 3.7

Exercise 3.7

$$E_{in}(W) = \frac{1}{N} \sum_{n=1}^N \ln(1 + e^{-y_n W^T x_n})$$

$$\nabla E_{in}(W) = \frac{1}{N} \sum_{n=1}^N \frac{-y_n x_n e^{-y_n W^T x_n}}{1 + e^{-y_n W^T x_n}} \rightarrow \begin{pmatrix} \frac{\partial}{\partial w_1} \\ \vdots \end{pmatrix} \text{ column vector.}$$

$$= -\frac{1}{N} \sum_{n=1}^N y_n x_n \frac{1}{1 + e^{y_n W^T x_n}} \quad \theta(s) = \frac{e^s}{1 + e^s}$$

$$= -\frac{1}{N} \sum_{n=1}^N y_n x_n \theta(-y_n W^T x_n) \quad 1 - \theta(s) = \theta(-s)$$

\downarrow
yes



if misclassified.

$$y_n W^T x_n < 0$$

$\theta(-y_n W^T x_n) > 0.5 \rightarrow \text{Contribute more to } \nabla E_{in}(W)$

correctly classified $\theta(-y_n W^T x_n) < 0.5$

Iteratively set the gradient of the error measure to zero – gradient descent

b) The linear model – Logistic Regression (model) – Gradient descent

Train learning models with smooth error measures – minimizing a twice-differentiable function

not necessarily rest in the global minima – end up at a local minimum (depend on starting weights)

but logistic regression with the cross-entropy error – only one (unique) global minimum – $E_{in}(w)$ convex

i. take a step in the direction of steepest descent 推导:

Taylor expansion:

$$\begin{aligned} \Delta E_{in} &= E_{in}(w(0) + \eta \hat{v}) - E_{in}(w(0)) \\ &= \eta \nabla E_{in}(w(0))^T \hat{v} + O(\eta^2) \\ &\geq \eta \|\nabla E_{in}(w(0))\|, \end{aligned}$$

For largest negative $\Delta E_{in} = \eta \|\nabla E_{in}(w(0))\|$ (ignore square term)

$$\hat{v} = -\frac{\nabla E_{in}(w(0))}{\|\nabla E_{in}(w(0))\|} \text{ leads to the largest decrease in } E_{in}$$

Exercise 3.8

\hat{v} is the direction which gives largest decrease in E_{in} only holds for small η , that's because when η is large, we can't ignore the squared term and smaller terms in the Taylor expansion. The lower bound can't be achieved.

ii. Learning rate η :

Fixed too small step size: insufficient when far from local minimum

Too large step size: bouncing around, possibly increasing E_{in}

Variable η : large step when far from the minimum

$$\eta_t = \eta \|\nabla E_{in}\|$$

$$\begin{aligned} \mathbf{w}(t+1) &= \mathbf{w}(t) - \underbrace{\eta \|\nabla E_{in}\|}_{\text{Variable } \eta_t} \frac{\nabla E_{in}(\mathbf{w}(t))}{\|\nabla E_{in}(\mathbf{w}(t))\|} \\ &= \mathbf{w}(t) - \underbrace{\eta}_{\text{Fixed } \eta_t} \nabla E_{in}(\mathbf{w}(t)) \end{aligned}$$

Fixed learning rate gradient descent:

- 1: Initialize the weights at time step $t = 0$ to $\mathbf{w}(0)$.
- 2: **for** $t = 0, 1, 2, \dots$ **do**
- 3: Compute the gradient $\mathbf{g}_t = \nabla E_{in}(\mathbf{w}(t))$.
- 4: Set the direction to move, $\mathbf{v}_t = -\mathbf{g}_t$.
- 5: Update the weights: $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \mathbf{v}_t$.
- 6: Iterate to the next step until it is time to stop.
- 7: Return the final weights.

$$E_{in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln \left(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n} \right).$$

Logistic regression algorithm:

- 1: Initialize the weights at time step $t = 0$ to $\mathbf{w}(0)$.
- 2: **for** $t = 0, 1, 2, \dots$ **do**
- 3: Compute the gradient

$$\mathbf{g}_t = -\frac{1}{N} \sum_{n=1}^N \frac{y_n \mathbf{x}_n}{1 + e^{y_n \mathbf{w}^T(t) \mathbf{x}_n}}.$$

- 4: Set the direction to move, $\mathbf{v}_t = -\mathbf{g}_t$.
- 5: Update the weights: $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \mathbf{v}_t$.
- 6: Iterate to the next step until it is time to stop.
- 7: Return the final weights \mathbf{w} .

iii. Initialization and termination

Initialization:

$\mathbf{w}(0) = \mathbf{0}$ - stuck on a perfectly symmetric hilltop.

safer: randomly - Choosing each weight independently from a Normal distribution with zero mean and small variance

Termination:

1. set an upper bound on the number of iterations (no guarantee on the quality)

2. (\because at minimum, gradient = 0) gradient below a certain threshold (if flat region, prematurely stop)
3. Require error itself small

Combination: a maximum number of iterations, marginal error improvement, coupled with small value for the error itself