



# วิทยานิพนธ์

ออนไลน์อัลกอริทึมสำหรับแก้ปัญหาการทำนายแบบมัลติคลาสที่แต่ละ  
เหตุการณ์มีน้ำหนักความสำคัญไม่เท่ากัน

ONLINE LEARNING ALGORITHM FOR MULTICLASS  
IMPORTANCE WEIGHTED PREDICTION PROBLEM

นายอดิศักดิ์ สุภิสุน

บัณฑิตวิทยาลัย มหาวิทยาลัยเกษตรศาสตร์

พ.ศ. 2551



ใบรับรองวิทยานิพนธ์  
บัณฑิตวิทยาลัย มหาวิทยาลัยเกษตรศาสตร์

วิศวกรรมศาสตรมหาบัณฑิต (วิศวกรรมคอมพิวเตอร์)

ปริญญา

วิศวกรรมคอมพิวเตอร์

วิศวกรรมคอมพิวเตอร์

สาขา

ภาควิชา

เรื่อง ออนไลน์อัลกอริทึมสำหรับแก้ปัญหการทำนายแบบมัลติคลาสที่แต่ละเหตุการณ์  
มีน้ำหนักความสำคัญไม่เท่ากัน

Online Learning Algorithm for Multiclass Importance Weighted Prediction Problem

นามผู้วิจัย นายอดิศักดิ์ สุภิสุน

ได้พิจารณาเห็นชอบโดย

อาจารย์ที่ปรึกษาวิทยานิพนธ์หลัก



( ผู้ช่วยศาสตราจารย์จিতร์ทัศน์ ฝักเจริญผล, Ph.D. )

อาจารย์ที่ปรึกษาวิทยานิพนธ์ร่วม



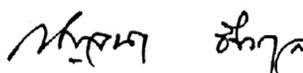
( รองศาสตราจารย์กฤษณะ ไวยมัย, D.U. )

หัวหน้าภาควิชา



( ผู้ช่วยศาสตราจารย์เข้มะชาติ วิภาตะวินิช, Ph.D. )

บัณฑิตวิทยาลัย มหาวิทยาลัยเกษตรศาสตร์รับรองแล้ว



( รองศาสตราจารย์กัญจนา วีระกุล, D.Agr. )

คณบดีบัณฑิตวิทยาลัย

วันที่ 29 เดือน พฤษภาคม พ.ศ. 2551

วิทยานิพนธ์

เรื่อง

ออนไลน์อัลกอริทึมสำหรับแก้ปัญหาการทำนายแบบมัลติคลาสที่แต่ละเหตุการณ์มีน้ำหนัก  
ความสำคัญไม่เท่ากัน

Online Learning Algorithm for Multiclass Importance Weighted Prediction Problem

โดย

นายอดิศักดิ์ สุทธิสุน

เสนอ

บัณฑิตวิทยาลัย มหาวิทยาลัยเกษตรศาสตร์  
เพื่อความสมบูรณ์แห่งปริญญาวิศวกรรมศาสตรมหาบัณฑิต (วิศวกรรมคอมพิวเตอร์)

พ.ศ. 2551

อดิศักดิ์ สุภีสุน 2551: ออนไลน์อัลกอริทึมสำหรับแก้ปัญหาการทำนายแบบมัลติคลาสที่แต่ละเหตุการณ์มีน้ำหนักความสำคัญไม่เท่ากัน วิทยุวิศวกรรมศาสตรมหาบัณฑิต (วิศวกรรมคอมพิวเตอร์) สาขาวิศวกรรมคอมพิวเตอร์ ภาควิชาวิศวกรรมคอมพิวเตอร์ อาจารย์ที่ปรึกษาวิทยานิพนธ์หลัก: ผู้ช่วยศาสตราจารย์จิตรีทัศน์ ฝักเจริญผล, Ph.D. 48 หน้า

ในวิทยานิพนธ์ฉบับนี้ เรานำเสนอการแก้ปัญหาการเรียนรู้ออนไลน์ที่ตัวอย่างมีน้ำหนักความสำคัญสองวิธีการ วิธีแรกใช้การเรียนรู้แบบลดรูป กล่าวคือ เราแสดงวิธีการลดรูปปัญหาการเรียนรู้ออนไลน์แบบมีน้ำหนักไปยังปัญหาการเรียนรู้ออนไลน์มาตรฐาน โดยการลดรูปดังกล่าวทำให้ขีดจำกัดด้านบนของความผิดพลาดของอัลกอริทึมเพิ่มขึ้นไม่เกิน  $c_{max}$  เท่า เมื่อ  $c_{max}$  คือน้ำหนักที่มากที่สุด ในส่วนของวิธีการที่สอง เรานำเสนออัลกอริทึมที่เป็นการปรับปรุงอัลกอริทึมเพอร์เซพตรอนอย่างง่าย และสามารถทำงานได้ในสถานการณ์เดียวกับเพอร์เซพตรอน อัลกอริทึมที่สองนี้สามารถใช้กับปัญหาการเรียนรู้แบบหลายประเภทได้ด้วย โดยใช้วิธีลดรูปจากปัญหาการเรียนรู้หลายประเภทไปเป็นปัญหาการเรียนรู้สองประเภท ซึ่งเราได้พิสูจน์ว่า การลดรูปการเรียนรู้ดังกล่าว 2 วิธีการที่แตกต่างกันนั้นสมมูลกัน

จากวิธีการที่เรานำเสนอ 2 วิธีข้างต้น วิธีที่สองมีการรับประกันเชิงทฤษฎีที่ดีกว่าวิธีแรกในหลายๆ กรณี รวมทั้งให้ผลการทดลองที่ดีกว่า อย่างไรก็ตาม วิธีแรกเป็นวิธีที่มีความเป็นทั่วไปมากกว่า เรายังได้ทำการทดลองเพื่อเปรียบเทียบประสิทธิภาพของอัลกอริทึมที่ได้แนะนำพร้อมกับชุดข้อมูลทดสอบจริงและชุดข้อมูลสังเคราะห์

อดิศักดิ์ สุภีสุน  
ลายมือชื่อนิสิต

อดิศักดิ์ สุภีสุน 26 พฤษภาคม 2551  
ลายมือชื่ออาจารย์ที่ปรึกษาวิทยานิพนธ์หลัก

Adisak Supeesun 2008: Online Learning Algorithm for Multiclass Importance Weighted Prediction Problem. Master of Engineering (Computer Engineering), Major Field: Computer Engineering, Department of Computer Engineering. Thesis Advisor: Assistant Professor Jittat Fakcharoenphol, Ph.D. 48 pages.

This thesis considers an online learning problem with importance weighted examples. We present two approaches. The first one uses learning reduction, i.e., we show how to reduce online importance weighted learning problems to standard online learning problems with a factor  $c_{\max}$  overhead on mistake bounds,  $c_{\max}$  is the maximum weight of examples. We also present another online algorithm based on perceptron learning algorithm. The second algorithm is a simple modification of perceptron algorithm and work in this same setting and by modifying the proof of gentile we can proof weighted loss bound of this algorithm. For the multiclass setting, the second algorithm can work in the setting by using multiclass-binary learning reduction. We also show two different method, the multi-vector method and its particular all-pair variant, for reducing multiclass problems to binary problems are equivalent.

While the second algorithm performs better on experiments and has a better theoretical guarantee in many cases, the first one is more general. We also perform experiments on real data sets and synthetic data sets.

Adisak Supeesun

Student's signature

Jittat Fakcharoenphol

Thesis Advisor's signature

26 / May / 2008

## กิตติกรรมประกาศ

วิทยานิพนธ์นี้สำเร็จลงได้ด้วยความช่วยเหลือจากผู้ช่วยศาสตราจารย์จิตรีทัศน์ ฝักเจริญผล ประธานกรรมการที่ปรึกษา ผู้ให้คำปรึกษาและแนวทางการทำวิจัย ตลอดจนข้อเสนอแนะที่เป็นประโยชน์ต่องานวิจัยนี้ รองศาสตราจารย์กฤษณะ ไวยมัย กรรมการที่ปรึกษาร่วม ที่กรุณาให้คำปรึกษา และข้อเสนอแนะที่มีคุณค่าเพื่อให้วิทยานิพนธ์นี้สมบูรณ์ยิ่งขึ้น

ข้าพเจ้าขอขอบคุณเพื่อนๆ นิสิตปริญญาโท, สมาชิกกลุ่มวิจัยเชิงทฤษฎีทุกท่านที่ช่วยเหลือและให้คำปรึกษาแก่ข้าพเจ้าจนสามารถทำงานวิจัยชิ้นนี้สำเร็จ ขอขอบคุณเจ้าหน้าที่ธุรการ ภาควิชาวิศวกรรมคอมพิวเตอร์ คณะวิศวกรรมศาสตร์ มหาวิทยาลัยเกษตรศาสตร์ที่คอยอำนวยความสะดวกด้านประสานงาน และการดำเนินงานต่างๆ

อดิศักดิ์ สุภีสุน

พฤษภาคม 2551

## สารบัญ

	หน้า
สารบัญ	(1)
สารบัญตาราง	(2)
สารบัญภาพ	(3)
คำนำ	1
วัตถุประสงค์	4
การตรวจเอกสาร	5
วิธีการ	15
ผลการวิจัย	21
สรุป	46
เอกสารและสิ่งอ้างอิง	47

## สารบัญญัตินำ

ตารางที่		หน้า
1	ผลการเปรียบเทียบประสิทธิภาพการทำนายของอัลกอริทึมแบบลดรูป, อัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญ และอัลกอริทึมเพอร์เซพตรอนกับชุดข้อมูลแบบสองประเภท	25
2	ผลการเปรียบเทียบประสิทธิภาพการทำนายของอัลกอริทึมแบบลดรูป, อัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญ และอัลกอริทึมเพอร์เซพตรอนเมื่อทำการลดรูปจากปัญหาการทำนายหลายประเภทเป็นปัญหาการทำนายหลายประเภทด้วยวิธีการหลายเวกเตอร์ กับชุดข้อมูลแบบหลายประเภท	40
3	ผลการเปรียบเทียบประสิทธิภาพการทำนายของอัลกอริทึมแบบลดรูป, อัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญ และอัลกอริทึมเพอร์เซพตรอนเมื่อทำการลดรูปจากปัญหาการทำนายหลายประเภทเป็นปัญหาการทำนายหลายประเภทด้วยวิธีการทุกคู่ กับชุดข้อมูลแบบหลายประเภท	40

## สารบัญญภาพ

ภาพที่		หน้า
1	ตัวอย่างการทำนายด้วยเพอร์เซพตรอน	8
2	อัลกอริทึมแบบลดรูป	16
3	อัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญ	17
4	เวกเตอร์น้ำหนักของวิธีการทุกคู่	19
5	การกระจายตัวของข้อมูลสังเคราะห์ชุดที่ 1	23
6	การกระจายตัวของข้อมูลสังเคราะห์ชุดที่ 2	24
7	การกระจายตัวของข้อมูลสังเคราะห์ชุดที่ 3	24

# ออนไลน์อัลกอริทึมสำหรับแก้ปัญหาการทำนายแบบมัลติคลาสที่แต่ละเหตุการณ์มีน้ำหนักความสำคัญไม่เท่ากัน

## Online Learning Algorithm for Multiclass Importance Weighted Prediction

### Problem

#### คำนำ

ในปัจจุบันมีการนำเอาวิธีการทางด้านการเรียนรู้ด้วยเครื่องจักร (machine learning) ไปประยุกต์ใช้กับการแก้ปัญหาของงานหลาย ๆ ประเภท กระบวนการทั่วไปจะเริ่มจากการหาตัวอย่างที่สุ่มจากการกระจาย (distribution) เดียวกับสถานการณ์ที่ใช้จริง แล้วนำตัวอย่างไปใช้เพื่อสร้างตัวทำนายที่มีคุณภาพ จากนั้นนำตัวทำนายที่ได้ไปใช้จริง ลักษณะของปัญหาดังกล่าวจัดเป็น *ปัญหาการเรียนรู้แบบออฟไลน์* กล่าวคือ ในการฝึกสอนตัวทำนายนั้น อัลกอริทึมจะได้ข้อมูลสำหรับการฝึกสอนทั้งหมดมาก่อน งานวิจัยนี้สนใจ *ปัญหาการเรียนรู้แบบออนไลน์* กล่าวโดยย่อคือ อัลกอริทึมจะสร้างตัวทำนายพร้อมๆ กับที่ได้รับข้อมูลทดสอบ ปัญหาหลักขณะนี้เกิดขึ้นในกรณีที่ผู้สร้างตัวทำนายไม่สามารถหาตัวอย่างการเรียนรู้ได้เนื่องจากไม่มีแหล่งข้อมูลหรือไม่ทราบการกระจายที่แท้จริงหรือใกล้เคียงกับที่ต้องการ หรือในกรณีที่การกระจายของตัวอย่างและคำตอบเปลี่ยนแปลงไม่แน่นอน

ในงานวิจัยนี้สนใจปัญหาการเรียนรู้แบบออนไลน์ในกรณีที่ตัวอย่างมีน้ำหนักความสำคัญ (Importance weighted online learning) ในปัญหานี้ อัลกอริทึมจะทำงานเป็นรอบๆ ในแต่ละรอบ อัลกอริทึมจะได้รับตัวอย่าง  $x$  และจะต้องทำนายประเภทของ  $x$  ก่อนที่จะได้รับค่าเฉลย  $y$  และน้ำหนักความสำคัญ  $c$  ของ  $x$  ถ้าอัลกอริทึมทำนายผิด อัลกอริทึมจะถูกปรับด้วยค่า  $c$  ประสิทธิภาพของอัลกอริทึมจะวัดจากผลรวมของการถูกปรับ เป้าหมายก็คือต้องการให้ผลรวมดังกล่าวมีน้อยที่สุด ส่วนปัญหาการเรียนรู้ที่ตัวอย่างมีน้ำหนักความสำคัญนั้นมีปรากฏขึ้นในการนำการเรียนรู้ด้วยเครื่องจักรไปใช้ในสถานการณ์จริง โดยมากมักพบในการนำไปประยุกต์ใช้ในทางธุรกิจที่ค่าใช้จ่ายกับผลตอบแทนของลูกค้าแต่ละรายไม่เท่ากัน สำหรับสถานการณ์แบบออนไลน์มักพบในการเสนอขายสินค้าให้กับลูกค้า เป็นต้น

ปัญหาการเรียนรู้ทั้งแบบที่ตัวอย่างไม่มีน้ำหนักความสำคัญและแบบที่ตัวอย่างมีน้ำหนักความสำคัญ ที่กล่าวมายังสามารถแบ่งได้อีก 2 กลุ่ม ได้แก่ ปัญหาการเรียนรู้สองประเภท โดยปัญหานี้จะมีคำตอบอยู่สองทางเลือก และปัญหาการเรียนรู้หลายประเภท ซึ่งมีคำตอบมีมากกว่าสองทางเลือก สำหรับปัญหาหลักของงานวิจัยนี้ คือ ปัญหาการเรียนรู้ออนไลน์แบบหลายประเภทที่ตัวอย่างมีน้ำหนักความสำคัญ โดยงานวิจัยนี้เสนอวิธีการสร้างอัลกอริทึมสำหรับปัญหาดังกล่าว 2 วิธีการ

วิธีแรก สร้างอัลกอริทึมสำหรับปัญหาการเรียนรู้แบบออนไลน์ที่มีน้ำหนักความสำคัญจากอัลกอริทึมแบบออนไลน์มาตรฐาน (นั่นคืออัลกอริทึมที่ใช้ในงานตัวอย่างไม่มีน้ำหนักความสำคัญ) วิธีนี้จะสามารถสร้างอัลกอริทึมสำหรับเซตของตัวอย่างและเซตของประเภทใดๆ ก็ได้ (ทำนายแบบสองประเภทหรือหลายประเภทก็ได้) ขึ้นกับอัลกอริทึมออนไลน์ที่ใช้เป็นฐาน แต่เซตของน้ำหนักของตัวอย่างจะต้องเป็นจำนวนเต็ม และได้พิสูจน์ประสิทธิภาพของอัลกอริทึมที่สร้างได้เทียบกับอัลกอริทึมฐาน เทคนิคที่ใช้ในการพิสูจน์คือการลดรูปการเรียนรู้ (learning reduction) กล่าวคือ เราได้พิสูจน์ว่าอัลกอริทึมที่ได้จะสามารถรับประกันความผิดพลาดได้ไม่แย่กว่า  $c_{\max}$  เท่าของความผิดพลาดของอัลกอริทึมฐาน เมื่อ  $c_{\max}$  คือค่าน้ำหนักที่มากที่สุดของข้อมูล

วิธีที่สอง เป็นการปรับปรุงอัลกอริทึมเพอร์เซพตรอน (Perceptron) เพื่อให้ทำงานกับข้อมูลที่มีน้ำหนักความสำคัญ เช่นเดียวกับอัลกอริทึมเพอร์เซพตรอน อัลกอริทึมนี้ทำงานได้กับตัวอย่างบน  $\mathbb{R}^n$  และสามารถทำนายแบบสองประเภท อย่างไรก็ตามในกรณีนี้ เราสามารถพิสูจน์ประสิทธิภาพการทำงานได้ดีกว่าการใช้วิธีแรก กล่าวคือ ถ้า  $c_{\max}$  และ  $c_{\min}$  คือค่าน้ำหนักมากที่สุดและน้อยที่สุดของตัวอย่าง (ไม่จำเป็นต้องเป็นจำนวนเต็มก็ได้) ให้  $\theta$  เป็นอัตราส่วน  $c_{\max}/c_{\min}$  เราสามารถแสดงได้ว่าความผิดพลาดที่ได้จะไม่เกิน  $\theta L + c_{\max} \theta^2 C + \theta^2 \sqrt{c_{\min} LC}$  เมื่อ  $L$  เป็นค่าความผิดพลาดแบบมีน้ำหนักความสำคัญของตัวทำนายเชิงเส้นใดๆ และ  $C$  เป็นพารามิเตอร์ที่ขึ้นกับความซับซ้อนของตัวทำนายเชิงเส้นนั้น ค่าความผิดพลาดนี้สามารถเทียบได้กับอัลกอริทึมเพอร์เซพตรอนในกรณีที่ตัวอย่างไม่มีน้ำหนักความสำคัญ (หรือกรณีที่  $c_{\max} = c_{\min} = 1$ ) ที่มีค่า  $L + C + \sqrt{LC}$  เนื่องจากค่าความผิดพลาดข้างต้นขึ้นกับค่า  $\theta$  เราจึงได้ทำการปรับปรุงวิธีการในการพิสูจน์ค่าดังกล่าวใหม่ ด้วยวิธีนี้เราสามารถพิสูจน์ประสิทธิภาพการทำงานได้ดีขึ้น นั่นคือ ค่าความผิดพลาดที่ได้จะไม่เกิน  $L + c_{\max} C + \sqrt{c_{\max} LC}$

อย่างไรก็ตามวิธีการที่สองที่นำเสนอใช้ได้กับกรณีตัวอย่างที่มีสองประเภท ในกรณีที่ตัวอย่างมีหลายประเภท งานวิจัยนี้สนใจการแก้ปัญหาด้วยวิธีการลดรูป กล่าวคือทำการลดรูปปัญหาการเรียนรู้แบบหลายประเภทไปเป็นปัญหาการเรียนรู้แบบสองประเภท โดยได้ทำการศึกษาวิธีการลดรูปสองวิธีการ ได้แก่ วิธีการหนึ่งต่อทั้งหมด (One-against-All) และวิธีการหนึ่งต่อหนึ่ง (One-against-One) และได้ทำการพิสูจน์ว่า ในกรณีที่ใช้อัลกอริทึมเพอร์เซพตรอนเป็นฐานการลดรูปด้วยวิธีการหนึ่งต่อทั้งหมดและวิธีการหนึ่งต่อหนึ่งนั้นสมมูลกัน

## วัตถุประสงค์

ศึกษาและปรับปรุงอัลกอริทึมการเรียนรู้แบบออนไลน์ สำหรับปัญหาการเรียนรู้แบบมี  
น้ำหนักความสำคัญ

## การตรวจเอกสาร

### นิยามพื้นฐาน

**นิยาม** การฝึกการเรียนรู้แบบมีการแนะนำ คือ คู่ลำดับ  $(K, Y, \ell)$  โดยที่  $K$  เป็นเซตของประเภท ที่ได้รับการแนะนำระหว่างขั้นตอนการฝึกสอน,  $Y$  เป็นเซตของประเภทที่ได้จากการทำนาย และ  $\ell: K \times Y \rightarrow [0, \infty)$  เป็นฟังก์ชันความสูญเสีย

**นิยาม** ปัญหาการเรียนรู้แบบมีการแนะนำ คือ คู่ลำดับ  $(D, X, T)$  โดยที่  $T = (K, Y, \ell)$  เป็น การฝึกการเรียนรู้แบบมีการแนะนำ,  $X$  เป็นเซตของลักษณะ (feature) และ  $D$  เป็นการกระจายตัวบน  $X \times K$

เป้าหมายของการแก้ปัญหการเรียนรู้แบบมีการแนะนำคือ การหาสมมติฐาน  $h: X \rightarrow Y$  ที่ทำให้ค่าเฉลี่ยของความสูญเสีย  $E_{(x,k) \sim D} \ell(k, h(x))$  น้อยที่สุด

**นิยาม** อัลกอริทึมการเรียนรู้แบบมีการแนะนำ สำหรับภารกิจ  $(K, Y, \ell)$  คือ กระบวนการ สร้างแผนที่สำหรับเซตจำกัดของกลุ่มตัวอย่าง  $(X \times K)^*$  ใดๆ ไปยังสมมติฐาน  $h: X \rightarrow Y$

เนื่องจากในงานวิจัยนี้สนใจเฉพาะการเรียนรู้แบบมีการแนะนำ ดังนั้นต่อไปจะเรียกการ เรียนรู้แบบมีการแนะนำ แบบย่อ ว่า “การเรียนรู้” และสมมติให้เซตของประเภทที่ได้รับการแนะนำ ระหว่างขั้นตอนการฝึกสอน  $K$  เท่ากับ เซตของประเภทที่ได้จากการทำนาย  $Y$

### การเรียนรู้แบบออฟไลน์

ให้  $X$  เป็นเซตของตัวอย่างทั้งหมดที่เป็นไปได้,  $Y$  เป็นเซตจำกัดของประเภทของตัวอย่าง ทั้งหมดที่เป็นไปได้,  $S$  เป็นเซตของกลุ่มตัวอย่างที่สุ่มโดยอิสระต่อกันจากการกระจายตัว  $D$  ซึ่งเป็น การกระจายตัวบน  $X \times Y$  สมมติฐานใดๆ คือฟังก์ชัน  $h: X \rightarrow Y$  ในการเรียนรู้หนึ่งๆ เราจะจำกัด การพิจารณาสมมติฐานให้อยู่ในเซตหนึ่งๆ เท่านั้น กล่าวคือ เราจะให้  $H$  เป็นเซตของสมมติฐาน ทั้งหมดที่สนใจ อัลกอริทึมการเรียนรู้แบบออฟไลน์ คืออัลกอริทึมที่รับกลุ่มตัวอย่าง  $S$  มาใช้สร้าง

สมมติฐาน  $h \in H$  ประสิทธิภาพของสมมติฐานดังกล่าวจะวัดโดยความผิดพลาดซึ่งมีนิยามเป็น  $\Pr_{(x,y) \sim D}[h(x) \neq y]$

จากนิยามข้างต้นจะเห็นว่าอัลกอริทึมการเรียนรู้แบบออฟไลน์จะมีเอาท์พุทเป็นสมมติฐาน หรือ *ตัวทำนาย* และตัวทำนายดังกล่าวสามารถนำไปใช้ในการทำนายประเภทตัวอย่างอื่นๆ ต่อไป

### การเรียนรู้แบบออนไลน์

ให้  $X$  เป็นเซตของตัวอย่างที่เป็นไปได้,  $Y$  เป็นเซตจำกัดของประเภทของตัวอย่างทั้งหมดที่เป็นไปได้ ในสถานการณ์การเรียนรู้แบบออนไลน์ อัลกอริทึมจะไม่ได้รับกลุ่มตัวอย่าง  $S$  ที่สมบูรณ์ก่อนเริ่มทำงาน หากแต่จะได้รับการตัวอย่างไปเรื่อยๆ กับการใช้งานจริง โดยมีการทำงานเป็นรอบๆ ดังนี้ เมื่อเริ่มต้นแต่ละรอบ อัลกอริทึมจะได้รับตัวอย่าง  $x \in X$  จากนั้นจะต้องทำนายประเภทของตัวอย่างนั้นทันที เมื่ออัลกอริทึมทำนายเสร็จ จะได้รับคำตอบ  $y \in Y$  เพื่อนำไปใช้ปรับปรุงความถูกต้องของการจำแนกในรอบถัดๆ ไป การวัดประสิทธิภาพของอัลกอริทึมจะวัดจากจำนวนครั้งที่ทั้งหมดที่ตอบผิด

โดยเนื้อแท้แล้ว อัลกอริทึมการเรียนรู้แบบออนไลน์ สามารถระบุโดยใช้ฟังก์ชัน  $f: H \times X \times Y \rightarrow H$  สำหรับปรับปรุงสมมติฐาน กล่าวคือในการทำงานแต่ละรอบอัลกอริทึมจะมีสมมติฐาน  $h_t$  สำหรับจำแนกประเภทของตัวอย่าง  $x_t \in X$  อัลกอริทึมจะตอบ  $\hat{y}_t = h_t(x_t)$  ถ้าอัลกอริทึมตอบผิด นั่นคือ  $\hat{y}_t \neq y_t$  อัลกอริทึมจะปรับปรุงสมมติฐานโดยใช้ฟังก์ชัน  $f$  กล่าวคือ  $h_{t+1} = f(h_t, x_t, y_t)$  ไม่เช่นนั้นจะให้  $h_{t+1} = h_t$

จากนิยามจะเห็นว่าในการเรียนรู้แบบออนไลน์จะมีสมมติฐานที่เปลี่ยนแปลงได้ตลอดกระบวนการการทำงาน ซึ่งต่างจากอัลกอริทึมการเรียนรู้แบบออฟไลน์ที่ใช้เพียงสมมติฐานเดียวในการตอบคำถาม

ในการพิจารณาการทำงานของอัลกอริทึมการเรียนรู้แบบออนไลน์ เราจะพิจารณาลำดับของข้อมูลป้อนเข้าและผลลัพธ์ของอัลกอริทึมบนลำดับป้อนเข้านั้นๆ กล่าวคือ จะเรียกลำดับ  $(x_1, y_1), (x_2, y_2) \dots (x_m, y_m)$  ว่า ลำดับข้อมูลป้อนเข้า ประสิทธิภาพการทำงานของอัลกอริทึมออนไลน์  $f$  ในลำดับดังกล่าวจะวัดโดยใช้ฟังก์ชันความสูญเสีย  $Loss: Y \times Y \rightarrow \mathbb{R}$  ซึ่งจะมีค่าเท่ากับ

$$\sum_{t=1}^m \text{Loss}(y_t, \hat{y}_t)$$

เมื่อ  $\hat{y}_t$  คือผลลัพธ์ของอัลกอริทึมในการทำงานรอบที่  $t$  โดยทั่วไปฟังก์ชันความสูญเสียที่นิยมใช้จำนวนครั้งที่ตอบผิดพลาด (mistake) คือ  $M(y, \hat{y}) = I[y \neq \hat{y}]$  เมื่อ  $I$  เป็นฟังก์ชันชี้วัด (indicator function)

### การเรียนรู้แบบออนไลน์แบบมีน้ำหนักความสำคัญ

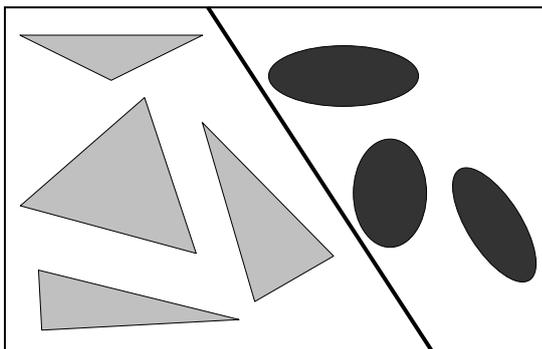
ในกรณีที่ตัวอย่างมีน้ำหนักความสำคัญ นั่นคือตัวอย่างจะระบุน้ำหนัก  $c$  จากเซต  $C \subseteq [0, \infty)$  อัลกอริทึมจะทราบน้ำหนักของตัวอย่างเมื่อได้ทำนายประเภทของตัวอย่างนั้นแล้ว อัลกอริทึมเรียนรู้ออนไลน์แบบมีน้ำหนักความสำคัญสามารถนิยามได้ในลักษณะเดียวกันกับอัลกอริทึมแบบออนไลน์ทั่วไป เพียงแต่ฟังก์ชัน  $f: H \times X \times Y \times C \rightarrow H$  จะรับน้ำหนักของตัวอย่างไปใช้ในการปรับสมมติฐานด้วย เรียกลำดับ  $(x_1, y_1, c_1), (x_2, y_2, c_2), \dots, (x_m, y_m, c_m)$  ว่า ลำดับข้อมูลป้อนเข้าแบบมีน้ำหนัก และจะวัดประสิทธิภาพของอัลกอริทึมออนไลน์  $f$  ในลำดับดังกล่าวเป็น

$$\sum_{t=1}^m c_t \cdot \text{Loss}(y_t, \hat{y}_t)$$

เมื่อ  $\hat{y}_t$  คือผลลัพธ์ของอัลกอริทึมในการทำงานรอบที่  $t$  และ  $\text{Loss}$  เป็นฟังก์ชันความสูญเสีย

### อัลกอริทึมเพอร์เซพตรอน (Perceptron Algorithm)

อัลกอริทึมเพอร์เซพตรอน Rosenblatt (1958) เป็นอัลกอริทึมการเรียนรู้ออนไลน์แบบหนึ่งที่ใช้ในการสร้างตัวทำนายเชิงเส้น สำหรับการทำนายประเภทของกลุ่มตัวอย่างแบบสองประเภท แนวคิดหลักของอัลกอริทึมคือพยายามหาเส้นตรงที่แบ่งแยกกลุ่มตัวอย่างทั้งสองประเภทออกจากกัน ดังแสดงในภาพที่ 1



ภาพที่ 1 ตัวอย่างการทำนายด้วยเพอร์เซพตรอน

อัลกอริทึมเพอร์เซพตรอนจะทำงานกับตัวทำนายเชิงเส้น ซึ่งนิยามโดยเวกเตอร์น้ำหนัก  $\mathbf{w} \in \mathcal{R}^n$  โดยอัลกอริทึมจะรับตัวอย่าง  $\mathbf{x}_i \in \mathcal{R}^n$  แล้วทำนาย  $\hat{y}_i = \text{sgn}\langle \mathbf{x}_i, \mathbf{w}_i \rangle$  เมื่อ  $\text{sgn}$  คือฟังก์ชันเครื่องหมาย จากนั้นจะได้รับเฉลย  $y_i$  ถ้าทำนายผิด  $\hat{y}_i \neq y_i$  อัลกอริทึมจะทำการปรับปรุงเวกเตอร์น้ำหนัก

$$\mathbf{w}_{i+1} \leftarrow \mathbf{w}_i + y_i \mathbf{x}_i$$

สำหรับการวิเคราะห์ประสิทธิภาพของอัลกอริทึมเพอร์เซพตรอนมีอยู่ 2 แนวทางได้แก่ การวิเคราะห์ในกรณีมีตัวทำนายเชิงเส้นใดๆ ที่สามารถแบ่งกลุ่มตัวอย่างทั้งสองกลุ่มออกจากกันได้ และการวิเคราะห์ในกรณีไม่มีตัวทำนายเชิงเส้นใดๆ ที่สามารถแบ่งกลุ่มตัวอย่างทั้งสองกลุ่มออกจากกันได้

**ทฤษฎีบทที่ 1** Novikoff (1962) สำหรับลำดับข้อมูลป้อนเข้า  $S = ((\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)) \in (\mathcal{R}^n \times \{-1, +1\})^m$  ใดๆ ถ้า  $R = \max_i \|\mathbf{x}_i\|_2$  และมีเวกเตอร์ขนาดหนึ่งหน่วย  $\mathbf{u}$  ( $\|\mathbf{u}\|_2 = 1$ ) ที่ทำให้สำหรับทุกๆ ตัวอย่างใน  $S$  มี  $y_i \cdot \langle \mathbf{u}, \mathbf{x}_i \rangle \geq \gamma$  แล้วจำนวนครั้งที่อัลกอริทึมเพอร์เซพตรอนทำนายผิดมีค่าไม่เกิน  $(R/\gamma)^2$

ทฤษฎีบทข้างต้นเป็นการวิเคราะห์เฉพาะกรณีมีตัวทำนายเชิงเส้นใดๆ ที่สามารถแบ่งกลุ่มตัวอย่างทั้งสองกลุ่มออกจากกันได้ ส่วนในกรณีไม่มีตัวทำนายเชิงเส้นใดๆ ที่สามารถแบ่งกลุ่มตัวอย่างทั้งสองกลุ่มออกจากกันได้ จะกล่าวถึงในส่วนถัดไป

## ประสิทธิภาพของตัวทำนายเชิงเส้น

ในการพิสูจน์ประสิทธิภาพของอัลกอริทึมการเรียนรู้แบบออนไลน์ นิยมพิสูจน์เทียบประสิทธิภาพกับสมมติฐาน  $h^*$  จากเซตของสมมติฐาน  $H'$  บางเซต ซึ่งไม่จำเป็นต้องเท่ากับ  $H$  เซตของสมมติฐานที่นิยมใช้ ในกรณีที่เซตของตัวอย่าง  $X = \mathcal{X}^n$  คือเซตของตัวจำแนกเชิงเส้น  $h^*$  ทั้งหมดที่นิยามด้วยเวกเตอร์  $\mathbf{u} \in \mathcal{X}^n$  กล่าวคือ  $h^*(\mathbf{x}) = \text{sgn}\langle \mathbf{x}, \mathbf{u} \rangle$  เมื่อ  $\text{sgn}$  คือฟังก์ชันเครื่องหมาย

ประสิทธิภาพของตัวจำแนกเชิงเส้นที่นิยามด้วยเวกเตอร์  $\mathbf{u}$  สามารถวัดได้หลายแบบ เช่น วัดโดยใช้จำนวนครั้งที่ตอบผิดพลาด ในที่นี้จะวัดโดยใช้ฟังก์ชันความสูญเสียแบบฮิง (Hinge loss) ที่มีนิยามเป็น

$$\ell(\mathbf{x}, y, \mathbf{u}) = \max\{0, 1 - y \cdot \langle \mathbf{u}, \mathbf{x} \rangle\} \quad (1)$$

ดังนั้นความสูญเสียแบบฮิง  $L(h^*)$  ทั้งหมดบนลำดับข้อมูลป้อนเข้าแบบไม่มีน้ำหนัก  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$  จะมีค่าเป็น  $\sum_{t=1}^m \ell(\mathbf{x}_t, y_t, \mathbf{u})$

ในกรณีที่ข้อมูลป้อนเข้ามีน้ำหนัก จะนิยามความสูญเสียทั้งหมดแบบฮิงที่มีน้ำหนัก  $L_w(h^*)$  ของตัวจำแนกเชิงเส้นที่นิยามด้วยเวกเตอร์  $\mathbf{u}$  บนลำดับ  $(\mathbf{x}_1, y_1, c_1), (\mathbf{x}_2, y_2, c_2), \dots, (\mathbf{x}_m, y_m, c_m)$  เป็น  $L_w(h^*) = \sum_{t=1}^m c_t \cdot \ell(\mathbf{x}_t, y_t, \mathbf{u})$

สำหรับการวิเคราะห์ประสิทธิภาพของอัลกอริทึมเพอร์เซพตรอนในกรณีไม่มีตัวทำนายเชิงเส้นใดๆ Gentile (2003) ได้ทำการวิเคราะห์ไว้ดังทฤษฎีบทต่อไปนี้

**ทฤษฎีบทที่ 2** Gentile (2003) สำหรับลำดับข้อมูลป้อนเข้า  $S = ((\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)) \in (\mathcal{X}^n \times \{-1, +1\})^m$  ใดๆ ถ้า  $R = \max_t \|\mathbf{x}_t\|_2$  และมีตัวทำนายเชิงเส้น  $h^*$  ที่นิยามด้วยเวกเตอร์  $\mathbf{u}$  โดยที่  $L(h^*) \leq L$  และให้  $C = R^2 \|\mathbf{u}\|_2^2$  แล้วจำนวนครั้งที่อัลกอริทึมเพอร์เซพตรอนทำนายผิดบน  $S$  มีค่าไม่เกิน  $L + C + \sqrt{LC}$

อันที่จริงแล้วทฤษฎีบทที่ 2 นี้เป็นรูปแบบเฉพาะของทฤษฎีบทที่ว่าด้วยขีดจำกัดบนของอัลกอริทึมพินอร์ม ( $p$ -Norm Algorithm) ทำนายคิดใน Gentile (2003) ซึ่งในบทพิสูจน์ดังกล่าวมีการวัดระยะห่างระหว่าง 2 เวกเตอร์ โดยใช้ฟังก์ชันดังนี้

$$d_r(\mathbf{u}, \mathbf{w}) = \frac{1}{2} \|\mathbf{u}\|_p^2 + \frac{1}{2} \|\mathbf{w}\|_q^2 - \langle \mathbf{u}, \mathbf{f}(\mathbf{w}) \rangle \quad (2)$$

เมื่อ  $\frac{1}{p} + \frac{1}{q} = 1$  และ ฟังก์ชัน  $\mathbf{f}: \mathfrak{R}^n \rightarrow \mathfrak{R}^n$  มีนิยามเป็น  $\mathbf{f} = (f_1, \dots, f_n)$  โดยที่

$$f_i(\mathbf{w}) = \frac{\text{sgn}(w_i) |w_i|^{q-1}}{\|\mathbf{w}\|_q^{q-2}}, \quad \mathbf{w} = (w_1, \dots, w_n) \in \mathfrak{R}^n$$

ฟังก์ชันระยะห่างตามสมการ (2) เรียกว่า การถู้ออกของเบรกแมน (Bregman divergence) สำหรับฟังก์ชัน  $\mathbf{f}$  มีอินเวอร์สฟังก์ชันเป็นฟังก์ชัน  $\mathbf{f}^{-1}: \mathfrak{R}^n \rightarrow \mathfrak{R}^n$  ซึ่งมีนิยามเป็น  $\mathbf{f}^{-1} = (f_1^{-1}, \dots, f_n^{-1})$  โดยที่

$$f_i^{-1}(\boldsymbol{\theta}) = \frac{\text{sgn}(\theta_i) |\theta_i|^{p-1}}{\|\boldsymbol{\theta}\|_p^{p-2}}, \quad \boldsymbol{\theta} = (\theta_1, \dots, \theta_n) \in \mathfrak{R}^n$$

สำหรับการพิสูจน์ทฤษฎีบทข้างต้นจะต้องใช้ทฤษฎีบทย่อย 2 ทฤษฎีบทดังนี้

ทฤษฎีบทย่อยที่ 1 ให้  $\mathbf{u}, \mathbf{w}, \mathbf{x} \in \mathfrak{R}^n$ ,  $a \in \mathfrak{R}$  และ  $\mathbf{w}' = \mathbf{f}^{-1}(\mathbf{f}(\mathbf{w}) + a\mathbf{x})$ ,  $\boldsymbol{\theta} = \mathbf{f}(\mathbf{w})$  แล้ว จะได้ว่า

$$a(\langle \mathbf{u}, \mathbf{x} \rangle - \langle \mathbf{w}, \mathbf{x} \rangle) = d_r(\mathbf{u}, \mathbf{w}) - d_r(\mathbf{u}, \mathbf{w}') + d_r(\mathbf{w}, \mathbf{w}')$$

ทฤษฎีบทย่อยที่ 2 ให้  $\mathbf{u}, \mathbf{w}, \mathbf{x} \in \mathfrak{R}^n$ ,  $a \in \mathfrak{R}$ ,  $p \geq 2$  และ  $\mathbf{w}' = \mathbf{f}^{-1}(\mathbf{f}(\mathbf{w}) + a\mathbf{x})$ ,  $\boldsymbol{\theta} = \mathbf{f}(\mathbf{w})$  แล้ว จะได้ว่า

$$d_r(\mathbf{w}, \mathbf{w}') \leq \frac{a^2}{2} (p-1) \|\mathbf{x}\|_2^2$$

สำหรับนิยามฟังก์ชันการลู่ออกของเบรกแมน, ทฤษฎีบทย่อยที่ 1 และทฤษฎีบทย่อยที่ 2 นั้นมีอยู่ใน Gentile (2003) ซึ่งในงานวิจัยนี้จะทำการปรับปรุงนิยามและทฤษฎีบทย่อยดังกล่าว เพื่อใช้สำหรับการพิสูจน์ทฤษฎีบทที่ว่าด้วยค่าจีคบังคับบนของความสูญเสียของอัลกอริทึมที่ได้นำเสนอ ซึ่งเป็นอัลกอริทึมที่ปรับปรุงมาจากเพอร์เซพตรอน

### การลดรูปการเรียนรู้

ในการแก้ปัญหาการเรียนรู้มีแนวทางการแก้ปัญหาอยู่ 2 แนวทาง แนวทางแรกคือ การแก้ปัญหาการเรียนรู้ต่างๆ โดยอิสระจากวิธีการอื่น หรือกล่าวคือพยายามแก้ปัญหาการเรียนรู้ต่างๆ โดยตรง ส่วนแนวทางที่สองคือ แก้ปัญหาโดยวิธีการแบบลดรูป นั่นคือการลดรูปจากปัญหาที่ต้องการจะแก้ไปเป็นปัญหาที่มีการแก้และวิเคราะห์ไปแล้ว ซึ่งมีข้อดีคือสามารถใช้อัลกอริทึมและทฤษฎีของปัญหาดังกล่าวกับปัญหาที่ต้องการจะแก้ได้ หรือกล่าวอีกนัยหนึ่ง การลดรูปการเรียนรู้  $R$  จากภารกิจ  $T = (K, Y, \ell)$  ไปยังภารกิจ  $T' = (K', Y', \ell')$  คือกระบวนการที่ใช้อัลกอริทึมการเรียนรู้สำหรับ  $T'$  เป็นกล่องดำในการแก้ภารกิจ  $T$  โดยรับประกันว่าถ้าแก้ปัญหาย่อยที่สร้างขึ้นโดย  $R$  ได้ดี แล้วจะแก้ปัญหาดังต้นได้ดีด้วย

นิยาม การลดรูปการเรียนรู้  $R(S, A)$  จากภารกิจ  $(K, Y, \ell)$  ไปยังภารกิจ  $(K', Y', \ell')$  คือกระบวนการที่รับข้อมูลที่ประกอบด้วย เซตจำกัดของข้อมูลป้อนเข้า  $S \in (X \times K)^*$  และ อัลกอริทึมการเรียนรู้  $A$  สำหรับภารกิจ  $(K', Y', \ell')$  และให้ผลลัพธ์เป็นสมมติฐาน  $h: X \rightarrow Y$  ซึ่งเป็นกลุ่มของสมมติฐานย่อยที่ได้จาก  $A$

### วิธีการหนึ่งต่อทั้งหมด (One-against-All Method)

วิธีการแบบหนึ่งต่อทั้งหมดเป็นวิธีการในการลดรูปการเรียนรู้จากปัญหาการเรียนรู้แบบหลายประเภทไปเป็นปัญหาการเรียนรู้แบบสองประเภทวิธีหนึ่ง ซึ่งจะทำการลดรูปการเรียนรู้  $k$  ประเภทไปเป็นการเรียนรู้สองประเภท สำหรับการเรียนรู้แบบหลายประเภท จะใช้เซตของตัวทำนายแบบสองประเภท  $B = \{b^r : r \in \{1, 2, \dots, k\}\}$  ซึ่งตัวทำนายแต่ละตัวถูกฝึกสอนโดยแผนที่  $(x, y) \rightarrow (x, I(y = r))$  จากตัวอย่างการเรียนรู้แบบหลายประเภทไปยังตัวอย่างการเรียนรู้แบบสองประเภท โดยวิธีการหนึ่งต่อทั้งหมดจะทำนาย  $r$  เมื่อ  $b^r(x) = 1$

เมื่อพิจารณาวิธีการหนึ่งต่อทั้งหมดที่กล่าวมาในข้างต้น จะเห็นว่าเวลาที่ใช้ในการทำนายแต่ละตัวอย่างมีค่าเป็น  $O(k)$  และวิธีการดังกล่าวสามารถสร้างตัวทำนายแบบหลายประเภทจากตัวทำนายแบบสองประเภทใดๆ ก็ได้ สำหรับอัลกอริทึมที่ใช้สร้างตัวทำนายแบบสองประเภทที่สนใจในงานวิจัยนี้คือ อัลกอริทึมเพอร์เซพตรอน ในการสร้างตัวทำนายแบบหลายประเภทโดยใช้เพอร์เซพตรอนเป็นฐานมีวิธีการหลักๆ ได้แก่ วิธีการหลายเวกเตอร์ และ วิธีการเวกเตอร์เดียว

### วิธีการหลายเวกเตอร์ (Multi-vector Method)

ให้  $Y$  เป็นเซตจำกัดของประเภทของตัวอย่างทั้งหมดที่เป็นไปได้ วิธีการหลายเวกเตอร์จะทำนายโดยใช้เซตของเวกเตอร์น้ำหนัก  $W = \{w^r : r \in Y\}$  ซึ่งประกอบด้วย  $|Y|$  เวกเตอร์ สำหรับตัวอย่าง  $\mathbf{x}_t$  จะทำนาย

$$\hat{y}_t = \arg \max_{r \in Y} \langle \mathbf{w}_t^r, \mathbf{x}_t \rangle$$

จากนั้นเมื่อได้รับเฉลย  $y_t$  ถ้าทำนายผิด  $\hat{y}_t \neq y_t$  จะทำการปรับปรุงเวกเตอร์น้ำหนัก

$$\mathbf{w}_{t+1}^{y_t} \leftarrow \mathbf{w}_t^{y_t} + \mathbf{x}_t, \mathbf{w}_{t+1}^{\hat{y}_t} \leftarrow \mathbf{w}_t^{\hat{y}_t} - \mathbf{x}_t$$

และ  $\mathbf{w}_{t+1}^r \leftarrow \mathbf{w}_t^r$  สำหรับทุกๆ  $r \in Y \setminus \{y_t, \hat{y}_t\}$

### วิธีการเวกเตอร์เดียว (Single-vector Method)

วิธีการเวกเตอร์เดียวจะทำนายโดยใช้เวกเตอร์น้ำหนัก  $\mathbf{w}_t$  เพียงตัวเดียวในการทำนายตัวอย่าง  $\mathbf{x}_t$  อย่างไรก็ตาม การทำนายด้วยวิธีการนี้จะต้องทำการสร้างแผนที่ตัวอย่าง ด้วยฟังก์ชันแผนที่  $\phi : \mathcal{R}^n \times Y \rightarrow \mathcal{R}^d$  นั่นคือ สำหรับตัวอย่าง  $\mathbf{x}_t$  จะทำนาย

$$\hat{y}_t = \arg \max_{r \in Y} \langle \mathbf{w}_t, \phi(\mathbf{x}_t, r) \rangle$$

จากนั้นเมื่อได้รับเฉลย  $y_t$  ถ้าทำนายผิด  $\hat{y}_t \neq y_t$  จะทำการปรับปรุงเวกเตอร์น้ำหนัก

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \phi(\mathbf{x}_t, y_t) - \phi(\mathbf{x}_t, \hat{y}_t)$$

สำหรับฟังก์ชันแผนที่ สามารถสร้างได้หลายวิธี ยกตัวอย่างวิธีง่ายๆ คือ ใช้ตัวอย่างแรกที่เจอในแต่ละประเภท เพื่อทำการสร้างแผนที่สำหรับตัวอย่าง  $\mathbf{x}_t$  และประเภทนั้นๆ หรือกล่าวคือ สำหรับแต่ละประเภท  $r \in Y$  ให้  $\mathbf{p}^r \in \mathcal{R}^n$  เป็นตัวอย่างแรกของประเภท  $r$  ที่เจอในลำดับข้อมูลป้อนเข้า, สามารถนิยามแผนที่  $\phi(\mathbf{x}_t, r)$  ซึ่งเป็นเวกเตอร์บน  $\mathcal{R}^n$  ได้เป็น

$$\phi_i(\mathbf{x}_t, r) = x_{t,i} p_i^r$$

เมื่อ  $\phi_i(\mathbf{x}_t, r)$ ,  $x_{t,i}$  และ  $p_i^r$  เป็นค่าในมิติที่  $i$  ของเวกเตอร์  $\phi(\mathbf{x}_t, r)$ ,  $\mathbf{x}_t$  และ  $\mathbf{p}^r$  ตามลำดับ

### วิธีการหนึ่งต่อหนึ่ง (One-against-One Method)

วิธีการแบบหนึ่งต่อหนึ่งเป็นวิธีการในการลดรูปการเรียนรู้การเรียนรู้แบบหลายประเภทไปเป็นภารกิจการเรียนรู้แบบสองประเภทอีกวิธีการหนึ่ง, ให้  $X$  เป็นเซตของตัวอย่างที่เป็นไปได้ และ  $Y = \{1, 2, \dots, k\}$  วิธีการแบบหนึ่งต่อหนึ่งจะทำนายโดยใช้เซตของตัวทำนายแบบสองประเภท  $B = \{b^{ij} : i \in Y, j \in Y, i \neq j\}$  โดยที่ตัวทำนาย  $b^{ij}$  เป็นฟังก์ชันจาก  $X$  ไปยัง  $\{i, j\}$  และ  $b^{ij}(x) = b^{ji}(x)$  วิธีการหนึ่งต่อหนึ่งจะทำนาย

$$\hat{y} = \arg \max_{r \in Y} \sum_{j \in Y: j \neq r} I[r = b^{ij}(x)]$$

เมื่อ  $I$  เป็นฟังก์ชันชี้วัด, พิจารณาวิธีการหนึ่งต่อหนึ่งที่กล่าวมาในข้างต้น จะเห็นว่าวิธีการนี้ใช้ตัวทำนายแบบสองประเภทเป็นจำนวน  $\binom{k}{2}$  ตัวทำนาย และใช้เวลาในการทำนายแต่ละตัวอย่างเป็น  $O(k^2)$  ซึ่งมากกว่าวิธีการแบบหนึ่งต่อทั้งหมด

อย่างไรก็ตามในกรณีการเรียนรู้แบบออฟไลน์การเปรียบเทียบประสิทธิภาพการทำนายของทั้งสองวิธีการยังไม่เป็นที่แน่ชัดว่าวิธีใดดีกว่ากัน และยังเป็นปัญหาเปิดที่สำคัญปัญหาหนึ่ง (Rifkin and Klautau (2004))

## งานวิจัยในอดีต

งานวิจัยชิ้นนี้เกี่ยวข้องกับงานวิจัยทางการเรียนรู้ด้วยเครื่องจักร 2 กลุ่มใหญ่ ๆ กล่าวคือ งานที่เกี่ยวข้องกับการเรียนรู้ที่ไวต่อค่าใช้จ่าย (Cost-sensitive learning) และงานที่เกี่ยวข้องกับการเรียนรู้แบบออนไลน์

*การเรียนรู้ที่ไวต่อค่าใช้จ่าย.* งานวิจัยด้านนี้มีมากมาย โดยเฉพาะในหลาย ๆ ปีที่ผ่านมา ได้มีความสนใจในปัญหาดังกล่าว เช่น Zadrozny *et al.* (2003); Elkan (2001); Abe *et al.* (2004); Domingos (1999); Liu (2006) โดยในการแก้ปัญหาที่มีหลายวิธี เช่น การพยายามปรับอัลกอริทึมหนึ่งๆ ให้ทำงานได้กับข้อมูลที่มีน้ำหนัก Liu (2006) หรืออาจจะเป็นการพยายามสร้างวิธีทั่วไปในการสร้างอัลกอริทึมสำหรับปัญหานี้ โดยอาจมีการประมาณการกระจายตัวของข้อมูล Zadrozny and Elkan (2001) หรืออาจจะใช้วิธีการลดรูปการเรียนรู้ไปยังปัญหาการทำนายแบบไม่มีน้ำหนัก Zadrozny *et al.* (2003) ซึ่งอัลกอริทึมแรกที่น่าเสนอในงานวิจัยนี้ก็ใช้แนวคิดของการลดรูปเช่นเดียวกัน

*การเรียนรู้แบบออนไลน์.* การเรียนรู้แบบออนไลน์เป็นสาขาที่มีการทำวิจัยอย่างกว้างขวางสำหรับเนื้อหาหลัก ดูหนังสือ Cesa-Bianchi and Lugosi (2006) อัลกอริทึมพื้นฐานของการเรียนรู้แบบออนไลน์มีหลายแบบ ที่ใช้ในงานวิจัยนี้คือเพอร์เซพตรอน Rosenblatt (1958) ซึ่ง Novikoff (1962) ได้ทำการวิเคราะห์ประสิทธิภาพของอัลกอริทึมนี้ในกรณีที่มีเส้นแบ่งระหว่างข้อมูลสองกลุ่มได้ ส่วนในงานวิจัยนี้เราใช้การวิเคราะห์ประสิทธิภาพในแบบ Gentile (2003) ที่สามารถพิจารณากรณีที่ไม่มีเส้นแบ่งระหว่างข้อมูลสองกลุ่มได้ อัลกอริทึมแบบออนไลน์สำหรับปัญหาการทำนายหลายประเภทที่ใช้เพอร์เซพตรอนมีมากมาย เช่น งานของ Fink *et al.* (2006) ที่สามารถใช้เป็นกล่องดำในการลดรูปในอัลกอริทึมแรกที่งานวิจัยนี้เสนอได้ด้วย

งานวิจัยของ Cammer *et al.* (2006) ได้นำเสนอและวิเคราะห์อัลกอริทึมที่เรียกว่า Passive-Aggressive สำหรับปัญหาการเรียนรู้แบบออนไลน์ อัลกอริทึมดังกล่าวสามารถจัดการกับปัญหาการเรียนรู้แบบหลายประเภทที่ไวต่อค่าใช้จ่ายได้ด้วย อย่างไรก็ตาม น้ำหนักความสำคัญของข้อมูลอัลกอริทึมนี้รองรับได้จะขึ้นกับประเภทของข้อมูลเท่านั้น ในขณะที่ในงานเช่นของ Zadrozny *et al.* (2003) และงานวิจัยนี้น้ำหนักความสำคัญจะขึ้นกับข้อมูลแต่ละตัว

## วิธีการ

งานวิจัยนำเสนอวิธีการในการแก้ปัญหาการเรียนรู้ที่ตัวอย่างมีน้ำหนักความสำคัญทั้งแบบสองประเภทและแบบหลายประเภท โดยใช้อัลกอริทึมการเรียนรู้แบบออนไลน์ สำหรับแนวคิดหลักของการแก้ปัญหาคือให้อัลกอริทึมการเรียนรู้แบบออนไลน์ทำการเรียนรู้ตามค่าน้ำหนักความสำคัญของแต่ละตัวอย่าง ในส่วนของการแก้ปัญหาการเรียนรู้สองประเภทแบบมีน้ำหนักความสำคัญได้นำเสนออัลกอริทึมการเรียนรู้แบบออนไลน์ 2 อัลกอริทึม โดยอัลกอริทึมแรกจะทำการทำซ้ำข้อมูลตามค่าน้ำหนักความสำคัญของตัวอย่าง ส่วนอีกอัลกอริทึมหนึ่งปรับปรุงมาจากอัลกอริทึมเพอร์เซพตรอนอัลกอริทึม กล่าวคือถ้าอัลกอริทึมทำนายผิดจะทำการปรับปรุงเวกเตอร์น้ำหนักโดยการถ่วงน้ำหนักของตัวอย่างด้วยค่าน้ำหนักความสำคัญของตัวอย่าง สำหรับการแก้ปัญหาการเรียนรู้หลายประเภทแบบมีน้ำหนักความสำคัญจะใช้วิธีการการลดรูปจากปัญหาการเรียนรู้หลายประเภทไปเป็นปัญหาการเรียนรู้สองประเภทโดยใช้วิธีการการลดรูปจากปัญหาการเรียนรู้หลายประเภทไปเป็นปัญหาการเรียนรู้สองประเภทโดยใช้อัลกอริทึมที่นำเสนอในส่วนของการแก้ปัญหาการเรียนรู้สองประเภทแบบมีน้ำหนักความสำคัญเป็นฐาน

### วิธีการแก้ปัญหาการเรียนรู้สองประเภทแบบมีน้ำหนักความสำคัญ

สำหรับวิธีการแก้ปัญหาการเรียนรู้สองประเภทแบบมีน้ำหนักความสำคัญเรานำเสนออัลกอริทึมการเรียนรู้แบบออนไลน์ 2 อัลกอริทึม ได้แก่อัลกอริทึมแบบลดรูปและอัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญ

#### อัลกอริทึมแบบลดรูป

ในหัวข้อนี้จะนำเสนอการลดรูปการเรียนรู้ในแบบออนไลน์ กล่าวคือเราจะทำการลดรูปปัญหาการเรียนรู้แบบมีน้ำหนักความสำคัญไปเป็นปัญหาการเรียนรู้แบบไม่มีน้ำหนักความสำคัญด้วยกระบวนการแบบออนไลน์ ซึ่งแนวคิดหลักของวิธีการการลดรูปการเรียนรู้แบบออนไลน์นี้คือ จะทำการแปลงลำดับข้อมูลป้อนเข้าแบบมีน้ำหนักไปเป็นลำดับข้อมูลป้อนเข้าแบบไม่มีน้ำหนัก และทำการทำนายตัวอย่างบนลำดับที่ได้จากการแปลงโดยใช้อัลกอริทึมการเรียนรู้แบบออนไลน์ใดๆ ที่ไม่สนใจค่าน้ำหนักความสำคัญของตัวอย่าง การลดรูปดังกล่าวทำให้ได้อัลกอริทึมลดรูปการเรียนรู้แบบออนไลน์ ซึ่งเราจะเรียกว่า อัลกอริทึมแบบลดรูป (Online-Reduction Algorithm) ซึ่งในที่นี้จะเขียนแทนด้วย  $A'$  สังเกตว่า อัลกอริทึมแบบลดรูป  $A'$  สามารถทำนายตัวอย่างแบบมีน้ำหนัก

ความสำคัญได้โดยเรียกใช้อัลกอริทึมการเรียนรู้แบบออนไลน์  $A$  นั่นคือ  $A$  เป็นเหมือนกล่องดำของ  $A''$

ภาพที่ 2 แสดงการทำงานของอัลกอริทึมแบบลดรูป  $A''$  ที่สร้างจากอัลกอริทึมการเรียนรู้ออนไลน์  $A$  ที่มีฟังก์ชัน  $f$  ในการปรับตัวทำนาย โดย  $h_t$  เป็นสมมติฐานที่  $A$  ใช้ในการทำนายรอบที่  $t$  จากอัลกอริทึมแบบลดรูป จะเห็นว่า  $A''$  จะเรียกใช้  $A$  ในการทำนายประเภทของ  $x_t$  จากนั้น  $A''$  จะได้รับเฉลย  $y_t$  และ  $c_t$  โดยสมมติให้  $c_t$  มีค่าเป็นจำนวนเต็มบวก (สามารถพิจารณาได้ว่า  $x_t, y_t$  และ  $c_t$  เป็นตัวอย่างหนึ่งบนลำดับข้อมูลป้อนเข้าแบบมีน้ำหนัก  $Q$ ) จากนั้น  $A''$  จะทำการทำซ้ำ  $x_t$  ออกมาอีก  $c_t$  ตัว แล้วเรียกใช้  $A$  เรียนรู้จากข้อมูลดังกล่าว ซึ่งเท่ากับว่า  $A''$  ทำการแปลงลำดับข้อมูลป้อนเข้าแบบมีน้ำหนัก  $Q$  ไปเป็นลำดับข้อมูลป้อนเข้าแบบไม่มีน้ำหนัก

**Algorithm Online-Reduction**

1. For  $t = 1, 2, 3, \dots$
2.     Receive an instance  $x_t$
3.     Predict  $\hat{y}_t \leftarrow h_t(x_t)$
4.     Receive correct label  $y_t$  and cost  $c_t$
5.      $h_{t+1} \leftarrow h_t$
6.     For  $i = 1$  to  $c_t$
7.          $x_{ti} = x_t$
8.         If  $h_{t+1}(x_{ti}) \neq y_t$

**ภาพที่ 2 อัลกอริทึมแบบลดรูป**

สังเกตว่า อัลกอริทึมแบบลดรูป  $A''$  สามารถใช้อัลกอริทึมการเรียนรู้แบบออนไลน์  $A$  ใดๆ เป็นกล่องดำก็ได้ ดังนั้น  $A''$  สามารถทำนายได้ในลักษณะเดียวกับอัลกอริทึมกล่องดำ  $A$  ที่นำมาใช้ ยกตัวอย่างเช่น ถ้า  $A$  เป็นอัลกอริทึมการเรียนรู้แบบหลายประเภท  $A''$  ก็จะสามารถทำนายแบบหลายประเภทด้วย

### อัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญ

ในหัวข้อนี้จะนำเสนอการปรับปรุงอัลกอริทึมเพอร์เซพตรอนเพื่อให้สามารถทำนายตัวอย่างที่มีน้ำหนักความสำคัญได้อย่างมีประสิทธิภาพ จากการทำงานของอัลกอริทึมเพอร์เซพตรอนที่กล่าวถึงในส่วนของ การตรวจเอกสาร จะเห็นว่าเพอร์เซพตรอนทำงานกับตัวอย่างซึ่งเวกเตอร์ ดังนั้นเราสามารถใช้นิวเคลียสหลักของงานวิจัยนี้ที่กล่าวไว้ว่า ให้อัลกอริทึมการเรียนรู้แบบออนไลน์ทำการเรียนรู้ตามค่าน้ำหนักความสำคัญของแต่ละตัวอย่าง โดยเราเสนอให้ทำการปรับปรุงอัลกอริทึมเพอร์เซพตรอนโดยใช้วิธีการปรับค่าเวกเตอร์น้ำหนักของเพอร์เซพตรอน ในรอบที่ทำนายผิดด้วยถ่วงน้ำหนักของตัวอย่างด้วยค่าน้ำหนักความสำคัญของตัวอย่าง โดยเราจะเรียกอัลกอริทึมที่ได้จากการปรับปรุงนี้ว่า อัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญ หรือเรียกโดยย่อว่า *IWP*

#### Algorithm *IWP*

1. INITIALIZE:  $\mathbf{w}_1 = 0$
2. For  $t = 1, 2, 3, \dots$
3.     Receive an instance  $\mathbf{x}_t$
4.     Predict  $\hat{y}_t \leftarrow \text{sgn}\langle \mathbf{w}_t, \mathbf{x}_t \rangle$
5.     Receive correct label  $y_t$  and cost  $c_t$
6.     If  $\hat{y}_t \neq y_t$
7.          $\mathbf{w}_{t+1} = \mathbf{w}_t + y_t c_t \mathbf{x}_t$

### ภาพที่ 3 อัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญ

ภาพที่ 3 แสดงการวิธีการทำงานของเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญ สังเกตว่าอัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญใช้ได้เฉพาะการทำนายแบบสองประเภท เช่นเดียวกับอัลกอริทึมเพอร์เซพตรอน ในหัวข้อถัดไปเราจะนำเสนอวิธีการที่ทำให้เพอร์เซพตรอนแบบมีน้ำหนักความสำคัญสามารถทำนายแบบหลายประเภทได้

## วิธีการแก้ปัญหาการเรียนรู้หลายประเภทแบบมีน้ำหนักความสำคัญ

ในการแก้ปัญหาการเรียนรู้หลายประเภทแบบมีน้ำหนักความสำคัญ เรานำเสนอให้ใช้วิธีการลดรูปจากปัญหาการเรียนรู้หลายประเภทไปเป็นปัญหาการเรียนรู้สองประเภท โดยใช้ อัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญที่ได้นำเสนอในส่วนที่ผ่านมาเป็นฐาน สำหรับวิธีการลดรูปจากปัญหาการทำนายการเรียนรู้หลายประเภทไปเป็นปัญหาการเรียนรู้สองประเภทที่เรานำมาใช้ ได้แก่ วิธีการหลายเวกเตอร์ และวิธีการทุกคู่

### วิธีการหลายเวกเตอร์ที่ใช้เพอร์เซพตรอนแบบมีน้ำหนักความสำคัญเป็นฐาน

วิธีการหลายเวกเตอร์เป็นวิธีการลดรูปจากปัญหาการเรียนรู้หลายประเภทไปเป็นปัญหาการทำนายสองประเภทซึ่งเป็นรูปแบบหนึ่งของวิธีการลดรูปแบบหนึ่งต่อทั้งหมด สำหรับการแก้ปัญหการเรียนรู้หลายประเภทแบบไม่มีน้ำหนักความสำคัญโดยใช้วิธีการหลายเวกเตอร์ที่ใช้เพอร์เซพตรอนเป็นฐาน เราได้กล่าวถึงแล้วในหัวข้อการตรวจเอกสาร ในส่วนนี้เราจะนำเสนอการปรับปรุงวิธีการดังกล่าวให้สามารถทำนายข้อมูลที่ตัวอย่างมีน้ำหนักความสำคัญได้ดี โดยเราจะเรียกวิธีการดังกล่าวว่า วิธีการหลายเวกเตอร์ที่ใช้เพอร์เซพตรอนแบบมีน้ำหนักความสำคัญเป็นฐาน (Multi-vector based on *IWP*) ซึ่งมีวิธีการดังต่อไปนี้

สำหรับปัญหาการเรียนรู้  $k$  ประเภท สมมติให้  $Y = \{1, 2, \dots, k\}$  เป็นเซตประเภทของตัวอย่าง วิธีการหลายเวกเตอร์ที่ใช้เพอร์เซพตรอนแบบมีน้ำหนักความสำคัญเป็นฐานจะทำนายโดยใช้เซตของเวกเตอร์น้ำหนัก  $W = \{\mathbf{w}^r : r \in Y\}$  ซึ่งประกอบด้วย  $k$  เวกเตอร์ สำหรับการทำนายรอบที่  $t$  ใดๆ จะใช้สมมติฐานการทำนาย  $h_t : \mathcal{X} \rightarrow \mathcal{Y}$  ในการคำนวณคะแนนของแต่ละประเภท  $r \in Y$  ซึ่งจะมีค่าคะแนนเป็น

$$h_t(\mathbf{x}_t, r) = \langle \mathbf{w}_t^r, \mathbf{x}_t \rangle \quad (3)$$

วิธีการหลายเวกเตอร์ที่ใช้เพอร์เซพตรอนแบบมีน้ำหนักความสำคัญเป็นฐานจะทำนาย

$$\hat{y}_t = \arg \max_{r \in Y} h_t(\mathbf{x}_t, r)$$

จากนั้นเมื่อได้รับเฉลย  $y_t$  ถ้าทำนายผิด  $\hat{y}_t \neq y_t$  จะทำการปรับปรุงเวกเตอร์น้ำหนัก

$$\mathbf{w}_{t+1}^{y_t} \leftarrow \mathbf{w}_t^{y_t} + c_t \mathbf{x}_t, \mathbf{w}_{t+1}^{\hat{y}_t} \leftarrow \mathbf{w}_t^{\hat{y}_t} - c_t \mathbf{x}_t$$

และ  $\mathbf{w}_{t+1}^r \leftarrow \mathbf{w}_t^r$  สำหรับทุกๆ  $r \in Y \setminus \{y_t, \hat{y}_t\}$

**วิธีการทุกคู่ที่ใช้เฟอร์เซพตรอนแบบมีน้ำหนักความสำคัญเป็นฐาน**

จากวิธีการหลายเวกเตอร์ซึ่งเป็นรูปแบบหนึ่งของวิธีการลดรูปจากปัญหาการเรียนรู้หลายประเภทไปเป็นปัญหาการเรียนรู้สองประเภทแบบหนึ่งต่อทั้งหมด เราได้ทำการนิยามวิธีการลดรูปแบบหนึ่งต่อหนึ่ง โดยใช้เฟอร์เซพตรอนแบบมีน้ำหนักความสำคัญเป็นฐาน เพื่อการแก้ปัญหาการเรียนรู้หลายประเภทแบบมีน้ำหนักความสำคัญ โดยเราจะเรียกวิธีการดังกล่าวว่า **วิธีการทุกคู่ที่ใช้เฟอร์เซพตรอนแบบมีน้ำหนักความสำคัญเป็นฐาน (All-pair based on IWP)** ซึ่งมีวิธีการดังต่อไปนี้

สำหรับปัญหาการเรียนรู้  $k$  ประเภท สมมติให้  $Y = \{1, 2, \dots, k\}$  เป็นเซตประเภทของตัวอย่าง วิธีการหลายเวกเตอร์ที่ใช้เฟอร์เซพตรอนแบบมีน้ำหนักความสำคัญเป็นฐานจะทำนายโดยใช้เซตของเวกเตอร์น้ำหนัก  $W = \{\mathbf{w}^{ij} : i \in Y, j \in Y, i \neq j\}$  ซึ่งประกอบด้วย  $k(k-1)$  นั่นคือ ในการทำนายรอบที่  $t$  สำหรับแต่ละคู่  $i, j \in Y$  จะมีเวกเตอร์น้ำหนัก  $\mathbf{w}_t^{ij}$  และ  $\mathbf{w}_t^{ji}$  เวกเตอร์ ดังที่แสดงในภาพที่ 4ก โดยที่  $\mathbf{w}_t^{ij} = -\mathbf{w}_t^{ji}$  ดังที่แสดงในภาพที่ 4ข

	$\mathbf{w}_t^{12}$	...	$\mathbf{w}_t^{1k}$
$\mathbf{w}_t^{21}$		...	$\mathbf{w}_t^{2k}$
$\vdots$	$\vdots$		...
$\mathbf{w}_t^{k1}$	$\mathbf{w}_t^{k2}$	...	

(ก)

	$\mathbf{w}_t^{12}$	...	$\mathbf{w}_t^{1k}$
$-\mathbf{w}_t^{12}$		...	$\mathbf{w}_t^{2k}$
$\vdots$	$\vdots$		...
$-\mathbf{w}_t^{1k}$	$-\mathbf{w}_t^{2k}$	...	

(ข)

ภาพที่ 4 เวกเตอร์น้ำหนักของวิธีการทุกคู่

ในการทำนายรอบที่  $t$  ใดๆ จะใช้สมมติฐานการทำนาย  $h_t : \mathfrak{X} \rightarrow \mathfrak{Y}$  ในการคำนวณคะแนนของแต่ละประเภท  $r \in Y$  ซึ่งจะมีค่าคะแนนเป็น

$$h_t(\mathbf{x}_t, r) = \sum_{j \in Y \setminus \{r\}} \langle \mathbf{w}_t^j, \mathbf{x}_t \rangle$$

นั่นคือ คะแนนของประเภท  $r$  จะเท่ากับผลรวมของคะแนนของประเภท  $r$  ที่ได้จากตัวทำนายแบบสองประเภทแต่ละตัว และจะทำนาย

$$\hat{y}_t = \arg \max_{r \in Y} h_t(\mathbf{x}_t, r)$$

จากนั้นเมื่อได้รับเฉลย  $y_t$  ถ้าทำนายผิด  $\hat{y}_t \neq y_t$  จะทำการปรับปรุงเวกเตอร์น้ำหนัก โดยสำหรับแต่ละ  $r \neq y_t$

$$\mathbf{w}_{t+1}^{y_t r} \leftarrow \mathbf{w}_t^{y_t r} + c_t \mathbf{x}_t$$

และ สำหรับแต่ละ  $r \neq \hat{y}_t$

$$\mathbf{w}_{t+1}^{\hat{y}_t r} \leftarrow \mathbf{w}_t^{\hat{y}_t r} - c_t \mathbf{x}_t$$

จากการทำงานของทั้งสองวิธีการที่กล่าวมา จะเห็นว่าวิธีการหลายเวกเตอร์ที่ใช้เพอร์เซพตรอนแบบมีน้ำหนักความสำคัญเป็นฐานใช้เวลาในการทำงานเป็น  $O(k)$  น้อยกว่าวิธีการทุกคู่ที่ใช้เพอร์เซพตรอนแบบมีน้ำหนักความสำคัญเป็นฐานซึ่งใช้เวลาในการทำงานเป็น  $O(k^2)$

## ผลการวิจัย

สำหรับผลการวิจัยเราจะแบ่งออกเป็น 2 ส่วน ได้แก่ ผลการวิจัยการแก้ปัญหาการเรียนรู้ออกสองประเภทแบบมีน้ำหนักความสำคัญ และผลการวิจัยการแก้ปัญหาการเรียนรู้ออกสองประเภทแบบมีน้ำหนักความสำคัญ

### ผลการวิจัยการแก้ปัญหาการเรียนรู้ออกสองประเภทแบบมีน้ำหนักความสำคัญ

สำหรับผลการวิจัยการแก้ปัญหาการเรียนรู้ออกสองประเภทแบบมีน้ำหนักความสำคัญแบ่งออกเป็น 2 ส่วน ได้แก่ ผลการวิจัยภาคทดลองและผลการวิจัยภาคทฤษฎี

#### ผลการวิจัยภาคทดลอง

เราได้ทำการทดลองเพื่อเปรียบเทียบประสิทธิภาพการทำนายของอัลกอริทึมที่ได้นำเสนอไปในส่วนวิธีการ ได้แก่ อัลกอริทึมแบบลดรูปที่ใช้เพอร์เซพตรอนเป็นกล่องดำ และอัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญ (JWP) รวมทั้งอัลกอริทึมเพอร์เซพตรอน โดยทดลองกับชุดข้อมูลแบบสองประเภท 5 ชุดข้อมูลดังนี้

#### ชุดข้อมูล KDD Cup 1998 Data

ชุดข้อมูล KDD Cup 1998 Data นำมาจาก Asuncion and Newman (2007) ชุดข้อมูลนี้จะประกอบไปด้วยข้อมูล 95412 รายการ ซึ่งเป็นรายละเอียดของบุคคลที่บริจาคและไม่บริจาคเงินเพื่อการกุศลรวมถึงจำนวนเงินที่ได้บริจาค จุดประสงค์ของการใช้ข้อมูลชุดนี้คือ เพื่อใช้ในการทดสอบอัลกอริทึมที่ใช้ในการตัดสินใจว่าจะต้องส่ง หรือ ไม่ส่งจดหมายขอรับบริจาคไปให้กับบุคคลใดบ้าง โดยเราจะให้กลุ่มบุคคลที่บริจาคจัดอยู่ในประเภทบวก และ กลุ่มบุคคลที่ไม่บริจาคจัดอยู่ในประเภทลบ สำหรับประเภทบวกค่าใช้จ่ายเมื่ออัลกอริทึมตอบผิดหรือค่าน้ำหนักความสำคัญ จะคิดจากยอดเงินที่บริจาคหักค่าแสตมป์ 0.68 ส่วนประเภทลบถ้าตอบผิดจะต้องเสียค่าแสตมป์ 0.68 ในการทดลองทำการสุ่มข้อมูลร้อยละ 10 ของรายการทั้งหมดมาใช้ และเลือกคุณลักษณะ (Features) จากคุณลักษณะที่มีค่าความสัมพันธ์ (Correlation) สูงสุด 10 อันดับแรกมาใช้ โดยกำหนดให้ค่าความสัมพันธ์ของคุณลักษณะที่  $i$  มีค่าดังนี้

$$\text{Correlation}^{(i)} = \sum_{t=1}^m x_t^i y_t$$

โดย  $x_t^i$  คือค่าของคุณลักษณะที่  $i$  ของตัวอย่าง  $x_t$

### ชุดข้อมูล Auto-MPG

ชุดข้อมูล Auto-MPG นำมาจาก Data Asuncion and Newman (2007) ชุดข้อมูลประกอบไปด้วยข้อมูล 398 รายการ แต่ละรายการมี 8 คุณลักษณะ ซึ่งเป็นข้อมูลเกี่ยวกับการใช้เชื้อเพลิงของรถยนต์แต่ละคันซึ่งมีหน่วยเป็นไมล์ต่อแกลลอน (MPG) โดยปกติชุดข้อมูลนี้ใช้กับงาน Regression นั่นคือการทำนายการใช้เชื้อเพลิงของรถยนต์แต่ละคัน แต่ในการทดลองนี้นำมาใช้กับงานจำแนกประเภท โดยใช้ข้อมูล 392 รายการ (เนื่องจากอีก 6 รายการมีคุณลักษณะไม่ครบ) สำหรับตัวอย่าง  $x_t$  กำหนดให้ประเภท  $y_t$  และ น้ำหนักความสำคัญ  $c_t$  มีค่าดังนี้

$$y_t = \text{sgn}(mpg^{(t)} - \overline{mpg}) \text{ และ } c_t = |mpg^{(t)} - \overline{mpg}|$$

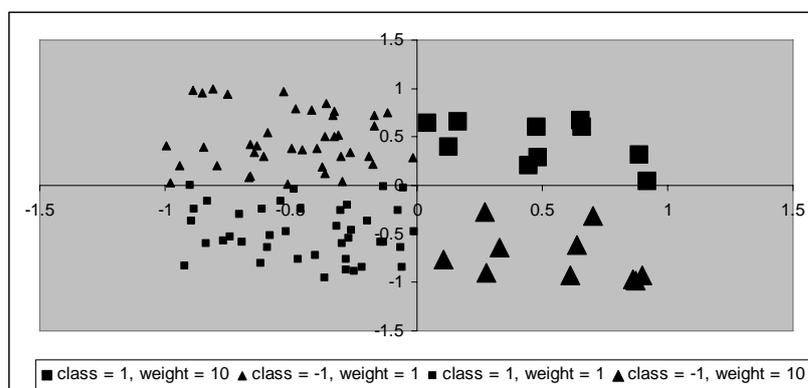
โดยที่  $mpg^{(t)}$  คือค่าคุณลักษณะ  $mpg$  ของตัวอย่าง  $x_t$  และ  $\overline{mpg}$  คือ ค่าเฉลี่ยของคุณลักษณะ  $mpg$

### ชุดข้อมูลสังเคราะห์

เพื่อให้เปรียบเทียบประสิทธิภาพของวิธีการต่างๆ เมื่อทำนายตัวอย่างที่มีน้ำหนักความสำคัญได้อย่างชัดเจน เราทำการสร้างชุดข้อมูลสังเคราะห์ที่มีการกระจายของน้ำหนักความสำคัญในชุดข้อมูลขัดแย้งกับการกระจายของประเภท นั่นคือสังเคราะห์ให้ข้อมูลที่มีน้ำหนักความสำคัญมากมีจำนวนตัวอย่างน้อย ส่วนข้อมูลที่มีน้ำหนักความสำคัญน้อยมีจำนวนตัวอย่างมาก กล่าวคือเราต้องการให้อัลกอริทึมที่สนใจน้ำหนักความสำคัญของตัวอย่าง มีจำนวนครั้งในการทำนายผิดมาก แต่เสียค่าใช้จ่ายรวมในการทำนายผิดน้อย ส่วนอัลกอริทึมที่ไม่ได้สนใจน้ำหนักความสำคัญของตัวอย่างเสียค่าใช้จ่ายรวมในการทำนายผิดมาก แต่มีจำนวนครั้งในการทำนายผิดน้อย เราทำการสังเคราะห์ข้อมูล 3 ชุดข้อมูล

### ข้อมูลสังเคราะห์ชุดที่ 1

ข้อมูลสังเคราะห์ชุดที่ 1 ประกอบไปด้วย 10000 ตัวอย่าง แบ่งเป็นข้อมูล 2 ประเภท ได้แก่ ประเภท -1 และประเภท +1 แต่ละประเภทมีข้อมูล 5000 ตัวอย่าง ซึ่งแบ่งเป็นตัวอย่างที่มีน้ำหนักที่มีน้ำหนักความสำคัญ 10 จำนวน 1000 ตัวอย่าง และตัวอย่างที่มีน้ำหนักที่มีน้ำหนักความสำคัญ 1 อีก 4000 ตัวอย่าง โดยมีการกระจายตัวดังแสดงในภาพที่ 5



ภาพที่ 5 การกระจายตัวของข้อมูลสังเคราะห์ชุดที่ 1

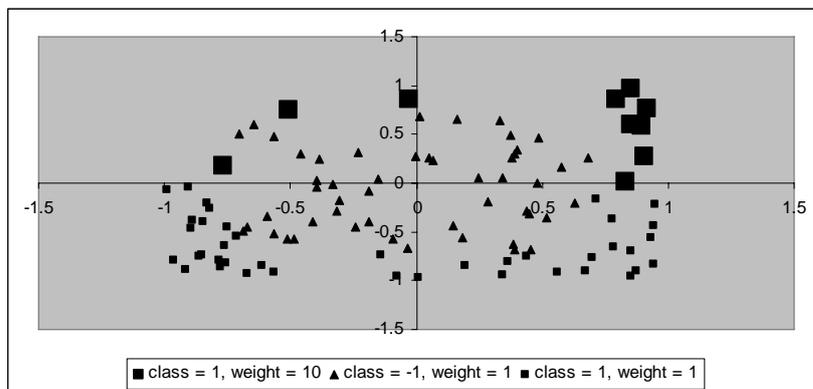
### ข้อมูลสังเคราะห์ชุดที่ 2

ข้อมูลสังเคราะห์ชุดที่ 2 ประกอบไปด้วยข้อมูลประกอบไปด้วย 10000 ตัวอย่าง แบ่งเป็น 2 ประเภท ได้แก่ประเภท -1 และประเภท +1 แต่ละประเภทมีข้อมูล 5000 ตัวอย่าง สำหรับประเภท +1 ประกอบด้วยตัวอย่างที่มีน้ำหนักที่มีน้ำหนักความสำคัญ 10 จำนวน 1000 ตัวอย่าง และตัวอย่างที่มีน้ำหนักที่มีน้ำหนักความสำคัญ 1 อีก 4000 ตัวอย่าง ส่วนประเภท -1 แต่ละตัวอย่างมีน้ำหนักความสำคัญเป็น 1 โดยมีการกระจายตัวดังแสดงในภาพที่ 6

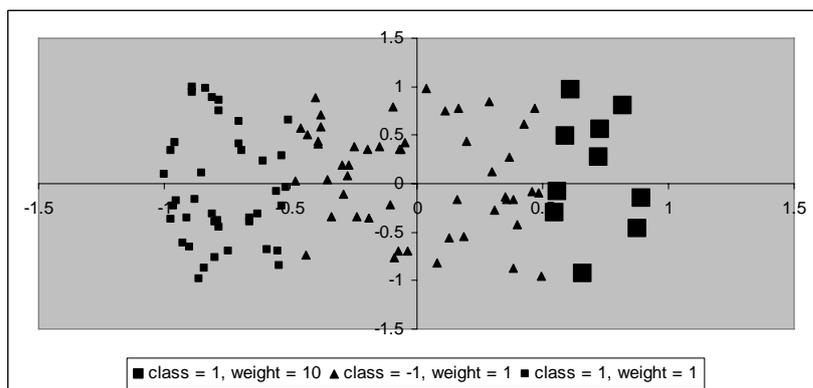
### ข้อมูลสังเคราะห์ชุดที่ 3

ข้อมูลสังเคราะห์ชุดที่ 3 ประกอบไปด้วยข้อมูลประกอบไปด้วย 10000 ตัวอย่าง แบ่งเป็น 2 ประเภท ได้แก่ประเภท -1 และประเภท +1 แต่ละประเภทมีข้อมูล 5000 ตัวอย่าง สำหรับประเภท +1 ประกอบด้วยตัวอย่างที่มีน้ำหนักที่มีน้ำหนักความสำคัญ 10 จำนวน 1000 ตัวอย่าง และตัวอย่างที่มี

น้ำหนักที่มีน้ำหนักความสำคัญ 1 อีก 4000 ตัวอย่าง ส่วนประเภท -1 แต่ละตัวอย่างมีน้ำหนักความสำคัญเป็น 1 โดยมีการกระจายตัวดังแสดงในภาพที่ 7



ภาพที่ 6 การกระจายตัวของข้อมูลสังเคราะห์ชุดที่ 2



ภาพที่ 7 การกระจายตัวของข้อมูลสังเคราะห์ชุดที่ 3

#### ผลการทดลอง

การเปรียบเทียบประสิทธิภาพการทำนายของ 3 อัลกอริทึม ได้แก่ อัลกอริทึมแบบลดรูปที่ใช้เพอร์เซพตรอนเป็นกลองคำ, อัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญ (*IWP*) และอัลกอริทึมเพอร์เซพตรอน สำหรับแต่ละชุดข้อมูลที่กล่าวมาในข้างต้น เราทำโดยการนับจำนวนครั้งและคิดค่าใช้จ่ายรวมที่แต่ละวิธีการทำนายผิดได้ผลดังที่แสดงในตารางที่ 1 เมื่อพิจารณาร้อยละของจำนวนที่แต่ละวิธีการ ทำนายผิด (%ตอบผิด) จะเห็นว่าโดยส่วนใหญ่อัลกอริทึมเพอร์เซพตรอนมี

จำนวนครั้งในการทำนายผิดใกล้เคียงกับอัลกอริทึมแบบลดรูปซึ่งน้อยกว่าอัลกอริทึม *IWP* แต่เมื่อพิจารณาร้อยละของค่าใช้จ่ายรวมในแต่ละวิธีการทำนายผิด (%ค่าใช้จ่าย) จะเห็นว่าอัลกอริทึม *IWP* กลับมีค่าใช้จ่ายรวมในการทำนายผิดน้อยกว่าอีกสองอัลกอริทึมที่ให้ผลการทดลองใกล้เคียงกัน จากผลการทดลองที่กล่าวมานี้ เราสามารถกล่าวได้ว่า อัลกอริทึม *IWP* ทำนายข้อมูลที่ตัวอย่างมีน้ำหนักความสำคัญได้ดีกว่าอัลกอริทึมแบบลดรูป และอัลกอริทึมเพอร์เซพตรอน

**ตารางที่ 1** ผลการเปรียบเทียบประสิทธิภาพการทำนายของอัลกอริทึมแบบลดรูป, อัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญ และอัลกอริทึมเพอร์เซพตรอน กับชุดข้อมูลแบบสองประเภท

ชุดข้อมูล	<i>IWP</i>		ลดรูป		เพอร์เซพตรอน	
	%ตอบผิด	%ค่าใช้จ่าย	%ตอบผิด	%ค่าใช้จ่าย	%ตอบผิด	%ค่าใช้จ่าย
KDD-98	10.33	54.90	52.66	51.18	<b>9.59</b>	54.92
Auto-Mpg	23.21	17.00	<b>20.15</b>	<b>13.02</b>	26.60	16.37
สังเคราะห์ที่1	35.70	56.25	61.31	<b>44.04</b>	<b>33.71</b>	57.79
สังเคราะห์ที่2	43.01	56.06	57.17	<b>40.03</b>	<b>41.81</b>	57.37
สังเคราะห์ที่3	40.64	56.31	56.80	<b>39.00</b>	<b>38.27</b>	58.57

สังเกตว่า อัลกอริทึมแบบลดรูปมีประสิทธิภาพการทำนายไม่แตกต่างจากอัลกอริทึมเพอร์เซพตรอนเท่าใดนัก ซึ่งประเด็นนี้เราจะทำการวิเคราะห์ในส่วนผลการวิจัยภาคทฤษฎี

### ผลการวิจัยภาคทฤษฎี

ในส่วนของการวิจัยภาคทฤษฎีเราจะทำการวิเคราะห์ประสิทธิภาพการทำนายของอัลกอริทึมแบบลดรูปที่ใช้เพอร์เซพตรอนเป็นกล่องดำ, อัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญ (*IWP*) และอัลกอริทึมเพอร์เซพตรอน

### การวัดประสิทธิภาพของอัลกอริทึมการเรียนรู้แบบออนไลน์

โดยทั่วไปการระบุประสิทธิภาพของอัลกอริทึมการเรียนรู้ออนไลน์  $A$  ใดๆ จะระบุโดยเทียบกับสมมติฐานจากบางเซตสมมติฐาน  $H$  กล่าวคือ ให้  $X$  เป็นเซตของตัวอย่างทั้งหมดที่เป็นไปได้ และ  $Y$  เป็นเซตของประเภทของตัวอย่างทั้งหมดที่เป็นไปได้ สำหรับบางฟังก์ชัน  $\ell$  ที่ระบุความสูญเสียของ  $h$  บนลำดับของข้อมูล  $S = ((x_1, y_1), \dots, (x_m, y_m)) \in (X \times Y)^m$  ใดๆ และฟังก์ชัน  $Loss$  ที่ระบุความสูญเสียของ  $A$  มีฟังก์ชัน  $g: \mathcal{H} \rightarrow \mathcal{H}$  ที่ถ้ามีสมมติฐาน  $h \in H$  ที่ทำนายได้ดีใน  $S$  นั่นคือ  $\ell(h, S) \leq L$  ค่าความสูญเสีย  $Loss(A, S)$  ของอัลกอริทึม  $A$  บนลำดับ  $S$  จะไม่เกิน  $g(L)$  โดยทั่วไปแล้ว ฟังก์ชัน  $\ell$  จะถูกนิยามให้เป็นผลรวมของค่าความสูญเสียบนข้อมูลแต่ละตัว กล่าวคือ  $\ell(h, S) = \sum_{i=1}^m \ell(h, (x_i, y_i))$  ในงานวิจัยนี้จะทำการวิเคราะห์โดยใช้ความจริงดังกล่าวด้วย

ในงานวิจัยนี้เราจะทำการวิเคราะห์ประสิทธิภาพของอัลกอริทึม โดยใช้ค่าความสูญเสียแบบมีน้ำหนัก  $Loss^w$  บนลำดับข้อมูลป้อนเข้าแบบมีน้ำหนัก  $Q = ((x_1, y_1, c_1), \dots, (x_m, y_m, c_m)) \in (X \times Y \times [0, \infty))^m$  ใดๆ เมื่อเทียบกับฟังก์ชันความสูญเสียแบบมีน้ำหนัก  $\ell^w(h, Q)$  ของ  $h \in H$  บน  $Q$  โดยที่  $Loss^w$  มีนิยามเป็นความสูญเสีย  $Loss$  คูณด้วยน้ำหนักของตัวอย่าง และ ฟังก์ชันความสูญเสียแบบมีน้ำหนักมีนิยามเป็น

$$\ell^w(h, x_i) = c_i \ell(h, x_i) \quad (4)$$

ต่อไปเราจะทำการวิเคราะห์ค่าความสูญเสียแบบมีน้ำหนักความสำคัญของแต่ละอัลกอริทึม

### การรับประกันค่าขีดจำกัดบนของความสูญเสียแบบมีน้ำหนักของอัลกอริทึมแบบลดรูป

จากการทำงานของอัลกอริทึมแบบลดรูปที่ได้นำเสนอในหัวข้อวิธีการ เราสามารถพิสูจน์ค่าความสูญเสียแบบมีน้ำหนัก  $Loss^w(A^w, Q)$  ของอัลกอริทึมแบบลดรูปซึ่งในที่นี้จะเขียนแทนด้วย  $A^w$  บนลำดับข้อมูลป้อนเข้าแบบมีน้ำหนัก  $Q$  ใดๆ เมื่อเทียบกับค่าความสูญเสียแบบมีน้ำหนักของตัวทำนาย  $h \in H$  ใดๆ ได้ด้วยการคิดเทียบค่าความสูญเสียของอัลกอริทึมการเรียนรู้แบบออนไลน์  $A$  ที่  $A^w$  ใช้เป็นกล่องดำ กับค่าความสูญเสียของ  $h$

ให้ฟังก์ชัน  $g: \mathfrak{R} \rightarrow \mathfrak{R}$  เป็นฟังก์ชันที่ใช้ระบุขีดจำกัดบนของความสูญเสียของ  $A$  และสมมติให้สำหรับตัวอย่างใดๆ บนลำดับ  $Q$  มีค่าเป็นจำนวนเต็มบวก เราสามารถรับประกันค่าขีดจำกัดบนของความสูญเสียแบบมีน้ำหนักของวิธีการแบบลดรูป  $A^w$  ได้ดังนี้

**ทฤษฎีบทที่ 3** สำหรับลำดับข้อมูลป้อนเข้าแบบมีน้ำหนัก  $Q = ((x_1, y_1, c_1), (x_2, y_2, c_2), \dots, (x_m, y_m, c_m)) \in (X \times Y \times I^+)^m$  ใดๆ ถ้ามีตัวทำนาย  $h \in H$  ที่มี  $\ell^w(h, Q) \leq L$  เราจะได้ว่า

$$\text{Loss}^w(A^w, Q) \leq c_{\max} \cdot g(L)$$

$$\text{เมื่อ } c_{\max} = \max_{t \in \{1, \dots, m\}} c_t$$

**บทพิสูจน์** พิจารณาลำดับลำดับข้อมูลป้อนเข้าแบบไม่มีน้ำหนักความสำคัญ

$S = (x_{11}, y_1), \dots, (x_{1c_1}, y_1), \dots, (x_{m1}, y_m), \dots, (x_{mc_m}, y_m)$  ที่  $A^w$  สร้างจาก  $Q$  เรียกลำดับย่อย  $(x_{t1}, y_t), \dots, (x_{tc_t}, y_t)$  บน  $S$  ว่า ลำดับย่อย  $X_t$  สังเกตว่าการทำงานของ  $A^w$  ในการปรับค่าตัวทำนายบน  $Q$  จะเหมือนกับการทำงานของ  $A$  บน  $S$

ในขั้นแรกเราจะเทียบ  $\ell(h, S)$  กับ  $\ell^w(h, Q)$  จะเห็นว่าถ้าตัวทำนาย  $h$  ทำนายตัวอย่างแรกบน  $X_t$  ผิด แล้ว  $h$  จะทำนายตัวอย่างอื่นๆ บน  $X_t$  และตัวอย่าง  $x_t$  บน  $Q$  ผิดด้วย และถ้า  $h$  ทำนายตัวอย่างแรกบน  $X_t$  ถูก แล้ว  $h$  จะทำนายตัวอย่างอื่นๆ บน  $X_t$  และตัวอย่าง  $x_t$  บน  $Q$  ถูกด้วย ดังนั้นจะได้ว่า

$$\ell(h, X_t) = \sum_{i=1}^{c_t} \ell(h, x_i) = c_t \ell(h, x_t) = \ell^w(h, x_t)$$

$$\text{นั่นคือ } \ell(h, S) = \ell^w(h, Q)$$

ถัดมาเราจะเทียบ  $\text{Loss}(A, S)$  กับ  $\text{Loss}^w(A^w, Q)$  พิจารณาตัวอย่างแบบมีน้ำหนัก  $(x_t, y_t, c_t)$  และลำดับย่อย  $X_t$  เนื่องจาก  $A^w$  ใช้คำตอบของ  $h_t(x_{t1})$  เป็นคำตอบของการทำนายตัวอย่าง  $x_t$  บน  $Q$  พิจารณากรณีที่  $A^w$  ทำนายตัวอย่าง  $x_t$  ผิด เราจะได้  $\text{Loss}^w(A^w, x_t)$  จะมีค่าเป็น  $c_t$  อย่างไรก็ตามสังเกตว่าถ้า  $A^w$  ทำนายผิดบนตัวอย่างดังกล่าว เราจะได้ว่า  $A$  ทำนายผิดในตัวอย่าง  $x_{t1}$  ด้วยเสมอ นั่นคือ  $\text{Loss}(A, X_t) \geq 1$  ดังนั้น จะได้

$$\text{Loss}^w(A^w, Q) \leq \sum_{t=1}^m c_t \cdot \text{Loss}(A, X_t) \leq c_{\max} \cdot \text{Loss}(A, S)$$

โดยการแทนค่าขีดจำกัดของความสูญเสียของ  $A$  เราจะได้ว่าทฤษฎีบทถูกต้อง เนื่องจาก  $\text{Loss}^w(A^w, Q) \leq c_{\max} \cdot \text{Loss}(A, S) \leq c_{\max} \cdot g(L)$  ■

### การรับประกันค่าขีดจำกัดบนของความสูญเสียแบบมีน้ำหนักของอัลกอริทึมการเรียนรู้แบบออนไลน์ใดๆ

ในส่วนนี้เราจะทำการวิเคราะห์ประเด็นสำคัญในส่วนของการวิจัยภาคทดลอง นั่นคือ เหตุใดอัลกอริทึมแบบลรูปซึ่งเรียนรู้โดยการทำซ้ำตัวอย่างที่ทำนายผิดตามค่าน้ำหนักความสำคัญของตัวอย่งนั้น จึงมีประสิทธิภาพการทำนายใกล้เคียงกับ อัลกอริทึมเพอร์เซพตรอนซึ่งไม่ได้สนใจน้ำหนักความสำคัญของตัวอย่าง โดยเราจะทำการวิเคราะห์ค่าขีดจำกัดบนของความผิดพลาดของอัลกอริทึมเพอร์เซพตรอน ซึ่งในที่นี้เราจะวิเคราะห์ในกรณีทั่วไป นั่นคือเราจะทำการวิเคราะห์ค่าขีดจำกัดบนของความสูญเสียแบบมีน้ำหนักของอัลกอริทึมการเรียนรู้แบบออนไลน์ใดๆ ที่ไม่สนใจค่าน้ำหนักความสำคัญของตัวอย่าง

สำหรับอัลกอริทึมการเรียนรู้แบบออนไลน์  $A$  ใดๆ ที่ไม่สนใจค่าน้ำหนักความสำคัญของตัวอย่าง ให้  $g: \mathcal{H} \rightarrow \mathcal{H}$  เป็นฟังก์ชันขีดจำกัดความสูญเสียของ  $A$  เมื่อเทียบกับตัวทำนาย  $h \in H$  ใดๆ และสำหรับลำดับข้อมูลป้อนเข้าแบบมีน้ำหนัก  $Q$  ใดๆ สมมติให้สำหรับตัวอย่างใดๆ บนลำดับ  $Q$  มีค่าเป็นจำนวนเต็มบวก เราสามารถรับประกันค่าขีดจำกัดบนของความผิดพลาดของ  $A$  ได้ดังนี้

**ทฤษฎีบทที่ 4** สำหรับลำดับข้อมูลป้อนเข้าแบบมีน้ำหนัก  $Q = ((x_1, y_1, c_1), (x_2, y_2, c_2), \dots, (x_m, y_m, c_m)) \in (X \times Y \times I^+)^m$  ใดๆ ถ้ามีตัวทำนาย  $h \in H$  ที่มี  $\ell^w(h, Q) \leq L$  จะได้ว่า

$$\text{Loss}^w(A, Q) \leq c_{\max} \cdot g(L)$$

เมื่อ  $c_{\max} = \max_{t \in \{1, \dots, m\}} c_t$

**บทพิสูจน์** เนื่องจาก  $A$  เป็นอัลกอริทึมที่ไม่สนใจน้ำหนักความสำคัญของตัวอย่าง ดังนั้นจะให้  $S = (x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$  เป็นลำดับข้อมูลป้อนเข้าแบบไม่มีน้ำหนัก ซึ่งเกิดจากลำดับข้อมูลป้อนเข้าแบบมีน้ำหนัก  $Q$  ที่ตัดค่าน้ำหนักความสำคัญของแต่ละตัวอย่างออกไป

ในขั้นแรกเราจะเทียบ  $\ell(h, S)$  กับ  $\ell^w(h, Q)$  เมื่อพิจารณาการทำนายของ  $h$  บน  $S$  และบน  $Q$  จะเห็นว่าถ้า  $h$  ทำนายตัวอย่าง  $(x_p, y_p)$  บน  $S$  ผิด แล้ว  $h$  จะทำนาย  $(x_p, y_p, c_p)$  บน  $Q$  ผิดด้วย และเนื่องจากสำหรับตัวอย่าง  $(x_p, y_p, c_p)$  ใดๆ มีค่า  $c_p$  เป็นจำนวนเต็มบวก ดังนั้นจะได้ว่า

$$\ell^w(h, Q) = \sum_{i=1}^m c_i \ell(h, x_i) \geq \ell(h, S)$$

ถัดมาเราจะเทียบ  $Loss(A, S)$  กับ  $Loss^w(A, Q)$  พิจารณาการทำนายของ  $A$  บน  $S$  และบน  $Q$  จะเห็นว่าถ้า  $A$  ทำนายตัวอย่าง  $(x_p, y_p)$  บน  $S$  ผิด แล้ว  $A$  จะทำนาย  $(x_p, y_p, c_p)$  บน  $Q$  ผิดด้วย นั่นคือ

$$Loss^w(A, Q) = \sum_{i=1}^m c_i \cdot Loss(A, x_i) \leq c_{\max} \cdot Loss(A, S)$$

เมื่อแทนค่าขีดจำกัดของความสูญเสียของ  $A$  เราจะได้ว่าทฤษฎีบทถูกต้อง เนื่องจาก  $Loss^w(A^w, Q) \leq c_{\max} \cdot Loss(A, S) \leq c_{\max} \cdot g(L)$  ■

จากทฤษฎีบทที่ 4 จะเห็นว่าค่าขีดจำกัดบนของความสูญเสียแบบมีน้ำหนักของอัลกอริทึมการเรียนรู้แบบออนไลน์ใดๆ ที่ไม่สนใจค่าน้ำหนักความสำคัญของตัวอย่างมีค่าเท่ากับค่าขีดจำกัดบนของของความสูญเสียแบบมีน้ำหนักของอัลกอริทึมแบบลดรูปในทฤษฎีบทที่ 3 ทำให้เราสามารถสรุปประเด็นที่กล่าวถึงในตอนต้นได้ว่า ไม่ว่าจะนำอัลกอริทึมการเรียนรู้แบบออนไลน์  $A$  ใดๆ ที่ไม่สนใจค่าน้ำหนักความสำคัญของตัวอย่าง ประสิทธิภาพการทำนายอย่างแน่นนั้น ไม่ได้แย่ไปกว่าประสิทธิภาพการทำนายของอัลกอริทึมแบบลดรูปที่ใช้  $A$  เป็นกล่องดำ

**การรับประกันค่าขีดจำกัดบนของความสูญเสียแบบมีน้ำหนักของอัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญแบบง่าย**

จากการทำงานของอัลกอริทึม  $IWP$  ที่ได้นำเสนอไปในหัวข้อวิธีการ จะเห็นว่า  $IWP$  ทำงานใกล้เคียงกับอัลกอริทึมเพอร์เซพตรอน แต่ปรับปรุงในส่วนของการปรับเวกเตอร์น้ำหนักในรอบที่มี

การทำนายผิด ซึ่งจะปรับด้วย  $c\mathbf{x}$  แทนที่จะเป็น  $\mathbf{x}$  ดังเช่นเพอร์เซพตรอน ดังนั้นสำหรับลำดับป้อนเข้าแบบมีน้ำหนัก  $Q = ((\mathbf{x}_1, y_1, c_1), (\mathbf{x}_2, y_2, c_2), \dots, (\mathbf{x}_m, y_m, c_m))$  สามารถพิจารณาได้ว่า  $IWP$  เป็นเพอร์เซพตรอน ที่มีลำดับข้อมูลป้อนเข้าแบบไม่มีน้ำหนักความสำคัญ  $S = (c_1\mathbf{x}_1, c_2\mathbf{x}_2, \dots, c_m\mathbf{x}_m)$  ซึ่งทำให้สามารถพิสูจน์ประสิทธิภาพของ  $IWP$  ได้ โดยใช้ทฤษฎีบทที่ 5 เป็นกล่องดำในการพิสูจน์สำหรับทฤษฎีบทที่ 5 นั้นเป็นรูปแบบเฉพาะของทฤษฎีบทที่ว่าด้วยขีดจำกัดบนของความสำเร็จอัลกอริทึมพินอร์ม ใน Gentile (2003) ซึ่งทฤษฎีบทดังกล่าว ให้ขีดจำกัดบนของความสำเร็จได้ทั้งกรณีที่มีและไม่มีตัวทำนายเชิงเส้นที่แบ่งกลุ่มตัวอย่างสองกลุ่มออกจากกันได้ การวัดประสิทธิภาพแบบนี้จะคิดเทียบกับสมมติฐาน  $h$  ซึ่งนิยามโดยเวกเตอร์  $\mathbf{u}$  โดยใช้ฟังก์ชันความเบี่ยงเบนซึ่งนิยามเป็น

$$D_\gamma(\mathbf{u}; (\mathbf{x}_i, y_i)) = \max\{0, \gamma - y_i \langle \mathbf{u}, \mathbf{x}_i \rangle\} \quad (5)$$

โดย  $\gamma$  เป็นค่าคงที่ใดๆ โดยที่  $\gamma > 0$  สำหรับฟังก์ชันความเบี่ยงเบนตามสมการ (5) นี้จะเป็นรูปแบบทั่วไปของฟังก์ชันความสูญเสียแบบฮิงก์ตามสมการ (1) ซึ่งได้กล่าวถึงในหัวข้อการตรวจเอกสาร เมื่อ  $\gamma = 1$

**ทฤษฎีบทที่ 5** Gentile (2003) ให้  $S = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)) \in (R^n \times \{1, -1\})^m$  เป็นลำดับข้อมูลป้อนเข้าของเพอร์เซพตรอน,  $M$  เป็นเซตของรอบที่เพอร์เซพตรอนตอบผิด และ  $R = \max_{i \in M} \|\mathbf{x}_i\|_2$  สำหรับเวกเตอร์  $\mathbf{u}$  ใดๆ ที่มี  $C = R^2 \|\mathbf{u}\|_2^2$  และ  $L = \sum_{i \in M} D_\gamma(\mathbf{u}; (\mathbf{x}_i, y_i))$  ถ้า  $\mathbf{w}_1 = 0$  จำนวนรอบที่เพอร์เซพตรอนตอบผิด  $|M| \leq \frac{L}{\gamma} + \frac{C}{\gamma^2} + \frac{1}{\gamma^2} \sqrt{\gamma LC}$

จากทฤษฎีบทที่ 5 จะเห็นว่าค่าความเบี่ยงเบนของตัวทำนายเชิงเส้น  $\mathbf{u}$  จะคิดเฉพาะครั้งที่เพอร์เซพตรอนตอบผิดเท่านั้น และจำนวนรอบที่เพอร์เซพตรอนตอบผิดก็จะคิดเทียบกับค่าดังกล่าว ในกรณีเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญ ประสิทธิภาพของอัลกอริทึมจะคิดเทียบกับสมมติฐาน  $h$  ซึ่งนิยามโดยเวกเตอร์  $\mathbf{u}$  จากสมการ (5) จะได้ว่า ค่าความเบี่ยงเบนของ  $h$  แบบมีน้ำหนักสำหรับตัวอย่าง  $\mathbf{x}_i$  จะมีค่าเป็น

$$D_c(\mathbf{u}; (c_i\mathbf{x}_i, y_i)) = \max\{0, c_i - y_i \cdot \langle c_i\mathbf{x}_i, \mathbf{u} \rangle\} \quad (6)$$

**ทฤษฎีบทที่ 6** สำหรับลำดับป้อนเข้าแบบมีน้ำหนัก  $Q = ((\mathbf{x}_1, y_1, c_1), (\mathbf{x}_2, y_2, c_2), \dots, (\mathbf{x}_m, y_m, c_m)) \in (\mathcal{R}^n \times \{-1, +1\} \times [0, \infty))^m$  ถ้ามีตัวทำนาย  $h \in H$  ใดๆ ที่นิยามโดย  $\mathbf{u}$  ที่มี  $\ell^w(h, Q) \leq L$  แล้ว จะได้ว่า

$$\text{Loss}^w(IWP, Q) \leq \theta L + c_{\max} \theta^2 C + \theta^2 \sqrt{c_{\min} L C}$$

$$\text{เมื่อ } C = \|\mathbf{u}\|_2^2 \cdot \max_{t \in \{1, \dots, m\}} \|\mathbf{x}_t\|_2^2, \quad c_{\max} = \max_{t \in \{1, \dots, m\}} c_t, \quad c_{\min} = \min_{t \in \{1, \dots, m\}} c_t \quad \text{และ} \quad \theta = \frac{c_{\max}}{c_{\min}}$$

**บทพิสูจน์** ให้ลำดับข้อมูลป้อนเข้าแบบไม่มีน้ำหนัก  $S = (c_1 \mathbf{x}_1, \dots, c_m \mathbf{x}_m)$  เป็นลำดับข้อมูลป้อนเข้าของเพอร์เซพตรอนและ ให้  $M$  แทนเซตของรอบที่  $IWP$  ตอบผิดพลาด พิจารณาตัวทำนาย  $h \in H$  ใดๆ ที่นิยามโดยเวกเตอร์  $\mathbf{u}$  ให้

$$L' = \sum_{t \in M} D_{c_{\min}}(u; (c_t \mathbf{x}_t, y_t))$$

แทนความเบี่ยงเบนที่ระดับ  $\gamma = c_{\min}$  ของ  $h$  จากทฤษฎีบทที่ 5 จะได้ว่า

$$|M| \leq \frac{L'}{c_{\min}} + \frac{C'}{c_{\min}^2} + \frac{1}{c_{\min}^2} \sqrt{c_{\min} L' C'}$$

$$\text{เมื่อ } C' = c_{\max}^2 \cdot \|\mathbf{u}\|_2^2 \cdot \max_{t \in M} \|\mathbf{x}_t\|_2^2$$

อย่างไรก็ตาม แต่ละครั้งที่  $IWP$  ตอบผิดบนตัวอย่าง  $\mathbf{x}_t$  จะมีค่าความสูญเสียเกิดขึ้น ไม่เกิน  $c_{\max}$  นั่นคือ

$$\text{Loss}(IWP, Q) \leq c_{\max} \left[ \frac{L'}{c_{\min}} + \frac{C'}{c_{\min}^2} + \frac{1}{c_{\min}^2} \sqrt{c_{\min} L' C'} \right]$$

ทฤษฎีบทเป็นจริง โดยสังเกตว่า

$$L \geq \ell^w(h, Q)$$

$$\begin{aligned}
&= \sum_{t=1}^m D_{c_t}(\mathbf{u}; (c_t \mathbf{x}_t, y_t)) && \text{จากสมการ (6)} \\
&\geq \sum_{t \in M} D_{c_t}(\mathbf{u}; (c_t \mathbf{x}_t, y_t)) && \blacksquare
\end{aligned}$$

ถ้าใช้อัลกอริทึมแบบลดรูปทำนายบนลำดับ  $Q$  โดยใช้เพอร์เซพตรอนเป็นกล่องค่า จะได้ขีดจำกัดบนของความสูญเสียแบบมีน้ำหนักเป็น  $c_{\max}(L + C + \sqrt{LC})$  สังเกตว่าอัลกอริทึม *IWP* จะให้ผลลัพธ์ที่ดีกว่า ในกรณีที่  $c_{\max}$  มีค่าใกล้เคียงกับ  $c_{\min}$  อย่างไรก็ตาม เราสามารถปรับปรุ่ค่าขีดจำกัดบนของความสูญเสียแบบมีน้ำหนักของ *IWP* ได้อีกดังจะกล่าวถึงในส่วนถัดไป

**การปรับปรุ่ค่าขีดจำกัดบนของความสูญเสียแบบมีน้ำหนักของอัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญ**

จากทฤษฎีบทที่ 6 จะเห็นว่าค่าขีดจำกัดบนของความสูญเสียแบบมีน้ำหนักของอัลกอริทึม *IWP* ขึ้นอยู่กับค่า  $c_{\max}/c_{\min}$  ซึ่งเราได้ทำการปรับปรุ่ค่าขีดจำกัดบนดังกล่าวและสามารถกำจัด  $c_{\max}/c_{\min}$  ออกไปได้ โดยการปรับปรุ่บทพิสูจน์ของทฤษฎีบทที่ว่าด้วยขีดจำกัดบนของความสูญเสียอัลกอริทึมพินอร์ม ใน Gentile (2003) ซึ่งได้ผลออกมาดังทฤษฎีบทต่อไปนี้

**ทฤษฎีบทที่ 7** สำหรับลำดับข้อมูลป้อนเข้าแบบมีน้ำหนักความ  $Q = ((\mathbf{x}_1, y_1, c_1), (\mathbf{x}_2, y_2, c_2), \dots, (\mathbf{x}_m, y_m, c_m)) \in (\mathcal{X}^n \times \{-1, +1\} \times [0, \infty))^m$  ใดๆ ถ้า  $R = \max_t \|\mathbf{x}_t\|_2$  และมีตัวทำนายเชิงเส้น  $h \in H$  ใดๆ ที่นิยามด้วยเวกเตอร์  $\mathbf{u}$  โดยที่  $\ell^w(h, Q) \leq L$  และให้  $C = R^2 \|\mathbf{u}\|_2^2$  แล้วจะได้ว่า

$$\text{Loss}^w(IWP, Q) \leq L + c_{\max} C + \sqrt{c_{\max} LC}$$

เมื่อ  $c_{\max} = \max_{t \in \{1, \dots, m\}} c_t$

สำหรับการพิสูจน์ทฤษฎีบทที่ 7 เราจะพิสูจน์ผ่านทางอีกทฤษฎีบทหนึ่งซึ่งเป็นการฉีกท้าวไปของทฤษฎีบทดังกล่าวโดยจะใช้วิธีการพิสูจน์ในแนวทางเดียวกับบทพิสูจน์ใน Gentile (2003) โดยทำการนิยามฟังก์ชันความเบี่ยงเบนแบบมีน้ำหนักใหม่ดังนี้

สำหรับตัวอย่าง  $(\mathbf{x}, y, c)$  ใดๆ, เวกเตอร์  $\mathbf{u}$  ใดๆ และค่าคงที่  $\gamma > 0$  ใดๆ ค่าความเบี่ยงเบนแบบมีน้ำหนักความสำคัญ

$$D_\gamma(\mathbf{u}; (\mathbf{x}, y, c)) = c(0, \gamma - y\langle \mathbf{u}, \mathbf{x} \rangle)_+ \quad (7)$$

ฟังก์ชันความเบี่ยงเบนแบบมีน้ำหนักในสมการ (7) นี้ เป็นกรณีทั่วไปของฟังก์ชันความเบี่ยงเบนแบบมีน้ำหนักในสมการ (6)

ในการพิสูจน์ทฤษฎีบทที่จะกล่าวต่อไป เราจะวัดค่าระยะห่างของเวกเตอร์  $\mathbf{u}$  และ  $\mathbf{w}$  ใดๆ ดังนี้

$$d(\mathbf{u}, \mathbf{w}) = \frac{1}{2} \|\mathbf{u}\|_2^2 + \frac{1}{2} \|\mathbf{w}\|_2^2 - \langle \mathbf{u}, \mathbf{w} \rangle \quad (8)$$

ฟังก์ชันนี้เป็นรูปแบบเฉพาะของการลู่ออกของเบรกแมนในสมการ (2) ซึ่งมี  $p = 2$  และใช้ทฤษฎีบทย่อย 2 ทฤษฎีบทดังนี้

**ทฤษฎีบทย่อยที่ 3** ให้  $\mathbf{u}, \mathbf{w}, \mathbf{x} \in \mathfrak{R}^n$ ,  $a \in \mathfrak{R}$  และ  $\mathbf{w}' = \mathbf{w} + a\mathbf{x}$  แล้ว จะได้ว่า

$$a(\langle \mathbf{u}, \mathbf{x} \rangle - \langle \mathbf{w}, \mathbf{x} \rangle) = d(\mathbf{u}, \mathbf{w}) - d(\mathbf{u}, \mathbf{w}') + d(\mathbf{w}, \mathbf{w}')$$

**ทฤษฎีบทย่อยที่ 4** ให้  $\mathbf{w}, \mathbf{x} \in \mathfrak{R}^n$ ,  $a \in \mathfrak{R}$  และ  $\mathbf{w}' = \mathbf{w} + a\mathbf{x}$  แล้ว จะได้ว่า

$$d(\mathbf{w}, \mathbf{w}') \leq \frac{a^2}{2} \|\mathbf{x}\|_2^2$$

สังเกตว่าทฤษฎีบทย่อยทั้งสองข้างต้นเป็นรูปแบบเฉพาะ เมื่อ  $p = q = 2$  ของทฤษฎีบทย่อยที่ 1 และทฤษฎีบทย่อยที่ 2 ซึ่งกล่าวไว้ในส่วนการตรวจเอกสารตามลำดับ ต่อไปเราจะทำการพิสูจน์ทฤษฎีบทซึ่งเป็นกรณีทั่วไปของทฤษฎีบทที่ 7 โดยในทฤษฎีบทนี้เราจะทำการพิสูจน์ค่าขีดจำกัดบนของความความสูญเสียแบบมีน้ำหนักของอัลกอริทึม *IWP* ที่มีเวกเตอร์น้ำหนักเริ่มต้น  $\mathbf{w}_1 \in \mathfrak{R}^n$  และ

สำหรับตัวอย่าง  $\mathbf{x}_t$  ใดๆ มีน้ำหนักความสำคัญ  $0 \leq c_t \leq 1$  โดยเราจะเรียกอัลกอริทึมดังกล่าวนี้ว่า  $IWP'$  ในการพิสูจน์เราจะใช้ฟังก์ชันความเวียนเบนแบบมีน้ำหนักดังที่นิยามในสมการ (7)

**ทฤษฎีบทที่ 8** สำหรับลำดับข้อมูลป้อนเข้าแบบมีน้ำหนัก  $Q = ((\mathbf{x}_1, y_1, c_1), (\mathbf{x}_2, y_2, c_2), \dots, (\mathbf{x}_m, y_m, c_m)) \in (\mathbb{R}^n \times \{-1, +1\} \times [0, 1])^m$  ใดๆ, ให้  $M$  เป็นเซตของรอบที่  $IWP'$  ทำนายผิดบน  $Q$  ถ้า  $R = \max_{t \in M} \|\mathbf{x}_t\|_2$  และสำหรับตัวทำนายเชิงเส้น  $h \in H$  ใดๆ ที่นิยามด้วยเวกเตอร์  $\mathbf{u} \in \mathbb{R}^n$  แล้วจะได้ว่า

$$\begin{aligned} \text{Loss}^w(IWP', Q) &\leq \frac{1}{\gamma} L_{\mathbf{u}, \gamma}(M) - \frac{\langle \mathbf{u}, \mathbf{w}_1 \rangle}{\gamma} + \frac{C}{2\gamma^2} \\ &\quad + \frac{1}{2\gamma^2} \sqrt{4\gamma C L_{\mathbf{u}, \gamma}(M) + C^2 - 4\gamma C \langle \mathbf{u}, \mathbf{w}_1 \rangle + 4\gamma^2 \|\mathbf{u}\|_2^2 \|\mathbf{w}_1\|_2^2} \end{aligned}$$

เมื่อ  $L_{\mathbf{u}, \gamma}(M) = \sum_{t \in M} D_\gamma(\mathbf{u}; (\mathbf{x}_t, y_t, c_t))$  และ  $C = R^2 \|\mathbf{u}\|_2^2$

**บทพิสูจน์** พิจารณาพจน์ในเครื่องหมายรากที่สองเราสามารถเขียนใหม่ได้เป็น

$$4\gamma C L_{\mathbf{u}, \gamma}(M) + (C - 2\gamma \langle \mathbf{u}, \mathbf{w}_1 \rangle)^2 + 4\gamma^2 (\|\mathbf{u}\|_2^2 \|\mathbf{w}_1\|_2^2 - (\langle \mathbf{u}, \mathbf{w}_1 \rangle)^2)$$

เนื่องจาก  $\|\mathbf{u}\|_2 \|\mathbf{w}_1\|_2 \geq \langle \mathbf{u}, \mathbf{w}_1 \rangle$  จึงทำให้ค่าข้างต้นมีค่าอย่างต่ำเท่ากับศูนย์ สังเกตว่า ถ้า

$\text{Loss}^w(IWP', Q)$  มีค่าไม่เกิน  $\frac{1}{\gamma} L_{\mathbf{u}, \gamma}(M) - \frac{c_{\max} \langle \mathbf{u}, \mathbf{w}_1 \rangle}{\gamma}$  แล้ว ทฤษฎีบทจะเป็นจริง เพราะฉะนั้นใน

การพิสูจน์ต่อจากนี้เราจะสมมติให้  $\text{Loss}^w(IWP', Q) > \frac{1}{\gamma} L_{\mathbf{u}, \gamma}(M) - \frac{c_{\max} \langle \mathbf{u}, \mathbf{w}_1 \rangle}{\gamma}$

สำหรับการทำนายรอบที่  $t \in M$  ใดๆ เวกเตอร์น้ำหนัก  $\mathbf{w}_t$  จะถูกปรับปรุงเมื่อใช้ทฤษฎีบทย่อยที่ 3 ในรอบนั้นๆ โดยให้  $a = c_t \gamma$  และแทน ด้วย  $\mathbf{u}$  ด้วย  $\lambda \mathbf{u}$  เมื่อ  $\lambda > 0$  จะได้ว่า

$$c_t y_t (\langle \lambda \mathbf{u}, \mathbf{x}_t \rangle - \langle \mathbf{w}_t, \mathbf{x}_t \rangle) = d(\lambda \mathbf{u}, \mathbf{w}_t) - d(\lambda \mathbf{u}, \mathbf{w}_{t+1}) + d(\mathbf{w}_t, \mathbf{w}_{t+1}) \quad (9)$$

สังเกตว่า ค่าของสมการข้างต้นแปรผันตรงกับค่าความแตกต่างของค่าทำนายของ  $h$  และ  $IWP'$  กล่าวคือ หากค่าดังกล่าวมีค่าน้อย แสดงว่า  $h$  ทำนายไม่ต่างจาก  $IWP'$  มากนักในรอบที่  $IWP'$

ทำนายผิด ดังนั้น เราต้องการให้ค่าดังกล่าวมีค่าน้อยๆ, จากสมการ (9) เนื่องจาก  $y_t \cdot \langle \mathbf{w}_t, \mathbf{x}_t \rangle \leq 0$  ดังนั้นจะได้ว่า

$$\lambda c_t y_t \langle \mathbf{u}, \mathbf{x}_t \rangle \leq d(\lambda \mathbf{u}, \mathbf{w}_t) - d(\lambda \mathbf{u}, \mathbf{w}_{t+1}) + d(\mathbf{w}_t, \mathbf{w}_{t+1})$$

เมื่อทำการหาผลรวมของสมการข้างต้น สำหรับทุกๆ  $t \in M$  จะได้ว่า

$$\lambda \hat{\gamma}_{\mathbf{u}, M} \text{Loss}^w(IWP', Q) \leq d(\lambda \mathbf{u}, \mathbf{w}_1) - d(\lambda \mathbf{u}, \mathbf{w}_{m+1}) + \sum_{t \in M} d(\mathbf{w}_t, \mathbf{w}_{t+1}) \quad (10)$$

เมื่อ  $\hat{\gamma}_{\mathbf{u}, M} = \frac{1}{\text{Loss}^w(IWP', Q)} \sum_{t \in M} c_t y_t \langle \mathbf{u}, \mathbf{x}_t \rangle$  และเมื่อใช้ทฤษฎีบทย่อยที่ 4 กับพจน์สุดท้ายของสมการ (10) จะได้ว่า

$$\sum_{t \in M} d(\mathbf{w}_t, \mathbf{w}_{t+1}) \leq \sum_{t \in M} \frac{c_t^2}{2} \|\mathbf{x}_t\|_2^2$$

อย่างไรก็ตาม เนื่องจาก  $c_t \leq 1$  ดังนั้น  $c_t^2 \leq c_t$ , เพราะฉะนั้นจะสามารถกำหนดขีดจำกัดของพจน์สุดท้ายของสมการ (10) ได้ดังนี้

$$\begin{aligned} \sum_{t \in M} \frac{c_t^2}{2} \|\mathbf{x}_t\|_2^2 &\leq \sum_{t \in M} \frac{c_t}{2} \|\mathbf{x}_t\|_2^2 \\ &= \frac{R^2}{2} \sum_{t \in M} c_t \\ &= \frac{R^2}{2} \text{Loss}^w(IWP', Q) \end{aligned} \quad (11)$$

จากสมการ (10) เนื่องจาก  $d(\lambda \mathbf{u}, \mathbf{w}_{t+1}) \geq 0$  จากนิยามการลู่เข้าของเบรกแมน ดังนั้นเมื่อแทน (8) และ (11) ในสมการ (10) จะได้ว่า

$$\lambda \hat{\gamma}_{\mathbf{u}, M} \text{Loss}^w(IWP', Q) \leq \frac{\lambda^2 \|\mathbf{u}\|_2^2 + \|\mathbf{w}_1\|_2^2}{2} - \lambda \langle \mathbf{u}, \mathbf{w}_1 \rangle + \frac{R^2}{2} \text{Loss}^w(IWP', Q) \quad (12)$$

พิจารณาสมการ  $\gamma - \frac{L_{u,\gamma}(M)}{\text{Loss}^w(IWP', Q)} \leq \hat{\gamma}_{u,M}$  ซึ่งเป็นจริง จากนิยามของความเบี่ยงเบนแบบมีน้ำหนักความสำคัญ  $D_\gamma(\mathbf{u}; (\mathbf{x}_i, y_i, c_i))$  และ  $\hat{\gamma}_{u,M}$  จะเห็นว่า หากค่าทางซ้ายของอสมการมีค่ามาก แสดงว่า ตัวทำนายเชิงเส้น  $h$  ทำนายได้ดีบน  $Q$  เมื่อพิจารณาเฉพาะครั้งที่  $IWP'$  ทำนายผิด หรือกล่าวได้ว่า เราเปรียบเทียบประสิทธิภาพการทำนายของ  $IWP'$  กับตัวทำนายที่ดี, เมื่อนำไปรวมกับอสมการ (12) จะได้

$$\frac{\lambda^2 \|\mathbf{u}\|_2^2 + \|\mathbf{w}_1\|_2^2}{2} - \lambda \langle \mathbf{u}, \mathbf{w}_1 \rangle + \frac{R^2}{2} \text{Loss}^w(IWP', Q) - \lambda \gamma \text{Loss}^w(IWP', Q) + \lambda L_{u,\gamma}(M) \geq 0 \quad (13)$$

พิจารณาค่าอสมการ (13) มีค่าน้อยที่สุดเมื่อ

$$\lambda = \frac{\langle \mathbf{u}, \mathbf{w}_1 \rangle - L_{u,\gamma}(M) + \gamma \text{Loss}^w(IWP', Q)}{\|\mathbf{u}\|_2^2}$$

ซึ่งมีค่าเป็นบวก ตามที่ข้อสมมติตอนต้นที่ว่า  $\text{Loss}^w(IWP', Q) > \frac{1}{\gamma} L_{u,\gamma}(M) - \frac{\langle \mathbf{u}, \mathbf{w}_1 \rangle}{\gamma}$  นั่นคือ

อสมการ (13) มีค่าน้อยที่สุดเป็น

$$\|\mathbf{w}_1\|_2^2 + R^2 \text{Loss}^w(IWP', Q) - \frac{(\langle \mathbf{u}, \mathbf{w}_1 \rangle - L_{u,\gamma}(M) + \gamma \text{Loss}^w(IWP', Q))^2}{\|\mathbf{u}\|_2^2} \quad (14)$$

เราจะแทน  $\text{Loss}^w(IWP', Q)$  ด้วย  $L_w$  ใน (14) และเรียกว่า (14') ซึ่งเป็นฟังก์ชันกำลังสองของ  $L_w$  ที่มีค่าสัมประสิทธิ์ของพจน์นำเป็นลบ และที่  $L_w = 0$  ค่าของ (14') มีค่าเป็นบวก เนื่องจาก  $\|\mathbf{u}\|_2 \|\mathbf{w}_1\|_2 \geq \langle \mathbf{u}, \mathbf{w}_1 \rangle$  ดังนั้นค่าของ  $\text{Loss}^w(IWP', Q)$  จะมีขีดจำกัดบนเป็นค่าขอบเขตบนที่น้อยที่สุดของ  $L_w > 0$  ที่ทำให้ (14') เป็นบวก นั่นคือค่ารากที่สองที่มากที่สุดของ (14') ซึ่งมีค่าเป็น

$$\begin{aligned} \text{Loss}^w(IWP', Q) &\leq \frac{1}{\gamma} L_{u,\gamma}(M) - \frac{\langle \mathbf{u}, \mathbf{w}_1 \rangle}{\gamma} + \frac{C}{2\gamma^2} \\ &\quad + \frac{1}{2\gamma^2} \sqrt{4\gamma C L_{u,\gamma}(M) + C^2 - 4\gamma C \langle \mathbf{u}, \mathbf{w}_1 \rangle + 4\gamma^2 \|\mathbf{u}\|_2^2 \|\mathbf{w}_1\|_2^2} \quad \blacksquare \end{aligned}$$

ทฤษฎีบทที่ 8 มีการกำหนดว่า  $\mathbf{w}_1 \in \mathfrak{R}^n$  และ  $c_i \leq 1$  ส่วนกรณีอัลกอริทึม *IWP* ที่มี  $\mathbf{w}_1=0$  และ  $c_i \geq 0$  สามารถเขียนอสมการข้างต้นใหม่ได้เป็น

$$\begin{aligned} \frac{Loss^w(IWP, Q)}{c_{\max}} &\leq \frac{L_{u,\gamma}(M)}{\gamma c_{\max}} + \frac{C}{2\gamma^2} + \frac{1}{2\gamma^2} \sqrt{4\gamma C \frac{L_{u,\gamma}(M)}{c_{\max}} + C^2} \\ &\leq \frac{L_{u,\gamma}(M)}{\gamma c_{\max}} + \frac{C}{2\gamma^2} + \frac{1}{2\gamma^2} \sqrt{\left( \sqrt{4\gamma C \frac{L_{u,\gamma}(M)}{c_{\max}}} + C \right)^2} \\ &= \frac{L_{u,\gamma}(M)}{\gamma c_{\max}} + \frac{C}{\gamma^2} + \frac{1}{\gamma^2} \sqrt{\gamma C \frac{L_{u,\gamma}(M)}{c_{\max}}} \end{aligned}$$

ทฤษฎีบทที่ 7 เป็นจริง เมื่อ  $\gamma = 1$

### ผลการวิจัยการแก้ปัญหาการเรียนรู้หลายประเภทแบบมีน้ำหนักความสำคัญ

สำหรับผลการวิจัยการแก้ปัญหาการเรียนรู้หลายประเภทแบบมีน้ำหนักความสำคัญแบ่งออกเป็น 2 ส่วน ได้แก่ ผลการวิจัยภาคทดลองและผลการวิจัยภาคทฤษฎี

#### ผลการวิจัยภาคทดลอง

งานวิจัยนี้ได้ทำการทดลองนี้ได้ทำการทดลองเพื่อเปรียบเทียบประสิทธิภาพการทำนายแบบหลายประเภทแบบมีน้ำหนักความสำคัญของ 3 อัลกอริทึม ได้แก่ อัลกอริทึมแบบลดรูปที่ใช้เพอร์เซพตรอนเป็นกลองคำ, อัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญ (*IWP*) และอัลกอริทึมเพอร์เซพตรอน เมื่อทำการลดรูปจากปัญหาการเรียนรู้แบบหลายประเภทไปเป็นปัญหาเรียนรู้สองประเภทด้วยวิธีการหลายเวกเตอร์และวิธีการทุกคู่ โดยทดลองกับชุดข้อมูลมาตรฐาน 4 ชุดข้อมูลใน Asuncion and Newman (2007) ดังนี้

#### ชุดข้อมูล Letter Recognition

ชุดข้อมูล Letter Recognition ประกอบด้วยข้อมูล 26 ประเภท และมีตัวอย่างทั้งสิ้น 20000 รายการ แต่ละรายการมี 16 คุณลักษณะ ซึ่งเป็นข้อมูลของภาพตัวอักษรพิมพ์ใหญ่ภาษาอังกฤษ แต่ละตัวอักษรจะแบ่งเป็น 20 รูปแบบตัวพิมพ์ (font) ในชุดข้อมูลนี้แต่ละตัวอย่างไม่มีน้ำหนัก

ความสำคัญกำกับ เราจึงทำการกำหนดค่าน้ำหนักความสำคัญของตัวอย่างใดๆ โดยใช้ค่าสุ่มจากศูนย์ถึงค่าความถี่ของตัวอักษรภาษาอังกฤษของตัวอย่างนั้นๆ เป็นน้ำหนักความสำคัญของตัวอย่างดังกล่าว

### ชุดข้อมูล Isolet

ชุดข้อมูล Isolet ประกอบด้วยข้อมูล 26 ประเภท และมีตัวอย่างทั้งสิ้น 6238 รายการ แต่ละรายการมี 617 คุณลักษณะ ซึ่งเป็นข้อมูลของเสียงของการอ่านตัวอักษรภาษาอังกฤษ โดยผู้อ่าน 50 คน ในชุดข้อมูลนี้แต่ละตัวอย่างไม่มีน้ำหนักความสำคัญกำกับ เราจึงทำการกำหนดค่าน้ำหนักความสำคัญของตัวอย่าง เช่นเดียวกับชุดข้อมูล Letter Recognition นั่นคือ ค่าน้ำหนักความสำคัญของตัวอย่างใดๆ มีค่าเป็นค่าสุ่มจากศูนย์ถึงค่าความถี่ของตัวอักษรภาษาอังกฤษของตัวอย่างนั้นๆ

### ชุดข้อมูล Car Evaluation

ชุดข้อมูล Car Evaluation ประกอบด้วยข้อมูล 4 ประเภท และมีตัวอย่างทั้งสิ้น 1728 รายการ แต่ละรายการมี 8 คุณลักษณะ ซึ่งเป็นข้อมูลการประเมินคุณภาพของรถยนต์ ในชุดข้อมูลนี้แต่ละตัวอย่างไม่มีน้ำหนักความสำคัญกำกับ เราจึงทำการกำหนดค่าน้ำหนักความสำคัญของตัวอย่างใดๆ เป็นค่าสุ่มจาก 0 ถึง  $\frac{1}{P(r)}$  โดยที่  $P(r)$  เป็นค่าความน่าจะเป็นที่จะสุ่มได้ตัวอย่างในประเภท

$$r \in Y$$

ในชุดข้อมูลนี้มีบางคุณลักษณะที่มีชนิดข้อมูลเป็นข้อความ เราได้ทำการแปลงคุณลักษณะเหล่านั้นเป็นกลุ่มของคุณลักษณะของ 0 หรือ 1 ในแนวทางเดียวกับ Rifkin and Klautau (2004)

### ชุดข้อมูล Abalone

ชุดข้อมูล Abalone ประกอบด้วยข้อมูล 28 ประเภท และมีตัวอย่างทั้งสิ้น 4177 รายการ แต่ละรายการมี 6 คุณลักษณะ ซึ่งเป็นข้อมูลที่ใช้ในการทำนายอายุของหอยทากจากข้อมูลทางกายภาพ ในชุดข้อมูลนี้แต่ละตัวอย่างไม่มีน้ำหนักความสำคัญกำกับ เราจึงทำการกำหนดค่าน้ำหนัก

ความสำคัญของตัวอย่างใดๆ เป็นค่าสุ่มจาก 0 ถึง  $\frac{1}{P(r)}$  โดยที่  $P(r)$  เป็นค่าความน่าจะเป็นที่จะสุ่ม  
ได้ตัวอย่างในประเภท  $r \in Y$

ในชุดข้อมูลนี้มีบางคุณลักษณะที่มีชนิดข้อมูลเป็นข้อความ เราได้ทำการแปลงคุณลักษณะ  
เหล่านั้นเป็นกลุ่มของคุณลักษณะของ 0 หรือ 1 ในแนวทางเดียวกับ Rifkin and Klautau (2004)

### ผลการทดลอง

เราได้ทำการเปรียบเทียบประสิทธิภาพการทำนายแต่ละอัลกอริทึมโดยใช้วิธีการเดียวกันใน  
ส่วนของการทดลองกับข้อมูลแบบสองประเภท นั่นคือทำการนับจำนวนครั้งและคิดค่าใช้จ่ายรวมทั้ง  
แต่ละวิธีการทำนายผิด สำหรับการทดลองกับชุดข้อมูลแบบหลายประเภทโดยการลดรูปด้วยวิธีการ  
หลายเวกเตอร์ได้ผลตามที่แสดงในตารางที่ 2 เมื่อพิจารณาร้อยละของจำนวนที่แต่ละวิธีการทำนาย  
ผิด (%ตอบผิด) จะเห็นว่าโดยส่วนใหญ่เพอร์เซพตรอนมีจำนวนครั้งในการทำนายผิดใกล้เคียงกับ  
วิธีการแบบลดรูปซึ่งน้อยกว่าวิธีการ *IWP* แต่เมื่อพิจารณาร้อยละของค่าใช้จ่ายรวมทั้งแต่ละวิธีการ  
ทำนายผิด (%ค่าใช้จ่าย) จะเห็นว่าอัลกอริทึม *IWP* กลับมีค่าใช้จ่ายรวมในการทำนายผิดน้อยกว่าอีก  
สองวิธีการที่ให้ผลการทดลองใกล้เคียงกัน ยกเว้นชุดข้อมูล Abalone ซึ่งทุกวิธีการผลการทดลอง  
ใกล้เคียงกัน นั่นคือทำนายผิดในเกือบทุกตัวอย่าง อย่างไรก็ตามในงานวิจัยของ Rifkin and Klautau  
(2004) ได้ชี้ให้เห็นว่าชุดข้อมูล Abalone เป็นชุดข้อมูลที่ยากต่อการทำนายให้ถูก เนื่องจาก  
อัลกอริทึมการเรียนรู้แบบออฟไลน์ที่สามารถทำนายได้ดิบบนตัวอย่างที่ไม่มีน้ำหนักความสำคัญยัง  
ทำนายผิดบนชุดข้อมูลดังกล่าวเกินร้อยละ 70

ส่วนการทดลองโดยการลดรูปด้วยวิธีการทุกคู่ได้ผลตามที่แสดงในตารางที่ 3 เมื่อ  
เปรียบเทียบประสิทธิภาพของแต่ละอัลกอริทึมที่ลดรูปด้วยวิธีการหลายเวกเตอร์ (ตารางที่ 2) และ  
วิธีการทุกคู่จาก (ตารางที่ 3) จะเห็นว่าการลดรูปด้วยวิธีการหลายเวกเตอร์และวิธีการทุกคู่มี  
ผลลัพธ์เท่ากัน ซึ่งเป็นประเด็นที่น่าสนใจว่าเหตุใดการลดรูป 2 วิธีการจึงให้ผลลัพธ์เท่ากัน โดยเรา  
จะทำการวิเคราะห์ประเด็นนี้ในส่วนของผลการวิจัยภาคทฤษฎี

**ตารางที่ 2** ผลการเปรียบเทียบประสิทธิภาพการทำนายของอัลกอริทึมแบบลดรูป, อัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญและอัลกอริทึมเพอร์เซพตรอน เมื่อทำการลดรูปจากปัญหาการทำนายหลายประเภทเป็นปัญหาการทำนายหลายประเภทด้วยวิธีการหลายเวกเตอร์ กับชุดข้อมูลแบบหลายประเภท

ชุดข้อมูล	IWP		ลดรูป		เพอร์เซพตรอน	
	%ตอบผิด	%ค่าใช้จ่าย	%ตอบผิด	%ค่าใช้จ่าย	%ตอบผิด	%ค่าใช้จ่าย
Letter	56.59	58.52	64.94	53.23	<b>55.74</b>	57.34
Isolet	<b>25.76</b>	26.56	43.52	<b>25.97</b>	25.92	26.67
Car	<b>20.27</b>	43.21	33.56	<b>39.84</b>	20.33	43.77
Abalone	85.10	93.60	92.99	94.41	<b>84.61</b>	<b>93.47</b>

**ตารางที่ 3** ผลการเปรียบเทียบประสิทธิภาพการทำนายของอัลกอริทึมแบบลดรูป, อัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญและอัลกอริทึมเพอร์เซพตรอน เมื่อทำการลดรูปจากปัญหาการทำนายหลายประเภทเป็นปัญหาการทำนายหลายประเภทด้วยวิธีการทุกคู่ กับชุดข้อมูลแบบหลายประเภท

ชุดข้อมูล	IWP		ลดรูป		เพอร์เซพตรอน	
	%ตอบผิด	%ค่าใช้จ่าย	%ตอบผิด	%ค่าใช้จ่าย	%ตอบผิด	%ค่าใช้จ่าย
Letter	56.59	58.52	64.94	53.23	<b>55.74</b>	57.34
Isolet	<b>25.76</b>	26.56	43.52	<b>25.97</b>	25.92	26.67
Car	<b>20.27</b>	43.21	33.56	<b>39.84</b>	20.33	43.77
Abalone	85.10	93.60	92.99	94.41	<b>84.61</b>	<b>93.47</b>

### ผลการวิจัยภาคทฤษฎี

ในส่วนของการวิจัยภาคทฤษฎีของการแก้ปัญหการเรียนรู้หลายประเภทแบบมีน้ำหนักความสำคัญ เราได้ทำการวิเคราะห์ประเด็นที่กล่าวถึงไว้ในส่วนของผลการวิจัยภาคทฤษฎีที่ว่า เหตุใดการลดรูปด้วยวิธีการหลายเวกเตอร์จึงให้ผลลัพธ์เท่ากับวิธีการทุกคู่ ซึ่งเราจะทำการพิสูจน์ว่าการ

ลดรูป 2 วิธีการดังกล่าวนี้สมมูลกัน และจากการสมมูลกันนี้ทำให้เราสามารถวิเคราะห์ค่าขีดจำกัดบนของความสูญเสียแบบมีน้ำหนักของวิธีการหลายเวกเตอร์ที่ใช้ *IWP* เป็นฐาน และวิธีการทุกคู่ที่ใช้ *IWP* เป็นฐานที่ได้นำเสนอไปในส่วนของวิธีการได้ในคราวเดียว

### การสมมูลกันระหว่างวิธีการหลายเวกเตอร์และวิธีการทุกคู่

ในส่วนนี้จะทำการพิสูจน์ว่าการลดรูปจากปัญหาการเรียนรู้อะไรหลายประเภทไปเป็นปัญหาการเรียนรู้อะไรสองประเภทด้วยวิธีการหลายเวกเตอร์และวิธีการทุกคู่ที่สมมูลกัน ในการพิสูจน์เพื่อให้ง่ายขึ้นเราจะพิสูจน์กรณีที่ตัวอย่างไม่มีน้ำหนักความสำคัญ นั่นคือ  $c_t = 1$  สำหรับทุกๆ  $t$  ส่วนในกรณีที่ตัวอย่างมีน้ำหนักความสำคัญ สามารถพิสูจน์ได้ด้วยวิธีการใกล้เคียงกัน

สำหรับปัญหาการเรียนรู้อะไร  $k$  ประเภท ให้  $Y = \{1, 2, \dots, k\}$  สำหรับการทำนายรอบที่  $t$  ใดๆ ให้  $h_t^M(\mathbf{x}, r)$  แทนคะแนนของประเภท  $r$  สำหรับวิธีการหลายเวกเตอร์ และ  $h_t^A(\mathbf{x}, r)$  แทนคะแนนของประเภท  $r$  สำหรับวิธีการทุกคู่ สังเกตว่า สำหรับ  $\mathbf{x} \in \mathcal{X}^n$  ใดๆ และสำหรับประเภท  $r$  ใดๆ ถ้า  $h_t^A(\mathbf{x}, r) = k \cdot h_t^M(\mathbf{x}, r)$  แล้วทั้งสองวิธีการจะทำนายเหมือนกันในทุกๆ ตัวอย่าง  $\mathbf{x}$  เนื่องจาก ถ้า  $h_t^M(\mathbf{x}, i) > h_t^M(\mathbf{x}, j)$  แล้ว  $h_t^A(\mathbf{x}, i) > h_t^A(\mathbf{x}, j)$  ต่อไปจะเป็นการพิสูจน์ทฤษฎีบทย่อยที่สำคัญสำหรับการพิสูจน์ความสมมูลของทั้งสองวิธีการ

**ทฤษฎีบทย่อยที่ 5** ถ้าการทำนายรอบที่  $t$ ,  $h_t^A(\mathbf{x}, r) = k \cdot h_t^M(\mathbf{x}, r)$  สำหรับ  $\mathbf{x} \in \mathcal{X}^n$  ใดๆ และสำหรับประเภท  $r$  ใดๆ แล้ว

$$h_{t+1}^A(\mathbf{x}, r) = k \cdot h_{t+1}^M(\mathbf{x}, r)$$

สำหรับ  $\mathbf{x} \in \mathcal{X}^n$  ใดๆ และสำหรับ  $r$  ใดๆ

**บทพิสูจน์** กรณีที่การทำนายรอบที่  $t$  ถูกต้องเวกเตอร์น้ำหนักจะไม่ถูกปรับปรุง ดังนั้นทฤษฎีบทตั้งเป็นจริง ส่วนในกรณีที่ทำนายผิดในรอบที่  $t$  สำหรับตัวอย่าง  $\mathbf{x}_t$  เนื่องจากทั้งสองวิธีการมีการทำนายเหมือนกัน นั่นคือทำนายประเภทที่ให้ค่าคะแนนมากที่สุด ดังนั้นทั้งสองวิธีการจะต้องปรับปรุงเวกเตอร์น้ำหนักดังนี้

ให้  $Y$  เป็นเซตของประเภททั้งหมดที่เป็นไปได้ และสำหรับการทำนายรอบที่  $t$  ให้  $\hat{y}_t \in Y$  เป็นประเภทที่ถูกทำนาย และ  $y_t \in Y$  เป็นประเภทที่ถูกต้อง พิจารณาการปรับปรุงเวกเตอร์น้ำหนักของวิธีการหลายเวกเตอร์ซึ่งจะปรับปรุงเวกเตอร์น้ำหนักดังนี้

$$\mathbf{w}_{t+1}^{y_t} \leftarrow \mathbf{w}_t^{y_t} + \mathbf{x}_t, \quad \mathbf{w}_{t+1}^{\hat{y}_t} \leftarrow \mathbf{w}_t^{\hat{y}_t} - \mathbf{x}_t$$

ส่วน  $\mathbf{w}_{t+1}^r$  สำหรับ  $r \in Y \setminus \{y_t, \hat{y}_t\}$  จะไม่ถูกปรับปรุง ดังนั้นสำหรับตัวอย่าง  $\mathbf{x}_t$  ใดๆ

$$\begin{aligned} h_{t+1}^M(\mathbf{x}, y_t) &= \langle \mathbf{w}_{t+1}^{y_t} + \mathbf{x}_t, \mathbf{x} \rangle \\ &= \langle \mathbf{w}_t^{y_t}, \mathbf{x} \rangle + \langle \mathbf{x}_t, \mathbf{x} \rangle \\ &= h_t^M(\mathbf{x}, y_t) + \langle \mathbf{x}_t, \mathbf{x} \rangle \end{aligned} \quad (15)$$

ในทำนองเดียวกัน

$$h_{t+1}^M(\mathbf{x}, \hat{y}_t) = h_t^M(\mathbf{x}, \hat{y}_t) - \langle \mathbf{x}_t, \mathbf{x} \rangle \quad (16)$$

ในส่วนของการปรับปรุงเวกเตอร์น้ำหนักของวิธีการทุกคู่ พิจารณากรณี  $r \in Y \setminus \{y_t, \hat{y}_t\}$

$$\begin{aligned} h_{t+1}^A(\mathbf{x}, y_t) &= \sum_{j \in Y \setminus \{r\}} \langle \mathbf{w}_{t+1}^j, \mathbf{x} \rangle \\ &= \left\langle \sum_{j \in Y \setminus \{r\}} \mathbf{w}_{t+1}^j, \mathbf{x} \right\rangle \end{aligned}$$

เวกเตอร์น้ำหนักที่ถูกปรับปรุงจะมีเพียง  $\mathbf{w}_{t+1}^{y_t}$  และ  $\mathbf{w}_{t+1}^{\hat{y}_t}$  นั่นคือ

$$\mathbf{w}_{t+1}^{y_t} \leftarrow \mathbf{w}_t^{y_t} + \mathbf{x}_t, \quad \mathbf{w}_{t+1}^{\hat{y}_t} \leftarrow \mathbf{w}_t^{\hat{y}_t} - \mathbf{x}_t$$

ซึ่งจะทำให้ผลรวมของเวกเตอร์น้ำหนักไม่เปลี่ยนแปลง นั่นคือ

$$\sum_{j \in Y \setminus \{r\}} \mathbf{w}_{t+1}^j = \sum_{j \in Y \setminus \{r\}} \mathbf{w}_t^j$$

ดังนั้นจะได้ว่า สำหรับ  $\mathbf{x}$  ใดๆ และ สำหรับ  $r \in Y \setminus \{y_t, \hat{y}_t\}$  ใดๆ

$$h_{t+1}^A(\mathbf{x}, r) = h_t^A(\mathbf{x}, r) \quad (17)$$

พิจารณาคณิประเภท  $y_t$  ซึ่งมี  $h_{t+1}^A(\mathbf{x}, y_t) = \sum_{j \in Y \setminus \{y_t\}} \langle \mathbf{w}_{t+1}^{y_t, j}, \mathbf{x} \rangle$  ดังนั้น สำหรับ  $r \in Y \setminus \{y_t, \hat{y}_t\}$  จะทำการปรับปรุงเวกเตอร์น้ำหนักดังนี้

$$\mathbf{w}_{t+1}^{y_t, r} \leftarrow \mathbf{w}_t^{y_t, r} + \mathbf{x}_t$$

และสำหรับประเภท  $\hat{y}_t$  เวกเตอร์น้ำหนัก  $\mathbf{w}_t^{y_t, \hat{y}_t}$  จะถูกปรับปรุงสองครั้ง นั่นคือ

$$\mathbf{w}_{t+1}^{y_t, \hat{y}_t} \leftarrow \mathbf{w}_t^{y_t, \hat{y}_t} + 2 \cdot \mathbf{x}_t$$

ผลรวมของเวกเตอร์น้ำหนักสำหรับประเภท  $y_t$  ของการทำนายรอบที่  $t+1$  มีค่าเป็น

$$\sum_{j \in Y \setminus \{y_t\}} \mathbf{w}_{t+1}^{y_t, j} = \sum_{j \in Y \setminus \{y_t\}} (\mathbf{w}_t^{y_t, j} + k \cdot \mathbf{x}_t)$$

ดังนั้นจะได้ว่า

$$\begin{aligned} h_{t+1}^A(\mathbf{x}, y_t) &= \sum_{j \in Y \setminus \{y_t\}} \langle \mathbf{w}_{t+1}^{y_t, j}, \mathbf{x} \rangle \\ &= \left\langle \sum_{j \in Y \setminus \{y_t\}} \mathbf{w}_t^{y_t, j} + k \cdot \mathbf{x}_t, \mathbf{x} \right\rangle \\ &= \left\langle \sum_{j \in Y \setminus \{y_t\}} \mathbf{w}_t^{y_t, j}, \mathbf{x} \right\rangle + k \cdot \langle \mathbf{x}_t, \mathbf{x} \rangle \\ &= h_t^A(\mathbf{x}, y_t) + k \cdot \langle \mathbf{x}_t, \mathbf{x} \rangle \end{aligned} \quad (18)$$

ในทำนองเดียวกัน

$$h_{t+1}^A(\mathbf{x}, \hat{y}_t) = h_t^A(\mathbf{x}, \hat{y}_t) - k \cdot \langle \mathbf{x}_t, \mathbf{x} \rangle \quad (19)$$

ทฤษฎีบทตั้งเป็นจริงตามสมการ (15) - (19) ■

จากทฤษฎีบทย่อยที่ 5 สังเกตว่า ถ้า  $h_t^A(\mathbf{x}, r) = k \cdot h_t^M(\mathbf{x}, r)$  แล้ววิธีการหลายเวกเตอร์และวิธีการทุกคู่จะทำนายเหมือนกันในทุกๆ รอบ ดังนั้นจะได้ว่า ทั้งสองวิธีการสมมูลกัน อย่างไรก็ตามเมื่อพิจารณาแง่ของเวลาในการทำงาน จะเห็นว่าวิธีการหลายเวกเตอร์ทำงานที่ใช้เวลาในการทำงานเป็น  $O(k)$  นั้นทำงานได้เร็วกว่าวิธีการทุกคู่ที่ใช้เวลาในการทำงานเป็น  $O(k^2)$

**การรับประกันค่าขีดจำกัดบนของความสูญเสียแบบมีน้ำหนักของวิธีการหลายเวกเตอร์และวิธีการทุกคู่ที่ใช้เพอร์เซพตรอนแบบมีน้ำหนักความสำคัญเป็นฐาน**

สำหรับการพิสูจน์ค่าขีดจำกัดบนของความสูญเสียแบบมีน้ำหนักของวิธีการหลายเวกเตอร์ที่ใช้ *IWP* เป็นฐานนั้น สามารถใช้แนวทางการพิสูจน์เดียวกับการพิสูจน์ทฤษฎีบทที่ว่าด้วยค่าขีดจำกัดบนของความสูญเสียแบบมีน้ำหนักของวิธีการหลายเวกเตอร์ที่ใช้เพอร์เซพตรอนเป็นฐานในงานวิจัยของ Fink *et al.* (2006) ซึ่งมีการนิยามฟังก์ชันความสูญเสียแบบอิงกันดังนี้

ให้  $Y = \{1, \dots, k\}$  เป็นเซตของประเภทของตัวอย่างและสมมติฐาน  $h_t$  ใดๆ ซึ่งมีนิยามตามสมการ (3) สำหรับตัวอย่าง  $(\mathbf{x}_t, y_t, c_t)$  ใดๆ ฟังก์ชันความสูญเสียแบบอิงกันสำหรับสมมติฐาน  $h_t$  มีค่าเป็น

$$\ell(h_t, \mathbf{x}_t) = \left( 1 - h_t(\mathbf{x}_t, y_t) + \max_{r \in Y \setminus \{y_t\}} h_t(\mathbf{x}_t, r) \right)_+ \quad (20)$$

เมื่อ  $y_t \in Y$  และ  $(\cdot)_+ = \max\{0, \cdot\}$  จากฟังก์ชันความสูญเสียแบบอิงกันในสมการ (20) และนิยามฟังก์ชันความสูญเสียแบบมีน้ำหนักในสมการ (4) เราสามารถนิยามฟังก์ชันความสูญเสียแบบอิงกันแบบมีน้ำหนัก สำหรับการทำนายรอบที่  $t$  ได้เป็น

$$\ell^w(h_t, \mathbf{x}_t) = c_t \left( 1 - h_t(\mathbf{x}_t, y_t) + \max_{r \in Y \setminus \{y_t\}} h_t(\mathbf{x}_t, r) \right)_+ \quad (21)$$

เมื่อ  $y_t \in Y$  และ  $(\cdot)_+ = \max\{0, \cdot\}$

เราสามารถรับประกันค่าขีดจำกัดบนของความสูญเสียแบบมีน้ำหนักของวิธีการหลาย  
 เวกเตอร์ที่ใช้ *MWP* เป็นฐานได้ โดยใช้ทฤษฎีบทที่ 7 เป็นบทตั้งต้นและใช้ความสูญเสียแบบมี  
 น้ำหนักตามที่นิยามในสมการ (21) ได้ดังบทสืบเนื่องต่อไปนี้

**บทสืบเนื่องที่ 1** สำหรับลำดับข้อมูลป้อนเข้าแบบมีน้ำหนัก  $Q = ((x_1, y_1, c_1), (x_2, y_2, c_2), \dots, (x_m, y_m, c_m)) \in (\mathcal{X}^n \times Y \times [0, \infty))^m$  ใดๆ ถ้ามีตัวทำนายเชิงเส้น  $h \in H$  ใดๆ ที่นิยามด้วยเวกเตอร์  $U = \{u^r : r \in Y\}$  โดยที่  $\ell^w(h, Q) \leq L$  แล้ว วิธีการหลายเวกเตอร์ที่ใช้เพอร์เซพตรอนแบบมีน้ำหนัก  
 ความสำคัญเป็นฐาน จะมีค่าความสูญเสียแบบมีน้ำหนักไม่เกิน

$$L + c_{\max} C + \sqrt{c_{\max} LC}$$

$$\text{เมื่อ } R = 2 \cdot \max_{t \in \{1, \dots, m\}} \|x_t\|_2, \quad C = R^2 \sum_{r \in Y} \|u^r\|_2^2 \quad \text{และ} \quad c_{\max} = \max_{t \in \{1, \dots, m\}} c_t$$

เนื่องจากวิธีการหลายเวกเตอร์และวิธีการทุกคู่สมมูลกัน ดังนั้นเราสามารถใส่บทสืบเนื่อง  
 ข้างต้นในการรับประกันค่าขีดจำกัดบนของความสูญเสียแบบมีน้ำหนักของวิธีการทุกคู่ที่ใช้เพอร์  
 เซพตรอนแบบมีน้ำหนักความสำคัญเป็นฐานได้ด้วย

## สรุป

งานวิจัยนี้เสนอวิธีการแก้ปัญหาการเรียนรู้แบบมีน้ำหนักความสำคัญทั้งแบบสองประเภท และแบบหลายประเภท สำหรับปัญหาการเรียนรู้สองประเภทแบบมีน้ำหนักความสำคัญนำเสนอ 2 อัลกอริทึม ได้แก่ อัลกอริทึมแบบลดรูป และ อัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญ ซึ่ง อัลกอริทึมที่สองให้ผลทั้งเชิงทฤษฎีและเชิงทดลองดีกว่าอัลกอริทึมแรก ส่วนปัญหาการเรียนรู้หลายประเภทแบบมีน้ำหนักความสำคัญเสนอให้ใช้อัลกอริทึมเพอร์เซพตรอนแบบมีน้ำหนักความสำคัญเป็นฐานและใช้การลดรูปจากปัญหาการเรียนรู้แบบหลายประเภทไปเป็นปัญหาการเรียนรู้แบบสองประเภทด้วย 2 วิธีการ ได้แก่ วิธีการหลายเวกเตอร์ และ วิธีการทุกคู่ และได้ทำการพิสูจน์ว่าการลดรูปสองวิธีดังกล่าวสมมูลกัน

อย่างไรก็ตาม งานวิจัยนี้ยังไม่ได้พิจารณากรณีการทำนายหลายประเภท ที่ความสูญเสียที่เกิดขึ้นจากการทำนายประเภทของข้อมูลผิดพลาด อาจไม่เท่ากันในทุกประเภท การเรียนรู้ในรูปแบบดังกล่าว เรียกว่าการเรียนรู้ที่มีความไวต่อค่าใช้จ่าย ผู้วิจัยหวังว่าแนวทางสร้างอัลกอริทึม และการพิสูจน์ที่นำเสนออาจนำไปใช้ในกรณีดังกล่าวได้ด้วย

## เอกสารและสิ่งอ้างอิง

- Asuncion, A. and D.J. Newman. 2007. **UCI Machine Learning Repository**. UCI Machine Learning Repository: Data Sets. Available Source: <http://www.ics.uci.edu/~mlern/MLRepository.html>, February 11, 2008.
- Abe, N., Zadrozny, B. and J. Langford. 2004. An Iterative Method for Multi-Class Cost-Sensitive Learning, pp. 3-11. *In Proceedings of the 10<sup>th</sup> ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, New York, NY, USA.
- Cammer, K., Dekel, O., Keshet, J., Shalev-Shwartz, S. and Y. Singer. 2006. Online Passive-Aggressive Algorithms. **The Journal of Machine Learning Research** (7): 551-585.
- Cesa-Bianchi, N. and G. Lugosi. 2006. **Prediction, Learning, and Games**. Cambridge University Press, New York, NY, USA.
- Domingos, P. 1999. MetaCost: A General Method for Making Classifiers Cost-Sensitive, pp. 155-164. *In Proceedings of the 5<sup>th</sup> ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, New York, NY, USA.
- Elkan, C. 2001. The Foundations of Cost-Sensitive Learning, vol. 2, pp. 973-973. *In Proceedings of the 7<sup>th</sup> International Joint Conference on Artificial Intelligence*, Seattle.
- Fink, M., Shalev-Shwartz, S., Singer, Y. and S. Ullman. 2006. Online Multiclass Learning by Interclass Hypothesis Sharing, pp. 313-320. *In Proceedings of the 23<sup>rd</sup> International Conference on Machine Learning*. ACM, New York, NY, USA.
- Gentile, C. 2003. The Robustness of the  $p$ -Norm Algorithms. **Machine Learning** 3 (10): 265-299.

- Novikoff, A.B. 1962. On Convergence Proofs on Perceptrons, pp. 615-622. *In Proceedings of the Symposium on the Mathematical Theory of Automata*. Polytechnic Institute of Brooklyn.
- Rifkin, R. and A. Klautau. 2004. In Defense of One-Vs-All Classification. **The Journal of Machine Learning Research** (5): 101–141.
- Rosenblatt, F. 1958. The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain. **Psychological Review** 6 (65): 386-408.
- X. Liu. 2006. Training Cost-Sensitive Neural Networks with Methods Addressing the Class Imbalance Problem. **IEEE Transactions on Knowledge and Data Engineering** 1 (18): pp. 63-77.
- Zadrozny, B. and C. Elkan. 2001. Learning and Making Decisions when Costs and Probabilities are both Unknown, pp. 204-213. *In Proceedings of the 7<sup>th</sup> ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, San Francisco, California.
- Zadrozny, B., Langford, J. and N. Abe. 2003. Cost-Sensitive Learning by Cost-Proportionate Example Weighting, pp. 435-442. *In Proceedings of the 3<sup>rd</sup> IEEE International Conference on Data Mining*. IEEE Computer Society, Washington, DC, USA.

## ประวัติการศึกษา และการทำงาน

ชื่อ –นามสกุล	อดิศักดิ์ สุทธิสุน
วัน เดือน ปี ที่เกิด	1 มีนาคม พ.ศ. 2524
สถานที่เกิด	จังหวัดสระบุรี
ประวัติการศึกษา	วศ.บ. (วิศวกรรมคอมพิวเตอร์) มหาวิทยาลัยเกษตรศาสตร์ วิทยาเขตศรีราชา (พ.ศ. 2546)
ตำแหน่งหน้าที่การงานปัจจุบัน	-
สถานที่ทำงานปัจจุบัน	-
ผลงานดีเด่นและรางวัลทางวิชาการ	-
ทุนการศึกษาที่ได้รับ	-