Principles of Machine Learning

(Special and a second second

EFLECTIVE.

IFS.

HARY.

ULTIES.

NG,KNOWING,

Fine subst

Intuition

Cuncism.

Comparison.

Acces intim

With

Po-Chen Wu

Media IC and System Lab Graduate Institute of Electronics Engineering National Taiwan University

Outline

- Introduction to Machine Learning
- Theory of Generalization
- Learning Algorithm
- Hazard of Overfitting
- Blending and Bagging

Outline

- Introduction to Machine Learning
- Theory of Generalization
- Learning Algorithm
- Hazard of Overfitting
- Blending and Bagging

Mars One Project



- A one-way ticket to Mars.
- There is a total of 2,782 applicants.
- The application consists of applicant's
 - General information
 - Motivational letter
 - Résumé
 - Video

Admission Ticket Approval

Applicant Information	
Age	37 years
Gender	Male
Occupation	Professor
Annual Salary	NTD 2,000,000
Year in Job	11 Years
Current Debt	NTD 110,000

Unknown target function to be learned:
 "Should we approve the admission ticket or not?"

Formalize the Learning Problem

Basic Notations

- input: $\mathbf{x} \in \mathcal{X}$ (volunteer application)
- output: $y \in \mathcal{Y}$ (good/bad after approving ticket)
- unknown target function to be learned $f: \mathcal{X} \rightarrow \mathcal{Y}$ (ideal ticket approval formula)
- − data ⇔ training examples: $\mathcal{D}: \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\} \text{ (historical records)}$
- hypothesis \Leftrightarrow skill with hopefully good performance: $g: \mathcal{X} \rightarrow \mathcal{Y}$ ('learned' ticket approval formula)

$$\{(\mathbf{x}_n, y_n)\} \text{ from } f \longrightarrow \mathsf{ML} \longrightarrow g$$

Learning Flow for Ticker Approval



- target *f* unknown
 (i.e. no programmable definition)
- hypothesis *g* hopefully ≈ *f* but possibly different from *f* (perfection 'impossible when *f* unknown'

What does *g* look like?

The Learning Model



Media IC & System Lab

Po-Chen Wu (吳柏辰)

Practical Definition of Machine Learning



Machine Learning:

use data to compute hypothesis g that approximates target f

Media IC & System Lab

Outline

- Introduction to Machine Learning
- Theory of Generalization
- Learning Algorithm
- Hazard of Overfitting
- Blending and Bagging

Sex Ratio of EE Students



Population





- Population
 - $\quad \text{girl proportion} = \mu$

μ: unknown

- boy proportion = 1μ
- Sample
 - $\quad \text{girl fraction} = \nu$

ν: known

- boy fraction = $1 - \nu$

Does in-sample ν say anything about out-of-sample μ ?

Hoeffding's Inequality



• In big sample (*N* large), ν is probably close to μ (within ϵ)

 $\mathbb{P}[|\nu - \mu| > \epsilon] \le 2\exp(-2\epsilon^2 N)$ Hoeffding's Inequality

the statement ' $\nu = \mu$ ' is probably approximately correct (PAC)

Connection to Learning

EE

- unknown girl prob. μ
- boy 🛉
- girl 🛉
- size-N sample from EE of i.i.d. students.

Learning

- fixed hypothesis $h(\mathbf{x}) \stackrel{?}{=} f(\mathbf{x})$
- $\mathbf{x} \in \mathcal{X}$
- $h \text{ is right} \Leftrightarrow h(\mathbf{x}) = f(\mathbf{x})$
- *h* is wrong $\Leftrightarrow h(\mathbf{x}) \neq f(\mathbf{x})$
- check h on \mathcal{D} : { (\mathbf{x}_n, y_n) } with i.i.d. \mathbf{x}_n

if large N & i.i.d. \mathbf{x}_n , can probably infer unknown $\llbracket h(\mathbf{x}) \neq f(\mathbf{x}) \rrbracket$ probability by known $\llbracket h(\mathbf{x}) \neq y \rrbracket$ fraction

 $\Rightarrow h(\mathbf{x}) = y$

 $\rightarrow h(\mathbf{x}) \neq \gamma$

Error Measure

- Classification error [...]
 - Often also called '0/1 error'
 - [[true]] = 1, [[false]] = 0

In-sample Error
$$E_{in}(h) = \frac{1}{N} \sum_{n=1}^{N} \llbracket h(\mathbf{x}_n) \neq y_n \rrbracket$$

Out-of-sample Error

$$E_{\text{out}}(h) = \mathop{\mathcal{E}}_{\mathbf{x} \sim P} \left[h(\mathbf{x}) \neq f(\mathbf{x}) \right]$$

 \mathcal{E} : expectation value $\mathbf{x} \sim P$, means the random variable \mathbf{x} has the probability distribution P

Find a Separation Line



The Formal Guarantee

• For any fixed h, in 'big' data (*N* large), in-sample error $E_{in}(h)$ is probably close to out-of-sample error $E_{out}(h)$ (within ϵ)

 $\mathbb{P}[|E_{\text{in}}(h) - E_{\text{out}}(h)| > \epsilon] \le 2\exp(-2\epsilon^2 N)$

 $\mathbb{P}[|\nu - \mu| > \epsilon] \le 2\exp(-2\epsilon^2 N)$

the statement $E_{in}(h) = E_{out}(h)$ is probably approximately correct (PAC)

the statement ' $\nu = \mu$ ' is probably approximately correct (PAC)

If ${}^{*}E_{in}(h) \approx E_{out}(h)$ ' and ${}^{*}E_{in}(h)$ small' $\Rightarrow E_{out}(h)$ small $\Rightarrow h \approx f$ with respect to P

Find a Separation Line



From Fixed h to Set \mathcal{H}

Vapnik-Chervonenkis (VC) bound:

 $\mathbb{P}[\exists h \in \mathcal{H} \text{ s.t. } |E_{\text{in}}(h) - E_{\text{out}}(h)| > \epsilon] \leq 4m_{\mathcal{H}}(2N) \exp\left(-\frac{1}{8}\epsilon^2N\right)$ Model Complexity $O(N^{d_{VC}})$ \downarrow d_{VC} : VC dimension of \mathcal{H}

 $\mathbb{P}[|E_{\text{in}}(h) - E_{\text{out}}(h)| > \epsilon] \le 2\exp(-2\epsilon^2 N)$

Theory of Generalization: $E_{in} \approx E_{out}$ if d_{VC} is small and N is large enough

Find a Separation Line



Media IC & System Lab

Noise & Model Complexity

 $\mathbb{P}[\exists h \in \mathcal{H} \text{ s.t. } |E_{\text{in}}(h) - E_{\text{out}}(h)| > \epsilon] \le 4m_{\mathcal{H}}(2N)\exp\left(-\frac{1}{8}\epsilon^2N\right)$

High Complexity

Low Complexity

Statistical Learning Flow

Outline

- Introduction to Machine Learning
- Theory of Generalization
- Learning Algorithm
- Hazard of Overfitting
- Blending and Bagging

A Simple Hypothesis Set : Perceptron

For x = (x₁, x₂,..., x_d), 'features of customer', compute a weighted 'score' and

Approve ticket if $\sum_{i=1}^{d} w_i x_1 > \text{threshold}$ Deny ticket if $\sum_{i=1}^{d} w_i x_1 < \text{threshold}$

- $\mathcal{Y}: \{+1(good), -1(bad)\}, 0 \text{ ignored}$
 - linear formula $h \in \mathcal{H}$ are

$$h(\mathbf{x}) = \operatorname{sign}\left(\left(\sum_{i=1}^{d} w_i x_1\right) - \operatorname{threshold}\right)$$

called 'perceptron' hypothesis historically

Media IC & System Lab

Vector Form of Perception Hypothesis

$$h(\mathbf{x}) = \operatorname{sign}\left(\left(\sum_{i=1}^{d} w_{i} x_{i}\right) - \operatorname{threshold}\right)$$
$$= \operatorname{sign}\left(\left(\sum_{i=1}^{d} w_{i} x_{i}\right) + (-\operatorname{threshold}) \cdot (+1)\right)$$
$$= \operatorname{sign}\left(\left(\sum_{i=0}^{d} w_{i} x_{i}\right)\right)$$
$$= \operatorname{sign}(\mathbf{w}^{T} \mathbf{x})$$

Each 'tall' w represents hypothesis h & is multiplied with 'tall' x, and we will use tall versions to simplify notation.

perceptron \Leftrightarrow linear (binary) classifiers

Media IC & System Lab

Select g from \mathcal{H}

- $\mathcal{H} =$ all possible perceptrons, g = ?
 - want: $g \approx f$ (hard when f unknown)
 - almost necessary: $g \approx f$ on \mathcal{D}
 - \succ ideally $g(\mathbf{x}_n) = f(\mathbf{x}_n) = y_n$
 - difficult: \mathcal{H} is of infinite size
 - idea: start from some g_0 , and 'correct' its mistakes on \mathcal{D}
 - will represent g_0 by its weight vector \mathbf{w}_0

Perceptron Learning Algorithm

- Start from some w_0 (say 0), and 'correct' its mistakes on $\mathcal D$

for t = 0, 1, ...① find a mistake of \mathbf{w}_t called $(\mathbf{x}_{n(t)}, y_{n(t)})$

 $\operatorname{sign}\left(\mathbf{w}_{t}^{T}\mathbf{x}_{n(t)}\right) \neq y_{n(t)}$

② (try to) correct the mistake by

 $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_{n(t)} \mathbf{x}_{n(t)}$

...until no more mistakes return last w (called w_{PLA}) as g

Line with Noise Tolerance

- Assume 'little' noise: $y_n = f(\mathbf{x}_n)$ usually
- If so, $g \approx f$ on $\mathcal{D} \iff y_n = g(\mathbf{x}_n)$ usually
- How about

$$\mathbf{w}_g \leftarrow \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{n=1}^{N} \llbracket y_n \neq \operatorname{sign}(\mathbf{w}^T \mathbf{x}_n) \rrbracket$$

NP-hard to solve, unfortunately

Can we modify PLA to get an 'approximately good' g?

Pocket Algorithm

 Modify PLA algorithm (black lines) by keeping best weights in pocket.

initialize pocket weights $\widehat{\mathbf{w}}$ for t = 0, 1, ...

① find a (random) mistake of \mathbf{w}_t called $(\mathbf{x}_{n(t)}, y_{n(t)})$ ② (try to) correct the mistake by

 $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_{n(t)} \mathbf{x}_{n(t)}$

③ if \mathbf{w}_{t+1} makes fewer mistakes than $\hat{\mathbf{w}}$, replace $\hat{\mathbf{w}}$ by \mathbf{w}_{t+1} ...until enough iterations return $\hat{\mathbf{w}}$ (called $\mathbf{w}_{\text{POCKET}}$) as g

Regression: $\mathcal{Y} = \mathbb{R}$

- Linear regression
 - find lines/hyperplanes with small residuals

Error Measure

- Popular/historical error measure:
 - square error $err(\hat{y}, y) = (\hat{y} y)^2$

Out-of-sample Error

How to minimize $E_{in}(\mathbf{w})$?

Matrix Form of $E_{in}(\mathbf{w})$

$$E_{in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{w}^{T} \mathbf{x}_{n} - y_{n})^{2} = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{x}_{n}^{T} \mathbf{w} - y_{n})^{2}$$
$$= \frac{1}{N} \left\| \begin{vmatrix} \mathbf{x}_{1}^{T} \mathbf{w} - y_{1} \\ \mathbf{x}_{2}^{T} \mathbf{w} - y_{2} \\ \dots \\ \mathbf{x}_{N}^{T} \mathbf{w} - y_{N} \end{vmatrix} \right|^{2}$$
$$= \frac{1}{N} \left\| \begin{vmatrix} -\mathbf{x}_{1}^{T} - - \\ -\mathbf{x}_{2}^{T} - - \\ \dots \\ -\mathbf{x}_{N}^{T} - - \end{vmatrix} \mathbf{w} - \begin{bmatrix} y_{1} \\ y_{2} \\ \dots \\ y_{N} \end{bmatrix} \right\|^{2}$$
$$= \frac{1}{N} \left\| |\mathbf{X}\mathbf{w} - \mathbf{y}||^{2} \quad \begin{array}{c} \mathbf{x} : N \times (d+1) \\ \mathbf{w} : (d+1) \times 1 \\ \mathbf{y} : N \times 1 \end{vmatrix} \right\|$$

How to Minimize $E_{in}(\mathbf{w})$?

- $E_{in}(\mathbf{w})$: continuous, differentiable, **convex**
- Necessary condition of 'best' w

$$\nabla E_{\rm in}(\mathbf{w}) \equiv \begin{bmatrix} \frac{\partial E_{\rm in}}{\partial w_0}(\mathbf{w}) \\ \frac{\partial E_{\rm in}}{\partial w_1}(\mathbf{w}) \\ \dots \\ \frac{\partial E_{\rm in}}{\partial w_d}(\mathbf{w}) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \dots \\ 0 \end{bmatrix} \qquad E_{\rm in}$$

not possible to 'roll down'

task: find \mathbf{w}_{LIN} such that $\nabla E_{\text{in}}(\mathbf{w}_{\text{LIN}}) = 0$

W

Recap: Matrix Calculus

- Several useful vector derivative formulas
 - $\mathbf{A} : n \times n$
 - $\mathbf{x} : n \times 1$
 - **y** : $n \times 1$

Denominator-layout notation

The Gradient $E_{in}(\mathbf{w})$

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = \frac{1}{N} (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y})$$

- $\nabla E_{in}(\mathbf{w}) = \frac{1}{N} (2\mathbf{X}^T \mathbf{X} \mathbf{w} 2\mathbf{X}^T \mathbf{y})$
- Task: find \mathbf{w}_{LIN} such that $\frac{2}{N}(\mathbf{X}^T\mathbf{X}\mathbf{w} \mathbf{X}^T\mathbf{y}) = \mathbf{0}$

invertible $\mathbf{X}^T \mathbf{X}$

 unique solution

 w_{LIN} = (X^TX)⁻¹X^T y pseudo-inverse X[†]

 often the case because

 $N \gg d + 1$

singular $\mathbf{X}^T \mathbf{X}$

- many optimal solutions
- one of the solutions $\mathbf{w}_{\text{LIN}} = \mathbf{X}^{\dagger} \mathbf{y}$

$$(\mathbf{X} \equiv \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T \longrightarrow \mathbf{X}^\dagger \equiv \mathbf{V} \mathbf{\Sigma}^\dagger \mathbf{U}^T)$$

Linear Regression Algorithm

① from \mathcal{D} , construct input matrix **X** and output vector **y** by

$$\mathbf{X} = \begin{bmatrix} --\mathbf{x}_{1}^{T} - - \\ --\mathbf{x}_{2}^{T} - - \\ \dots \\ --\mathbf{x}_{N}^{T} - - \end{bmatrix}, \qquad \mathbf{y} = \begin{bmatrix} y_{1} \\ y_{2} \\ \dots \\ y_{N} \\ \dots \\ N \times 1 \end{bmatrix}$$

② calculate pseudo-inverse $\underbrace{X^{\dagger}}_{(d+1)\times N}$

Logistic Hypothesis

For x = (x₀, x₁, x₂,..., x_d), 'features of customer', compute a weighted 'score':

$$\mathbf{s} = \sum_{i=0}^{d} \mathbf{w}_i x_i = \mathbf{w}^T \mathbf{x}$$

logistic hypothesis: $h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x})$

S

 $\theta(s)$

Logistic Function

$$\theta(-\infty) = 0$$

$$\theta(0) = \frac{1}{2} \implies \theta(s) = \frac{e^s}{1 + e^s} = \frac{1}{1 + e^{-s}}$$

$$\theta(\infty) = 1$$

smooth monotonic sigmoid function of

smooth, monotonic, sigmoid function of s

Logistic regression: use

$$h(\mathbf{x}) = \theta(s) = \theta(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

to approximate target function $f(\mathbf{x}) = P(y|\mathbf{x})$

Cross-Entropy Error

$$g = \underset{\mathbf{w}}{\operatorname{argmax}} \prod_{n=1}^{N} \theta(y_{n} \mathbf{w}^{T} \mathbf{x}_{n}) \equiv \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{n=1}^{N} -\ln\theta(y_{n} \mathbf{w}^{T} \mathbf{x}_{n})$$

• $E_{in}(\mathbf{w})$ for logistic regression:

$$\frac{1}{N}\sum_{n=1}^{N}\ln\left(1+\exp(-y_{n}\mathbf{w}^{T}\mathbf{x}_{n})\right) = \frac{1}{N}\sum_{n=1}^{N}\exp(\mathbf{w},\mathbf{x}_{n},y_{n})$$

$$\operatorname{err}(\mathbf{w}, \mathbf{x}_n, y_n) = \ln(1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n))$$

cross-entropy error

Minimizing $E_{in}(\mathbf{w})$

$$\underset{\mathbf{w}}{\operatorname{argmin}} E_{\operatorname{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^{N} \ln(1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n))$$

• $E_{in}(\mathbf{w})$: continuous, differentiable, **convex**

• How to minimize? Locate valley $\Rightarrow \nabla E_{in}(\mathbf{w}) = 0$

The Gradient $\nabla E_{in}(\mathbf{w})$

$$\operatorname{argmin}_{\mathbf{w}} E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{\substack{n=1 \\ N}}^{N} \ln(1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n))$$
$$\Rightarrow \nabla E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{\substack{n=1 \\ n=1}}^{N} \theta(-y_n \mathbf{w}^T \mathbf{x}_n)(-y_n \mathbf{x}_n) = 0$$

- Scaled θ -weighted sum of $-y_n \mathbf{x}_n$
 - $\text{ all } \theta(\cdot) = 0 ?$
 - > only if $y_n \mathbf{w}^T \mathbf{x}_n \gg 0$ (linear separable \mathcal{D})
 - weighted sum = 0?
 - non-linear

No closed-form solution!

 $\theta(s)$

Iterative Optimization

for t = 0, 1, ...

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta \mathbf{v}$$

when stop, return last w as g

- Smooth $E_{in}(\mathbf{w})$ for logistic regression: choose **v** to get the ball roll 'downhill'?
 - direction v (assumed of unit length)
 - step size η (assumed positive)
- A greedy approach for some given $\eta > 0$:

$$\min_{\|\mathbf{v}\|=1} E_{\mathrm{in}}\left(\underbrace{\mathbf{w}_t + \eta \mathbf{v}}_{\mathbf{w}_{t+1}}\right)$$

Gradient Descent

 $\min_{\|\mathbf{v}\|=1} E_{\mathrm{in}}(\mathbf{W}_t + \boldsymbol{\eta}\mathbf{v})$

• If η is small, then by Taylor expansion:

$$E_{\rm in}(\mathbf{w}_t + \eta \mathbf{v}) \approx \underbrace{E_{\rm in}(\mathbf{w}_t)}_{\rm known} + \underbrace{\eta}_{\rm given \ positive} \mathbf{v}^T \underbrace{\nabla E_{\rm in}(\mathbf{w}_t)}_{\rm known}$$

• Optimal **v** : opposite direction of $\nabla E_{in}(\mathbf{w}_t)$

$$\mathbf{v} = -\frac{\nabla E_{\text{in}}(\mathbf{w}_t)}{\|\nabla E_{\text{in}}(\mathbf{w}_t)\|}$$

• Gradient descent: for small η

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta \frac{\nabla E_{\text{in}}(\mathbf{w}_t)}{\|\nabla E_{\text{in}}(\mathbf{w}_t)\|}$$

Gradient descent: a simple & popular optimization tool

Choice of η

• η should better be monotonic of $\|\nabla E_{in}(\mathbf{w}_t)\|$

• If red $\eta \propto \|\nabla E_{in}(\mathbf{w}_t)\|$ by ratio purple η

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta \frac{\nabla E_{\text{in}}(\mathbf{w}_t)}{\|\nabla E_{\text{in}}(\mathbf{w}_t)\|} \implies \mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta \nabla E_{\text{in}}(\mathbf{w}_t)$$

$$\uparrow$$
the fixed learning rate

Logistic Regression Algorithm

initialize \mathbf{w}_0

for t = 0, 1, ...

① compute

$$\nabla E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^{N} \theta(-y_n \mathbf{w}^T \mathbf{x}_n) (-y_n \mathbf{x}_n)$$

② update by

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta \nabla E_{\text{in}}(\mathbf{w}_t)$$

...until $\nabla E_{in}(\mathbf{w}_t) = 0$ or enough iterations return last \mathbf{w}_t as g

Stochastic Gradient Descent (SGD)

SGD logistic regression:

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta \theta (-y_n \mathbf{w}^T \mathbf{x}_n) (y_n \mathbf{x}_n)$$

over uniform choice of *n*

- Idea: replace true gradient by stochastic gradient
- After enough steps,

average true gradient ≈ average stochastic gradient

- pros: simple & cheaper computation
 useful for big data or online learning
- cons: less stable in nature

Three Linear Models

• Linear weighted sum: $\mathbf{s} = \mathbf{w}^T \mathbf{x}$

Outline

- Introduction to Machine Learning
- Theory of Generalization
- Learning Algorithm
- Hazard of Overfitting
- Blending and Bagging

Circular Separable

 D is not linear separable but circular separable by a circle of radius <u>3</u> centered at origin:

$$h(\mathbf{x}) = \operatorname{sign}(-x_1^2 - x_2^2 + 9)$$

= sign(9 \cdot 1 + (-1) \cdot x_1^2 + (-1) \cdot x_2^2)
 $\widetilde{w}_0 \ z_0 \ \widetilde{w}_1 \ z_1 \ z_1 \ \widetilde{w}_2 \ z_2$
= sign($\widetilde{w}^T \mathbf{z}$)

Circular Separable and Linear Separable

• $\{(\mathbf{x}_n, y_n)\}$ circular separable $\Rightarrow \{(\mathbf{z}_n, y_n)\}$ linear separable

• Nonlinear Feature Transform: $\mathbf{x} \in \mathcal{X} \xrightarrow{\Phi} \mathbf{z} \in \mathcal{Z}$

$$(z_0, z_1, z_2) = \mathbf{z} = \mathbf{\Phi}(\mathbf{x}) = (1, x_1^2, x_2^2)$$

 $h(\mathbf{x}) = \tilde{h}(\mathbf{z}) = \operatorname{sign}(\widetilde{w}^T \Phi(\mathbf{x})) = \operatorname{sign}(\widetilde{w}_0 + \widetilde{w}_1 x_1^2 + \widetilde{w}_2 x_2^2)$

Linear Hypothesis in **Z**-Space

General quadratic hypothesis set: A 'bigger' Z-Space

$$\mathbf{\Phi}_2(\mathbf{x}) = (1, x_1, x_2, x_1^2, x_1 x_2, x_2^2)$$

- perceptron in \mathcal{Z} -Space \Leftrightarrow quadratic hypotheses in \mathcal{X} -Space
- *Q*-th order polynomial transform:

 $\mathbf{\Phi}_Q(\mathbf{x}) = \left(1, x_1, x_2, \dots, x_d, x_1^2, x_1 x_2, \dots, x_d^2, \dots, x_1^Q, x_1^{Q-1} x_2, \dots, x_d^Q\right)$

$$- \underbrace{1}_{\widetilde{w}_0} + \underbrace{\widetilde{d}}_{\text{others}} \text{ dimensions} = O(Q^d)$$

- number of free parameters $\widetilde{w}_i = 1 + \widetilde{d} \approx d_{VC}(\mathcal{H}_{\Phi_Q})$ $\mathcal{H}_{\Phi_Q} = \left\{ h(\mathbf{x}) : h(\mathbf{x}) = \widetilde{h}\left(\Phi_Q(\mathbf{x})\right) \text{ for some linear } \widetilde{h} \text{ on } \mathbf{Z} \right\}$

 $Q \text{ large} \Rightarrow \text{large } d_{\text{VC}}$

Po-Chen Wu (吳柏辰)

Hazard of Overfitting

• Vapnik-Chervonenkis (VC) bound (remember?):

 $\mathbb{P}[\exists h \in \mathcal{H} \text{ s.t. } | \underline{E_{\text{in}}(h)} - \underline{E_{\text{out}}(h)} | > \epsilon] \le 4m_{\mathcal{H}}(2N) \exp\left(-\frac{1}{8}\epsilon^2N\right)$

Model Complexity $O(N^{d_{VC}})$

Regularization: The Magic

• *Q*-th order polynomial transform for $x \in \mathbb{R}$ + linear regression:

• Idea: 'step back' from \mathcal{H}_{10} to \mathcal{H}_{2}

Stepping Back as Constraint

- hypothesis w in \mathcal{H}_{10} : $w_0 + w_1 x + w_2 x^2 + w_3 x^3, \dots, + w_{10} x^{10}$
- hypothesis **w** in \mathcal{H}_2 : $w_0 + w_1 x + w_2 x^2$
- \Rightarrow that is, $|\mathcal{H}_2| = |\mathcal{H}_{10}$ & constraint that $w_3 = w_4 = \cdots = w_{10} = 0$

$\mathcal{H}_{10} \equiv \{ \mathbf{w} \in \mathbb{R}^{10+1} \}$	$\mathcal{H}_2 \equiv \{ \mathbf{w} \in \mathbb{R}^{10+1} \text{ while } w_3 = w_4 = \dots = w_{10} = 0 \}$
regression with \mathcal{H}_{10} :	regression with \mathcal{H}_2 :
$\min_{\mathbf{w}\in\mathbb{R}^{10+1}}E_{\mathrm{in}}(\mathbf{w})$	$\min_{\mathbf{w}\in\mathbb{R}^{10+1}} E_{in}(\mathbf{w}) \text{ s. t. } w_3 = w_4 = \dots = w_{10} = 0$

step back = constrained optimization of E_{in}

In fact, we can just use $\mathbf{w} \in \mathbb{R}^{2+1}$ in this case.

Regression with Looser Constraint

- regression with $\mathcal{H}_2 \equiv \{\mathbf{w} \in \mathbb{R}^{10+1}, \text{ while } w_3 = \cdots = w_{10} = 0\}$: $\min_{\mathbf{w} \in \mathbb{R}^{10+1}} E_{\text{in}}(\mathbf{w}) \text{ s.t. } w_3 = \cdots = w_{10} = 0$
- regression with $\mathcal{H}'_2 \equiv \{\mathbf{w} \in \mathbb{R}^{10+1}, \text{ while } \geq 8 \text{ of } w_q = 0 \}$: $\min_{\mathbf{w} \in \mathbb{R}^{10+1}} E_{\text{in}}(\mathbf{w}) \text{ s.t.} \sum_{q=0}^{10} \left[w_q \neq 0 \right] \leq 3$
 - more flexible than \mathcal{H}_2 : $\mathcal{H}_2 \subset \mathcal{H}'_2$
 - $\text{ less risky than } \mathcal{H}_{10}^{\prime} : \qquad \qquad \mathcal{H}_2^{\prime} \subset \mathcal{H}_{10}$

Bad news for sparse hypothesis set \mathcal{H}'_2 : NP-hard to solve.

Regression with Softer Constraint

- regression with $\mathcal{H}_{2}' \equiv \{\mathbf{w} \in \mathbb{R}^{10+1}, \text{ while } \geq 8 \text{ of } w_{q} = 0\}$: $\min_{\mathbf{w} \in \mathbb{R}^{10+1}} E_{\text{in}}(\mathbf{w}) \text{ s.t.} \sum_{q=0}^{10} \llbracket w_{q} \neq 0 \rrbracket \leq 3$ • regression with $\mathcal{H}(C) \equiv \{\mathbf{w} \in \mathbb{R}^{10+1}, \text{ while } \|\mathbf{w}\|^{2} \leq C\}$: $\min_{\mathbf{w} \in \mathbb{R}^{10+1}} E_{\text{in}}(\mathbf{w}) \text{ s.t.} \sum_{q=0}^{10} w_{q}^{2} \leq C$
 - $\mathcal{H}(\mathcal{C})$: overlaps but not exactly the same as \mathcal{H}_2'
 - soft and smooth strudcture over $C \ge 0$ $\mathcal{H}(0) \subset \mathcal{H}(1) \subset \mathcal{H}(2) \subset \cdots \subset \mathcal{H}(\infty) = \mathcal{H}_{10}$

regularized hypothesis w_{REG} :

optimal solution from regularized hypothesis set $\mathcal{H}(C)$

The Lagrange Multiplier

•
$$\min_{\mathbf{w}\in\mathbb{R}^{10+1}} E_{\mathrm{in}}(\mathbf{w}) = \frac{1}{N} \underbrace{\sum_{n=1}^{N} (\mathbf{w}^{T} \mathbf{z}_{n} - \mathbf{y}_{n})^{2}}_{(\mathbf{Z}\mathbf{w}-\mathbf{y})^{T}(\mathbf{Z}\mathbf{w}-\mathbf{y})} \text{ s.t. } \underbrace{\sum_{q=0}^{Q} w_{q}^{2}}_{\mathbf{w}^{T}\mathbf{w}} \leq C$$
$$\Rightarrow \min_{\mathbf{w}\in\mathbb{R}^{10+1}} E_{\mathrm{in}}(\mathbf{w}) = \frac{1}{N} (\mathbf{Z}\mathbf{w} - \mathbf{y})^{T} (\mathbf{Z}\mathbf{w} - \mathbf{y}) \text{ s.t. } \mathbf{w}^{T}\mathbf{w} \leq C$$

 Assume w^Tw = C, find Lagrange multiplier λ > 0 and w_{REG} such that

$$\frac{\partial}{\partial \mathbf{w}} (E_{in}(\mathbf{w}) + \mathbf{\lambda} \mathbf{w}^T \mathbf{w}) = 0$$

$$\Rightarrow \nabla E_{in}(\mathbf{w}_{REG}) + 2\mathbf{\lambda} \mathbf{w}_{REG} = 0$$

$$\Rightarrow \frac{2}{N} (\mathbf{Z}^T \mathbf{Z} \mathbf{w}_{REG} - \mathbf{Z}^T \mathbf{y}) + 2\mathbf{\lambda} \mathbf{w}_{REG} = 0$$

$$\lambda/N$$

Ridge Regression

$$\frac{\partial}{\partial \mathbf{w}} \left(E_{\rm in}(\mathbf{w}) + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w} \right) = 0$$

- Augmented error $E_{aug}(\mathbf{w})$: $E_{in}(\mathbf{w}) + \frac{\lambda}{N} \underbrace{\mathbf{w}^T \mathbf{w}}_{regularizer}$
- Regularization with E_{aug} instead of constrained E_{in} :

$$\mathbf{w}_{\text{REG}} \leftarrow \underset{\mathbf{w}}{\operatorname{argmin}} E_{\operatorname{aug}}(\mathbf{w}) \text{ for given } \lambda > 0$$

optimal solution:

$$\frac{\partial}{\partial \mathbf{w}_{\text{REG}}} E_{\text{aug}}(\mathbf{w}_{\text{REG}}) = 0 \implies \mathbf{w}_{\text{REG}} \leftarrow (\mathbf{Z}^T \mathbf{Z} + \lambda \mathbf{I})^{-1} \mathbf{Z}^T \mathbf{y}$$

- called ridge regression in Statistics

minimizing unconstrained E_{aug} effectively minimizes some *C*-constrained E_{in}

Media IC & System Lab

Po-Chen Wu (吳柏辰)

The Results

Seeing is believing

- philosophy: a little regularization goes a long way!
- call $\frac{\lambda}{N} \mathbf{w}^T \mathbf{w}$ weight-decay regularization:

➤ larger $\lambda \iff$ prefer shorter w ⇔ effectively smaller C

The VC Message

• The best d_{VC}^* is in the middle.

Model Selection Problem

• Which one is better?

selecting by E_{val}

Validation Set \mathcal{D}_{val}

- $\mathcal{D}_{val} \subset \mathcal{D}$: called validation set
 - 'on-hand' simulation of test set
- Make sure \mathcal{D}_{val} is 'clean'
 - feed only $\mathcal{D}_{\text{train}}$ to \mathcal{A}_m for model selection

Model Selection by Best E_{val}

Po-Chen Wu (吳柏辰)

V-fold Cross Validation

- V-fold cross-validation: random-partition of D to V equal parts
 - take V 1 for training and 1 for validation orderly

$$E_{\rm CV}(\mathcal{H},\mathcal{A}) = \frac{1}{V} \sum_{\nu=1}^{V} E_{\rm VAL}^{(\nu)}(g_{\nu}^{-})$$

- selection by $E_{\rm CV}$:

$$m^* = \underset{1 \le m \le M}{\operatorname{argmin}} \left(E_m = E_{\text{CV}}(\mathcal{H}_m, \mathcal{A}_m) \right)$$

practical rule of thumb: V = 10

training

validation

Outline

- Introduction to Machine Learning
- Theory of Generalization
- Learning Algorithm
- Hazard of Overfitting
- Blending and Bagging

Blending (Aggregation)

- Blending (or aggregation) : mix or combine hypotheses for better performance
 - Uniform blending (voting) for classification:

$$G(\mathbf{x}) = \operatorname{sign}\left(\sum_{m=1}^{M} 1 \cdot g_m(\mathbf{x})\right)$$

– Uniform blending regression:

$$G(\mathbf{x}) = \frac{1}{T} \sum_{m=1}^{M} 1 \cdot \frac{g_m}{\mathbf{x}}(\mathbf{x})$$

– Linear blending:

$$G(\mathbf{x}) = \operatorname{sign}\left(\sum_{m=1}^{M} \alpha_t \cdot g_m(\mathbf{x})\right) \text{ with } \alpha_m \ge 0$$

Bagging (Bootstrap Aggregation)

- Bootstrapping
 - Bootstrap sample $\widetilde{\mathcal{D}}_m$: re-sample N examples from \mathcal{D} uniformly with replacement can also use arbitrary N' instead of original N
- Bootstrap Aggregation (BAGging)
 - Consider a iterative process that for m = 1, 2, ..., M
 - ① Request size-N' data $\widetilde{\mathcal{D}}_m$ from bootstrapping
 - ② Obtain g_t by $\mathcal{A}(\widetilde{\mathcal{D}}_m)$, and $G = \text{Uniform}(\{g_t\})$

Bootstrap aggregation (Bagging): a simple meta algorithm on top of base algorithm ${\cal A}$

Reference

 Machine learning slides by Prof. Hsuan-Tien Lin <u>http://www.csie.ntu.edu.tw/~htlin/course/ml14fall/</u>