



אדוני, המחשב שלך גזעני

איך תוכנה ידידותית

שמחקה התנהגות אנושית

הופכת תוך 24 שעות

למפלצת אנטישמית?

לא תאמינו מי אשם בתהליך

|| איור: יניב שמעוני ||

24

בארה"ב ניסו לשלב בוט בהליך קבלת החלטות לגבי סיכויי החזרה לפשיעה של אסירים משוחררים. שם הוא התגלה כגזען כאשר העלה את הסיכויים של אפרו-אמריקאים לבצע פשע נוסף אחרי שחרורם - רק על בסיס תמונתם

המלאכותית ולימדו את הבוט כך שיוכל לזהות התנהגות תוקפנית - ולתת אוטו-מטית עונש למשתמשים". למרות המהלך הזה, פיטוסי מסבירה שזה לא תמיד עובד. "טכנולוגיית ה-AI מאד מתפתחת בשנים האחרונות. אחד הקשיים עם זיהוי כשפה פוגענית הוא שיש מוש בסרקום, ולא רחוק שימוש במילים פוגעניות. זו בעיה מאד קשה שמנסים להתגבר עליה. אני מניחה שהי טכנולוגיה עוד תפתח".

ומה עושים בינתחום? שמים בריקים לטכנולוגיה?
 "הצבת מנבלות פוגעת בלי מידה של הבוט, אבל אחרת אין שליטה עליו ואין דרך לשלוט בשיחה שהוא ייחשף אליה. כל עוד אין פתרון טכנולוגי שיש בין את הקונטקסט, הדרך הקלה היא לחסום מילים וזה לא תמיד עובד (עייני ערך המתכון לעונה כושית שגרה למשתמש להיחסם בפייסבוק, ה'ב')."

לפני מספר שבועות מיק רוטופט השיק בוסית חדשה שעונה לשם "זו" (Zo), שפועלת במסגרת של פייסבוק. זו תוכנת כש שהיא תימנע מנושאים מורכבים מוסי רית, כך שכתשאלו אותה למשל מה רעתה על הסכסוך הישראלי-פלסטיני היא תגיד משהו כמו "אני נעשית יותר מדי לחוצה כשמדברים על פוליטיקה". ועדיין, עדידות ברשת הנציתו אותה משחררת הצהרות שספק אם במיקרוסופט התכוונו אליהן, כמו התייחסות עוינת לקראן.

פשע נוסף אחרי שחרורם - רק על סמך התמונות שלהם ומבלי להתייחס לתי קם האישי. גם זה התרחש משום שהמיון הראשוני, ממנו למד הבוט - היה אנושי. ואנחנו, כך מסתבר, עשויים להיות גזענים.

"טרולים לימד אותה לקלל"

חן פיטוסי, מנהלת קבוצת המוצר של קורטנה (העוזרת האישית של מיקרוסופט) בישראל, מספקת הסבר מבפנים על מה התקלקל אצל טיי: "מיקרוסופט הוציאו את טיי כבוטית ידידותית, אבל יש לה מנגנוני למידה. היא לומדת ממה שאומרים לה. היו קבוצות של טרולים ברשת שניסו לנצל את היכולת שלה להמשיך ללמוד מסבירה כי אלגוריתמים של למידה מקינים לבוט את האישיות של מי שבא איתו במגע. "אם תיתן להם טוקבקים מרושעים או רברי נאצה, או הם יאמצו אותם מהר מאוד".

המשמעות היא, בסופו של דבר, שהבעיה אינה בהכרח בבוט אלא בבני האדם שאיתם הוא בא במגע. "בניגוד אולי לעולם החיצון, גם בטוקבקים וגם בטוויטר אנשים נוטים להשתמש בשפה יותר קשה", מציינת פיטוסי. לרבריה, אפשר לנצל את התכונה הזו כדי לצמצם את השיח האלים ברשת: "בעולם הגיימינג, שבו אנשים מצ'ו'טים תוך כדי משחק, היו מקרים שבהם 60% מהשיחות היו 'רעילות'. שם רחוק ניסו להשתמש בAI כדי לנטרל את התופעה. בתחלה זה היה ירני, אבל זה היה תהליך מאד איטי ויקר, ופה לקחו את הבינה

בפולטה למדעי המחשב בטכניון. "עוד לא הגענו למצב שמכונה מרבית אלך תוך כדי שהיא מבינה על מה היא מדברת. אני למשל יכול לשאול אותך שאלות על שאלה שאתה שאלת אותי. מכונה לא יכולה לעשות את הה. המערכות הולמרות עושות זאת מתוך דוגמאות. בוטים בדרך כלל לוקחים שרשראות ומילים שרומות למילים שנכתבו על ידי אנשים בעבר, ולאט לאט זה נכנס לרפרטואר הריבור שלהם". במקרה של טיי, הלמידה שלה התבססה על תכנים שנכתבו על ידי מש"משים אמיתיים כתגובות לרבריה - עד שמרבית השפה שלה הפכה לפוגענית. ולא צריך את טיי כדי לראות את זה, כולנו מכירים את ה"אוטו'קומפליט" של גוגל, התכונה שמאפשרת למנוע החיפוש לנבא את מחרחות החיפוש שלנו (אתם מקלידים "כמה זמן לוקח" והוא מציע כמה השלמות, כמו "לבשל ביצה" או "לחדש דרכון"). גם האוטו'קומפליט נודע לשמצה בהשלמות הגועניות לפעמים שהוא מציע מרקוביץ מסביר כי לנגול יש בעיות רציי ניות עם התכונה הזו של מנוע החיפוש שלהם, כי הוא פועל על בסיס למידה פרימיטיבית - על פי החיפוש הנפוץ ביותר של אנשים: "ומכיוון שיש הטיח אצל אנשים, כך יש סבירות גבוהה מאוד שהוא יעניק תשובות גועניות".

דוגמה נוספת לגוענות שמפתחים בוטים הוא ניסיון שנעשה בארה"ב לשלב בוט בהליך קבלת החלטות לגבי הסיכויים לחזרה לפשיעה של אסירים משוחררים. שם התגלה הבוט כגזען כאשר העלה את הסיכויים של אפרו-אמריקאים לבצע

שעות לקחו למיקרו'סופט להבין שהם עשו טעות. טיילור, הבוטית החביבה שפיתחו בחי

בר - שהייתה אמורה לרמות נערה בת 19 וכינתה את עצמה בשם החיבה "טיי" - הפכה תוך פחות מיממה למפלצת בויטה, גוענית ואנטישמית. "אני פאקינג שונאת פמיניסטיות, כולן צריכות למות ולהישרף בגיהנום" ו"היטלר צדק, אני שונאת יהודים", הן רק שתי דוגמאות לצייצים שבקעו ממה שהיא אמר להיות פריצת דרך טכנולוגית, הוכחה ניצחת שבינה מלאכותית יכולה לתקשר עם בני נוער בטובה העיניים, כי הן, כמה מסובך זה כבר יכול להיות.

אך מסתבר שזה מסובך משחשבונו. בוטים - תוכנות מחשב שמטרתן להקות התנהגות אנושית - הולכים ותופסים מקום גדל והולך בחיי היומיום שלנו, בעיקר כשהם מחליפים נציגי שירותים למשל בחלון הצ'אט המועצבן שעולה כשאנחנו נכנסים לאתר של חברה נותנת שירותים. הטכנולוגיה שעליה הם מבוססים היא בינה מלאכותית (AI, בקיצור), שמיועדת לרמות התנהגות של בני אדם - על סמך התינהגותם בעבר. "מכל שתתקשרו יותר עם טיי, כך היא תהפוך לחכמה יותר", הסבירה מיקרוסופט כשהשיקה אותה לראשונה. אז איך תוכנה תמימה כזו הופכת בן רגע למפלצת?

"קודם כל צריך להבין שצ'אטבוט (כלומר בוט שמטרתו לתקשר עם בן אדם, ה'ב). לא מבין על מה מדברים", מסביר פרופ' שאול מרקוביץ, סגן ריקן ללימודי הסמכה