

Abordări recente în evaluarea formativă a utilizabilității

Costin Pribeanu, Ruxandra Dora Marinescu, Irina Cristescu, Lucreția Pîrvu, Dragoș Daniel Iordache

Institutul Național de Cercetare-Dezvoltare în Informatică – ICI București

Bd. Mareșal Averescu nr.8-10, București

{pribeanu, doruma, irina.cristescu, lpîrvu, iordache}@ici.ro

REZUMAT

Creșterea gradului de utilizare cât și creșterea calitativă a informației digitale necesită o abordare centrată pe utilizator a procesului de proiectare, care să conducă la o îmbunătățire substanțială a utilizabilității aplicațiilor informatice. Scopul acestui articol este de a prezenta sintetic cele mai recente abordări în evaluarea formativă a utilizabilității sistemelor interactive. Problematika prezentată este diversă, incluzând indicatori folosiți în prezent pentru gradul de încredere și validitatea rezultatelor obținute cu diferite metode.

Cuvinte cheie

Utilizabilitate, evaluare formativă, evaluare euristică.

Clasificare ACM

D.2.2: Design tools and techniques. H5.2 User interfaces.

INTRODUCERE

Dezvoltarea societății informaționale moderne se bazează pe accesul larg la tehnologia informatică într-o economie din ce în ce mai competitivă. În acest sens, creșterea gradului de utilizare a tehnologiei informatice și creșterea cantitativă și calitativă a informației digitale sunt două priorități majore. Atât creșterea gradului de utilizare cât și creșterea calitativă a informației digitale necesită o abordare centrată pe utilizator a procesului de proiectare, care să conducă la o îmbunătățire substanțială a utilizabilității aplicațiilor informatice.

În funcție de momentul și scopul evaluării, evaluarea utilizabilității poate fi formativă sau sumativă. Evaluarea formativă a utilizabilității este efectuată iterativ, pe parcursul ciclului de dezvoltare cu scopul identificării și eliminării problemelor de utilizabilitate cât mai devreme (Theofanos & Quesenberry, 2005).

Deși pe plan internațional s-au făcut eforturi considerabile de cercetare în ultima decadă, probleme serioase de utilizabilitate apar în continuare cu o frecvență ridicată la toate tipurile de produse software. Acest fapt denotă că evaluarea utilizabilității nu este încă integrată în ciclul de dezvoltare a sistemelor software.

Pe plan național, evaluarea utilizabilității este practică sporadic, în cadrul unor proiecte de cercetare. Utilizabilitatea produselor informatice este scăzută, întrucât nu există încă o practică curentă de evaluare a produselor informatice înainte de a fi instalate la utilizator.

Scopul acestui articol este de a prezenta sintetic cele mai recente abordări în evaluarea formativă a utilizabilității sistemelor interactive. Problematika prezentată este diversă, incluzând indicatori folosiți în prezent pentru

gradul de încredere și validitatea rezultatelor obținute cu diferite metode.

SCOP ȘI METODE UTILIZATE

Scopul evaluării formative

Evaluarea formativă a utilizabilității are ca scop identificarea problemelor de utilizabilitate, analiza acestora și formularea de sugestii pentru remediere. O preocupare de interes actual este analiza comparativă a problemelor de utilizabilitate care pot fi identificate cu diferite metode.

O problemă de utilizabilitate a fost definită de Nielsen (1993) ca orice aspect al interfeței cu utilizatorul despre care se presupune (se observă) că ar crea utilizatorului dificultăți / nemulțumiri raportat la un indicator important al utilizabilității (ușurință în învățare, ușurință în operare, rată de erori, satisfacție subiectivă) și care poate fi atribuit unui singur aspect de proiectare.

Problemele de utilizabilitate se pot clasifica după gradul de severitate în trei categorii:

- Severe: probleme care nu permit ca utilizatorul să îndeplinească sarcina sau care rezultă în pierderi importante de date sau timp.
- Moderate: probleme care au un impact semnificativ asupra îndeplinirii sarcinii dar pentru care utilizatorul poate găsi o soluție.
- Mînore: probleme care sunt iritante pentru utilizator dar care nu au un impact semnificativ asupra îndeplinirii sarcinii.

Pierderea de date este considerată catastrofală fie dacă datele nu pot fi refăcute, fie dacă refacerea necesită costuri foarte mari. Pierderea de timp este relativă la durata sarcinii.

Practica standard cere remedierea problemelor severe de utilizabilitate înainte de furnizarea produsului / aplicației software. În cel mai rău caz, acestea nu trebuie tolerate în sarcinile critice sau în cele cu o frecvență ridicată de execuție.

Problemele de utilizabilitate pot fi identificate atât prin testare cu utilizatori cât și prin evaluare euristică. Așa cum arată Law & Hvannberg (2002), cele două metode sunt complementare.

Evaluarea euristică

Evaluarea euristică – HE (Heuristic Evaluation) este o metodă de inginerie a utilizabilității realizată de un număr redus de evaluatori care examinează o interfață cu utilizatorul, judecând respectarea unui set de principii de

utilizabilitate (euristici), elaborează o listă de probleme de utilizabilitate (UP) clasificate pe categorii de severitate, corespunzător impactului estimat asupra performanțelor utilizatorului sau acceptanței (Molich & Nielsen, 1990; Nielsen, 1994).

Evaluarea euristică este, în esență, o inspecție de utilizabilitate. Se recomandă utilizarea acestei metode înaintea testării cu utilizatori, astfel încât să se elimine un număr cât mai mare de probleme de utilizabilitate cu costuri cât mai reduse.

De regulă, în evaluarea euristică se colectează:

- Măsuri cantitative: numărul de probleme de utilizabilitate pe două niveluri de severitate: major și minor.
- Măsuri calitative: descrierea detaliată a UP individuale și a contextului în care apar.

HE este cea mai răspândită metodă de evaluare, folosită de 76% din specialiștii din domeniul utilizabilității (UPA Survey, 2005).

Este necesară luarea în considerație a expertizei evaluatorilor pentru a estima corect gradul de încredere. Nu numai expertiza generală în utilizabilitate este utilă, dar și specializarea acestora pe anumite aspecte ale sistemelor interactive, atât la nivel general (probleme cognitive, design grafic, sisteme web, intraneturi) cât și de detaliu (arhitectura informației, navigare, procesarea tranzacțiilor). În prezent există preocupări de atestare a evaluatorilor pe niveluri de expertiză, în scopul perfecționării și alocării optime a resursei umane.

Testarea cu utilizatori (UT)

Evaluarea formativă a utilizabilității prin testarea cu utilizatori (UT – User Testing) este o metodă care necesită un număr relativ restrâns de utilizatori (minim 5, ideal 15-25), care testează aplicația pe baza unor sarcini predefinite.

Testarea cu utilizatori este definită ca o metodă empirică intensivă de evaluare a utilizabilității implicând participanți având caracteristici apropiate de cele ale utilizatorilor reali ai produsului care va fi evaluat.

Prin testarea cu utilizatori se înregistrează comportamentul utilizatorului cu ajutorul unor tehnici specifice cum sunt observarea, protocoalele de « gândire cu voce tare » (TA – Think Aloud) sau înregistrările video. Avantajul principal al tehnicii « gândire cu voce tare » o constituie mai buna înțelegere a modelului mental al utilizatorului și interacțiunea cu produsul. În funcție de logistica existentă, se recomandă combinarea mai multor tehnici.

Principala diferență între evaluarea euristică și testarea cu utilizatori, o reprezintă absența utilizatorului în primul caz.

De regulă, în testarea cu utilizatori măsurile colectate sunt:

- Măsuri ale eficacității și eficienței : durata, erorile (numărul de alegeri greșite din meniu, numărul de selecții greșite, alte erori), solicitarea ajutorului (numărul de cereri explicite de asistență din partea experimentatorului și numărul de solicitări de ajutor online)

- Măsuri ale satisfacției : expresiile emoționale (observarea frustrării, observarea confuziei), percepțiile, opiniile și aprecierile participanților.

Măsurile eficacității și eficienței sunt colectate cu ajutorul unui stopwatch și/sau al fișierelor de log. Măsurile satisfacției subiective sunt colectate cu : a) chestionare post-test b) Protocoale de « gândire cu voce tare ».

UT este o metodă de evaluare utilizată frecvent, fiind folosită de 75% din specialiștii în domeniul utilizabilității (UPA Survey, 2005).

PREOCUPĂRI ACTUALE ÎN EVALUAREA FORMATIVĂ

Descrierea problemelor de utilizabilitate

Hvanberg & Law (2003) au propus o schemă de clasificare și un model de descriere a problemelor de utilizabilitate pornind de la schema de clasificare a defectelor (DCS – Defect Classification Scheme), considerând că o problemă de utilizabilitate este similară unui defect (deficiență) a unui sistem software. Schema cuprinde 10 atribute, pentru fiecare dintre ele fiind menționate motivul pentru care a fost ales precum și un set de valori posibile, după cum urmează : identificatorul problemei, descrierea problemei, metoda de identificare (HE / UT), declanșator (descrierea contextului), impact, etapa de dezvoltare în care se anticipează că va fi remediată, tipul problemei, etapa în care a fost remediată (pentru comparație cu etapa anticipată), sugestii pentru prevenire.

Într-o lucrare ulterioară, Vilbersdottir et al. (2006) au extins / modificat schema de clasificare cu câteva atribute. A fost introdus contextul (partea de interfață) separat de declanșator și frecvența (număr de evaluatori care au identificat problema). Au fost adăugate atribute specifice UT : eficiență (timp de execuție sarcină), eficacitate (rată de succes), număr de frustrări și număr de solicitări asistență.

Deși aceste clasificări sunt utile pentru studiul utilizabilității, inclusiv din punct de vedere metodologic, acestea au o valoare practică redusă : majoritatea practicienilor raportează problemele de utilizabilitate într-un format mult mai succint.

Hvanberg et al. (2007), într-o lucrare mai recentă, au utilizat un format structurat care a fost adaptat după Cockton & Woolrich (2001) și care conține : identificatorul UP, descriere, dificultățile întâmpinate de utilizator, contextul specific (localizarea problemei în interfață), cauzele posibile (ce anume este greșit în proiect), euristiciile utilizate și gradul de severitate (impactul, pe trei niveluri : major, moderat și minor).

Consolidarea problemelor de utilizabilitate

Consolidarea problemelor de utilizabilitate este parte integrantă din evaluarea formativă în care sunt implicați mai mulți evaluatori. Conform cu Law & Hvanberg (2008), procesul de consolidare a problemelor de utilizabilitate identificate de mai mulți utilizatori (experimentul se referă la testarea cu utilizatori) se desfășoară în două etape:

- *filtrarea*: eliminarea duplicatelor din lista de probleme identificate de către un singur evaluator

- *integrarea*: combinarea problemelor de utilizabilitate dintre diferite liste identificate de mai mulți evaluatori, cu scopul de a reține problemele unice, relevante.

Prima etapă are rolul de consolidare individuală: fiecare evaluator compară listele de probleme identificate de diferiți utilizatori pentru aceeași sarcină pentru a produce o listă de probleme unice pentru acea sarcină.

A doua etapă are rolul de consolidare colaborativă: doi sau mai mulți evaluatori, care au utilizat aceleași rapoarte de testare, cad de acord asupra unei liste unice de probleme de utilizabilitate pentru fiecare sarcină.

Un studiu asemănător, tot pe baza testării cu utilizatori, a fost raportat de Hornbaek & Frokjaer (2008), Scopul studiului a fost compararea tehnicilor utilizate pentru consolidarea problemelor de utilizabilitate. Tehnicile utilizate au fost următoarele:

- *Modificări similare* în aplicație, în vederea remedierii. Două probleme sunt identice, dacă remedierea uneia o rezolvă și pe cealaltă.
- *Priorizare practică* pentru proiectanți sau dezvoltatori. Presupune o listă ordonată după prioritate și are o valoare practică (așa se procedează).
- *Caracteristicile de descriere a unei probleme de utilizabilitate*: cauză, problemă, rezultat și modificare sugerată.
- *Cadrul de lucru al acțiunii utilizator* (UAF - User Action Framework) care constă într-o schemă de clasificare bazată pe 382 de categorii organizate pe 6 niveluri.

O concluzie a autorilor este că utilizarea de tehnici diferite produce rezultate diferite, atunci când sunt implicați evaluatori începători. Prima metodă a produs un număr mai mare de probleme de utilizabilitate. Cea de a doua metodă a întrunit cel mai mare coeficient de agreere între participanți, în timp ce cea de a treia metodă a întrunit cel mai scăzut consens.

GRADUL DE ÎNCREDERE ȘI VALIDITATE

Indicatori de încredere

Încrederea se referă la măsura în care evaluările independente produc același rezultat. Validitatea se referă la măsura în care problemele detectate în cadrul evaluării sunt aceleași cu cele care apar în practica utilizării sistemului.

Pentru măsurarea gradului de încredere, Hertzum & Jacobsen (2001) recomandă doi indicatori: rata de detecție și coeficientul de agreere între oricare doi evaluatori (any-two-agreement).

Rata de detecție se calculează după formula:

$$\text{Rata detecției} = \text{Media } |P_i| / |PU| \\ \text{pentru toți } n \text{ evaluatori}$$

În această ecuație, P_i este setul de probleme detectate de evaluatorul E_i , PU este setul de probleme detectate colectiv de toți evaluatorii, iar n este numărul de evaluatori. Calculul se face pe baza listei de probleme

unice de utilizabilitate, adică după eliminarea duplicatelor care apar în listele individuale și între aceste liste.

Coeficientul de agreere se calculează după formula:

$$\text{Coeficientul de agreere} = \text{Media} \\ |P_i \cap P_j| / |P_i \cup P_j| \text{ pentru toate} \\ \text{perechile } 1/2n*(n-1) \text{ de evaluatori}$$

În această ecuație, P_i și P_j sunt seturile de probleme detectate de evaluatorii E_i și E_j iar n este numărul de evaluatori. Calculul se face pe baza listei de probleme unice de utilizabilitate, adică după eliminarea duplicatelor care apar în listele individuale și între aceste liste.

Efectul evaluatorului

Practica evaluării arată că atât problemele de utilizabilitate identificate, cât și modul de descriere a acestora pot varia semnificativ de la un evaluator la altul. Această diferență este cunoscută sub numele de *efect al evaluatorului* și este întâlnit atât în evaluarea euristică cât și în testarea cu utilizatori. În mod asemănător, lista problemelor de utilizabilitate diferă semnificativ și de la un utilizator la altul, diferență cunoscută sub numele de *efect al utilizatorului*.

Efectul evaluatorului în evaluarea utilizabilității este definit de Hertzum & Jacobsen (2001) ca situația în care mai mulți evaluatori care apreciază aceeași interfață cu aceeași metodă de evaluare, identifică seturi de UP diferite și grade de severitate diferite pentru aceeași problemă. Efectul evaluatorului este o măsură a gradului de încredere cu care pot fi privite rezultatele evaluării.

Analiza a 11 studii de utilizabilitate efectuate cu 3 metode de evaluare diferite: inspecție cognitivă (CW – Cognitive Walkthrough), evaluare euristică (HE) și «gândire cu voce tare», făcută de Hertzum & Jacobsen (2001), arată că efectul evaluatorului se manifestă pentru:

- Evaluatori experți și începători.
- Probleme de utilizabilitate minore și majore.
- Identificare de probleme și apreciere grad de severitate.
- Evaluare de sisteme simple și sisteme complexe.

Hertzum & Jacobsen (2001) arată că efectul utilizatorului nu poate fi atribuit slăbiciunii unei metode, particularităților unui studiu / artefact și nici nu este un incident întâmplător. Ca atare, nu se pune problema existenței acestui efect ci a cauzelor și căilor în care poate fi diminuat. Principala cauză este faptul că o evaluare a utilizabilității presupune interpretare, care depinde, la rândul ei, de diferențele individuale între evaluatori în ceea ce privește abilitățile cognitive, expertiza, motivația și personalitatea.

Referitor la aplicarea metodelor de evaluare, Hertzum & Jacobsen (2001) au propus trei recomandări pentru a reduce efectul evaluatorului:

- Scopul analizei și selectarea sarcinilor trebuie să fie explicite.
- Este important să fie găsite cele mai multe dintre problemele unui sistem și se recomandă în mod expres să fie implicați mai mulți evaluatori.

- Accent pe procedurile de evaluare și pe criteriile de selecție a problemelor în fiecare caz.

Intr-un alt studiu, Jacobsen et al. (1998) arată că 4 evaluatori au identificat în medie 52% din numărul total de probleme de utilizabilitate.

Indicatori de validitate

Evaluarea euristică produce un set de probleme de utilizabilitate care sunt considerate *anticipate* în timp ce testarea cu utilizatori identifică probleme de utilizabilitate *reale*. Diferența între cele două seturi o reprezintă problemele *fals pozitive* (alarme false, identificate de HE dar neconfirmate de UT) și *fals negative* (problemele de utilizabilitate care nu au fost identificate de HE).

Pentru măsurarea validității, eficacității și eficienței evaluării, se folosesc 3 indicatori (Hartson et al., 2001) : validitate, completitudine (thoroughness) și eficacitate generală (overall effectiveness).

Validitatea se definește după formula:

$$\text{Validitate} = \frac{UP \text{ anticipate}}{UP \text{ anticipate} + UP \text{ fals pozitive}}$$

Completitudinea se definește după formula:

$$\text{Completitudine} = \frac{UP \text{ anticipate}}{UP \text{ reale}}$$

Eficacitatea generală se definește după formula:

$$\text{Eficacitate} = \text{Validitate} * \text{Completitudine}$$

Așa cum se observă, calcularea acestor indicatori presupune evaluarea cu ambele metode, pe baza aceluiași scenariu de utilizare, fapt care mărește costurile evaluării.

CONCLUZII

În ultimii ani asistăm la tendințe de extindere a cadrului de lucru în evaluarea utilizabilității sistemelor informatice, ca urmare a provocărilor societății informaționale:

- Extinderea caracterului multidisciplinar al activității de evaluare determinat de cerințele unor tehnologii specifice, cum sunt cele de e-commerce, e-government și e-learning.
- Studiul măsurilor adecvate pentru evaluarea utilizabilității și pentru înțelegerea relațiilor între diferitele măsuri utilizate în practica evaluării.
- Integrarea evaluării utilizabilității în metodologia de proiectare centrată pe om.
- Creșterea interesului pentru validitatea rezultatelor și adecvarea metodelor de evaluare, inclusiv prin compararea rezultatelor obținute cu diferite metode și utilizarea de metode complementare.

Un aspect important, tot mai des referit în literatura de specialitate din ultimii doi ani, este utilitatea „în aval” a metodelor de evaluare formativă (downstream utility): măsura în care plusul sau minusul de utilizabilitate (îmbunătățirea sau deteriorarea) a unui sistem poate fi atribuit în mod direct remedierilor solicitate de evaluările efectuate (Law et al., 2007).

Confirmare

Această lucrare este finanțată din proiectul de cercetare PS MCSI 49/2008.

REFERINȚE

1. Cockton, G., Woolrych, A. (2001). Understanding inspection methods: lessons from an assessment of heuristic evaluation. Blandford, A., Vanderdonck, J., Gray, P.D. (Eds.), *Proceedings of People and Computers XV*. Springer-Verlag, 171–182
2. Hartson, H.R., Andre, T.S., Williges, R.C., (2001). Criteria for evaluating usability evaluation methods. *International Journal of Human-Computer Interaction* 13, 373–410
3. Hertzum, M., Jacobsen, N.E. (2001) „The evaluator effect: A chilling fact about usability evaluation methods”. *International Journal of Human-Computer Interaction* 13. 421-443.
4. Hornbaek, K & Frokjaer, E. (2008). Comparison of techniques for matching of usability problem descriptions. *Interacting with Computers* 20. 505-514.
5. Hvannberg, E.T. and Law, E.L.-C. (2003) “Classification of Usability Problems (CUP) Scheme”, *Proceedings of Interact 2003*.
6. Hvannberg, E.T. and Law, E.L.-C., Larusdotir, M.C. (2007) „Heuristic Evaluation: Comparing ways of finding an reporting usability problems, *Interacting with Computers* 19, 255-240.
7. Jacobsen, N.E., Hertzum, M., John, B. (1998) „The evaluator effect in usability tests” *Proceedings of ACM CHI’98*. ACM Press. 255–256.
8. Law, L-C. & Hvannberg, E.T. (2002) “Complementarities and convergence of heuristic evaluation and usability test: A case study of UNIVERSAL brokerage platform”. *Proceedings NordiCHI Conference 2002*. ACM. 71-79
9. Law, E.C, Lárusdóttir, M.C, Norgaard, M. (Eds) (2007) *Downstream Utility 2007 : The Good, the Bad, and the Utterly Useless Usability Evaluation Feedback*. IRIT Press – Toulouse.
10. Molich, R., and Nielsen, J. (1990). Improving a human-computer dialogue, *Communications of the ACM* 33, 3, 338-348.
11. Nielsen, J. (1993). *Usability Engineering*. Academic Press, New York.
12. Nielsen, J. (1994). “Heuristic evaluation”. Nielsen, J., and Mack, R.L. (Eds.), *Usability Inspection Methods*, John Wiley & Sons, New York.
13. Theofanos, M. & Quesenbery, W. (2005) “Towards the Design of Effective Formative Test Reports”. *Journal of Usability Studies, Issue 1, Vol.1*. 27-45
14. *UPA 2005 Survey*. Usability Professionals Association.
15. Vilbergtottir, S.G., Law, E.L.-C., Hvannberg, E.T.. (2006) „Classification of Usability Problems (CUP) Scheme: Augmentation and Exploitation, *Proceedings of NordiCHI 2006*, 281-290.