

Support Vector Regression을 이용한 소프트웨어 개발비 예측

박 찬 규*

Estimating Software Development Cost using Support Vector Regression

Chan-Kyoo Park*

■ Abstract ■

The purpose of this paper is to propose a new software development cost estimation method using SVR (Support Vector Regression). SVR, one of machine learning techniques, has been attracting much attention for its theoretic clearness and good performance over other machine learning techniques. This paper may be the first study in which SVR is applied to the field of software cost estimation.

To derive the new method, we analyze historical cost data including both well-known overseas and domestic software projects, and define cost drivers affecting software cost. Then, the SVR model is trained using the historical data and its estimation accuracy is compared with that of the linear regression model. Experimental results show that the SVR model produces more accurate prediction than the linear regression model.

Keyword : Software Cost Estimation, Support Vector Regression, Machine Learning

1. 서 론

정보시스템이 조직의 운영 효율성과 경쟁력 향

상을 위한 필수적인 수단으로 정착됨에 따라 정보 시스템을 구축·유지보수·운영하는 비용이 지속적으로 증가하고 있다. 2004년 기준으로 국내 소프

트웨어산업 시장규모는 연간 약 180억 달러 이상으로 추정될 정도로 급속히 성장하였다. 정보시스템 비용이 차지하는 비중이 높아짐에 따라, 공공기관 및 일반 기업의 정보화 비용의 효과적 관리가 중요한 문제로 대두되었고 이를 위해서는 정보화사업의 소요비용을 정확히 예측할 필요가 있다.

정보시스템 구축사업의 경우 전체 사업비용은 크게 소프트웨어개발비와 하드웨어/패키지 소프트웨어 구입비로 구분된다. 하드웨어/패키지 소프트웨어 구입비는 시장가격 조사를 통해 쉽게 예측할 수 있으므로 본 연구에서는 논외로 한다. 소프트웨어 개발비는 요구사항의 불확실성과 모호성, 소프트웨어 생산성 측정의 어려움, 과거 데이터의 부족, 신기술에 의한 위험 및 개발환경 의존성 등으로 인해 예측이 쉽지 않다[24].

소프트웨어 개발비 예측기법에는 크게 알고리즘 방법(algorithmic method), 전문가예측방법, 유추(analogy)법, 기계학습(machine learning)에 의한 예측방법 등으로 구분될 수 있다[10, 20]. 알고리즘 방법으로는 1965년 SDC(System Development Corporation)모델이 개발된 이후 COCOMO II[6], ISBSG모형[12], ESTIMACS[13], PRICE-S[17], KnowledgePLAN[21] 등 최근까지 많은 공개 또는 상용 소프트웨어 비용산정 모형들이 개발되었다. 또한, 국내에서도 기능점수를 활용한 알고리즘 기반의 소프트웨어 개발비용 예측에 관한 연구들이 김현수[2], 이양규[5], 박찬규[3, 4], 김우제[1] 등에 의해 수행된 바 있다. 유추방법은 과거의 유사한 사업으로부터 유추하여 비용을 예측하는 방법이고, 기계학습에 의한 예측방법은 신경망(neural network), 사례기반추론(case-based reasoning), 의사결정나무(decision tree) 등의 기계학습기법 등을 이용하여 소프트웨어 비용을 예측하는 방법을 말한다.

기계학습에 의한 소프트웨어 개발비용 예측방법으로 가장 먼저 사용된 기법은 신경망에 의한 비용 예측이다. Shepperd[19]는 문제복잡도, 어플리케이션 특이성(novelty of application), 설계도구 사용,

팀크기(team size) 등을 입력요소로 하는 feed-forward 신경망을 통해 소프트웨어 개발비를 예측하는 기법을 제안하였다. Kitchenham[14]은 사례기반 추론방법을 도입하여 소프트웨어 개발비를 예측하였고, 기존의 회귀분석을 기반으로 하는 알고리즘 방법보다 예측 정확성이 향상됨을 보였다. 또한 나무기반(tree-based) 방법으로 회귀나무(regression tree), 의사결정나무(decision tree) 등을 활용하여 소프트웨어 비용예측을 시도한 연구들도 수행되었다[16, 22]. 신경망, 사례기반추론 또는 회귀모형을 이용한 소프트웨어 비용예측의 정확도를 비교하면 신경망에 의한 비용예측이 가장 정확도가 높은 것으로 보고된 바 있다[9].

본 연구에서는 기계학습 분야에서 새롭게 주목을 받고 있는 SVM(Support Vector Machine)을 소프트웨어 개발비 예측에 적용하여 그 정확도를 기존의 방법과 비교하고자 한다. SVM은 Vapnik [23]에 의해 제안된 이후 이론적인 명확성과 우수한 성능이 입증되면서 패턴인식, 텍스트 분류, 생체정보학(bioinformatics), 주가 예측, 침입탐지 등 활용 범위를 점차 넓혀가고 있다[8]. SVM은 주로 분류(classification) 문제에 사용되며 SVM을 회귀분석에 사용할 수 있도록 확장한 방법이 SVR(Support Vector Regression)이다. 아직까지 SVR을 소프트웨어 비용예측 분야에 적용한 국내외 연구가 없었다는 점에서 본 연구의 첫 번째 의의가 있다. 또한, 국내 소프트웨어 시장의 급격한 성장에도 불구하고 국내 상황에 적합한 소프트웨어 개발비 예측을 위한 연구 환경은 미비한 실정이다. 개발비 예측을 위해 필수적으로 요구되는 국내 소프트웨어 사업의 비용자료도 매우 부족한 실정이고, 개발비 예측도 거의 소프트웨어 사업대가기준에 의존하고 있을 뿐이다. 이러한 상황에서 국내 소프트웨어 개발사업의 비용자료와 해외에서 널리 통용되고 있는 비용자료를 종합적으로 분석하여 소프트웨어 개발비 예측을 위한 새로운 도구를 제시한다는 면에서도 본 연구의 의의를 찾을 수 있다.

이후 전개될 본 논문의 구성은 2장에서 SVM과

SVR을 간단히 소개하고, 3장에서는 국내외 소프트웨어 개발비 예측모형을 소개한다. 4장에서는 소프트웨어 비용에 영향을 미치는 비용요소(cost driver)를 정의하고, 5장에서는 전통적인 회귀식을 통한 예측결과와 SVR를 이용한 비용예측 결과를 비교함으로써 SVR이 회귀분석보다 더 정확한 비용예측 방법이 될 수 있음을 보인다. 마지막으로 6장에서 본 연구의 결론과 추후 연구과제 등을 제시한다.

2. SVM(Support Vector Machine)과 SVR(Support Vector Regression)

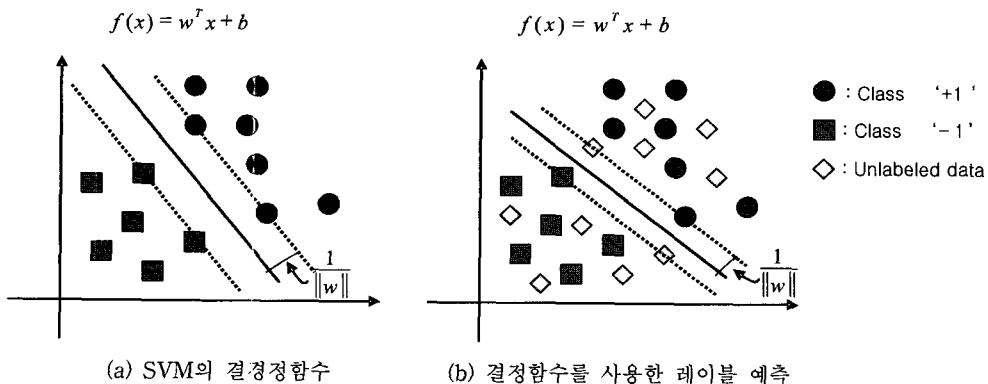
먼저 SVM에 대해 알아보자. 1장에 언급했듯이 SVM은 기계학습에 사용되는 기법이다. SVM은 다차원 상의 점들로 표현되는 학습데이터를 두 개의 그룹으로 분할하는 초평면(hyperplane)을 구하고, 이 초평면을 미지의 데이터가 어떤 그룹에 속할 지를 예측하는 결정함수(decision function)로 사용한다. l 개의 학습데이터 $(x^1, y_1), \dots, (x^l, y_l)$ 가 있고 $x^i \in R^n$, $y_i \in \{1, -1\}$ 이며 모든 y_i 의 값을 알고 있다고 하자. 즉 x^i 는 i 번째 학습데이터의 다차원 상의 좌표이고 y_i 는 i 번째 학습데이터가 '+1' 클래스(class)에 속하는지 또는 '-1' 클래스(class)에 속하는 지를 나타내는 레이블(label)이다. <그림 1>

(a)에서 보는 것처럼 SVM은 학습데이터를 그 레이블 y_i 에 따라 두 개의 클래스로 분할하는 초평면 중에서 두 개의 클래스가 가장 멀리 떨어지도록 분할하는 초평면을 찾는다. 이 때 초평면 $f(x) = w^T x + b$ 가 SVM의 결정함수가 되며, 레이블이 알려지지 않은 데이터 x 에 대해 $f(x) \geq 0$ 이면 '+1' 클래스로, $f(x) < 0$ 이면 '-1' 클래스로 예측한다. 여기서 $(2/\|w\|)$ 를 마진(margin)이라 부르는데 초평면에 의해 분할된 두 클래스간의 최단거리를 나타낸다. SVM은 개념적으로 마진을 최대화 하는 초평면을 결정함수로 사용한다.

SVM의 결정함수 $f(x) = w^T x + b$ 의 w 와 b 를 구하는 문제를 수리계획법으로 모형화하면 다음 (SP_0) 와 같이 이차계획법 문제로 모델링된다.

$$\begin{aligned} \min w^T w \\ (SP_0): \quad \text{s.t. } y_i(w^T x^i + b) \geq 1, \text{ for } i = 1, \dots, l \end{aligned}$$

그러나 현실 문제에서는 <그림 2>(a)처럼 레이블이 다른 학습데이터들이 초평면에 의해 두 그룹으로 완벽하게 분할되지 않고 보통 중첩(overlap)되는 경우가 많으므로 (SP_0) 의 각 제약식을 완화한 모형인 (SP) 를 사용한다. (SP) 의 쌍대문제는 (SD) 와 같은데, 계산상의 여러 가지 장점 때문에 실제로는 (SD) 를 푼다[18]. (SD) 에서 $e = (1, \dots, 1)^T \in R^l$ 이고, $Q_{ij} = y_i y_j (x^i)^T x^j$ 이다.



<그림 1> SVM의 결정함수와 예측

$$\begin{aligned} \min \quad & w^T w + C \sum_i \xi_i \\ (SP) : \quad & \text{s.t. } y_i(w^T x^i + b) \geq 1 - \xi_i, \forall i \\ & \xi_i \geq 0, \forall i \end{aligned}$$

$$\begin{aligned} \max \quad & e^T \alpha - \frac{1}{2} \alpha^T Q \alpha \\ (SD) : \quad & \text{s.t. } y^T \alpha = 0, \\ & 0 \leq \alpha_i \leq C, \forall i \end{aligned}$$

(SP)의 목적함수는 두 개의 항으로 구성되어 있다. 첫 번째 항은 $w^T w = \|w\|^2$ 으로 마진의 역수이다. 결과적으로 $\|w\|^2$ 을 최소화하는 것은 마진을 최대화하는 w 를 찾는 것을 의미한다. 두 번째 항의 $\sum \xi_i$ 는 경험오류(empirical error)로 구해진 결정함수가 주어진 학습데이터에 대해 얼마만큼 오류(error)를 보이고 있는 가를 나타낸다. 여기서 상수 C 는 마진과 경험오류간의 상대적 가중치를 나타낸다. 종합해보면 (SP)의 목적함수는 마진을 최대화하고 경험오류를 최소화하는 결정함수 $f(x)$ 를 찾는 것을 의미한다. 또 (SD)의 최적해를 α^* 라 하면 (SP)와 (SD)의 최적조건(optimality condition)으로부터 w 와 결정함수 $f(x)$ 는 각각

$$\begin{aligned} w &= \sum_i \alpha_i^* y_i x^i \\ f(x) &= w^T x + b = \sum_i \alpha_i^* y_i (x^i)^T x + b \end{aligned}$$

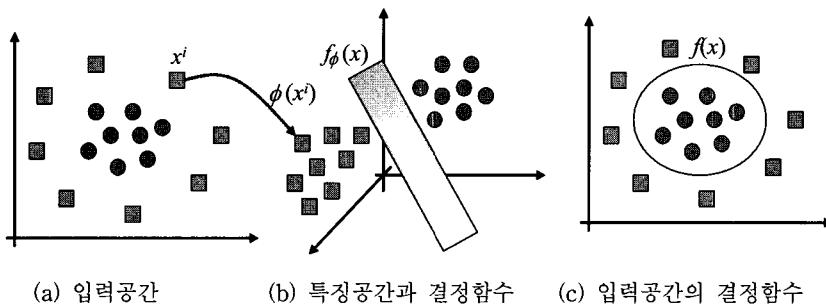
로 표현된다. 이때 $\alpha_i^* > 0$ 인 x^i 를 지지벡터(support vector)라 부른다.

SVM은 학습데이터들을 분할하는 초평면을 결

정함수로 사용하는데 <그림 1>에 나타난 바와 같이 결정함수가 단순히 학습데이터 x^i 들의 선형함수로만 제약된다면 SVM은 성능 면에서 많은 한계를 가지게 될 것이다. 이러한 문제점을 극복하기 위해 <그림 2>(a)와 같이 입력공간(input space)에 있는 학습데이터 x^i 를 고차원의 특징공간(feature space)의 점 $\phi(x^i)$ 로 변환한 다음 (b)와 같이 특징공간의 점들을 그 레이블에 따라 두 개의 그룹으로 분할하는 결정함수 $f_\phi(x)$ 를 구한다. 특징공간의 결정함수 $f_\phi(x)$ 을 원래의 입력공간으로 변환하면 (c)와 같이 비선형함수 형태를 띠게 된다.

학습데이터 x^i 를 특징공간의 점으로 변환하는 과정은 함수 $\phi: x \rightarrow \phi(x^i)$ 를 통해 이루어진다. 그 결과로 특성공간에서 정의되는 SVM의 쌍대문제 (SD)를 풀기 위해서는 $Q_{ij} = \phi(x^i)^T \phi(x^j)$ 를 계산해야 한다. 만약 $\phi(x^i)$ 를 명시적으로 구하여 Q_{ij} 를 계산하려면 $\phi(x^i)$ 의 차원이 매우 높거나 또는 무한 차원인 경우도 있으므로 많은 계산시간이 소모될 것이다. 따라서 이러한 문제점을 해결하기 위해 함수 ϕ 를 직접 구할 필요 없이 커널함수(kernel function)를 통해 특징공간으로의 변환이 이루어진다. 커널함수는 $\phi(x^i)$ 와 $\phi(x^j)$ 를 명시적으로 구하지 않고도 $Q_{ij} = \phi(x^i)^T \phi(x^j)$ 를 쉽게 계산할 수 있게 해주는데, 커널함수로 가장 많이 사용되는 함수는 RBF(Radial Basis Function)커널, 다항(polynomial)커널 등이 있다. RBF 커널은

$$k(x^i, x^j) = \exp\left(-\frac{\|x^i - x^j\|^2}{\sigma^2}\right)$$



<그림 2> 입력공간과 특징공간

이고, 다항 커널은

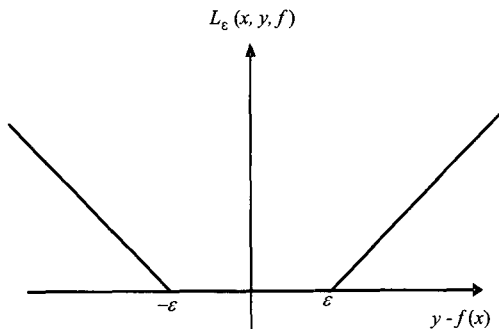
$$k(x^i, x^j) = ((x^i)^T x^j)^p$$

이다. RBF 커널은 학습데이터 $x^i \in R^n$ 를 무한 차원의 특징공간의 점으로 변환하고, 다항커널은 입력데이터를 $\binom{n+d-1}{d}$ 차원의 특징공간의 점으로 변환한다.

SVM은 학습데이터를 '+1' 클래스와 '-1' 클래스로 구분하는 데 사용되지만, 임의의 실수값을 예측할 수 있도록 SVM를 일반화한 방법이 SVR이다. (SP)의 목적함수에 있는 $\sum_i \xi_i$ 대신에 SVR에서는 다음과 같은 ϵ -무감도 손실함수(ϵ -insensitive loss function) $L_\epsilon(x, y, f)$ 를 사용한다.

$$\begin{aligned} L_\epsilon(x, y, f) &= \max(0, |y - f(x)| - \epsilon) \\ &= \max(0, |y - (w^T x + b)| - \epsilon) \end{aligned}$$

ϵ -무감도 손실함수는 실제값 y 와 예측된 값 $f(x) = w^T x + b$ 간의 오차가 ϵ 이하면 $L_\epsilon(x, y, f) = 0$ 이고 오차가 ϵ 보다 크면 오차의 절대값에서 ϵ 만큼 차감한 값인 $L_\epsilon(x, y, f) = |y - f(x)| - \epsilon$ 을 갖게 된다. 이를 그림으로 나타내면 <그림 3>과 같다.



<그림 3> ϵ -무감도 손실함수

결국 SVR에서는 실제값 y ,와 예측값 $f(x) = w^T x + b$ 의 값을 가능한 한 ϵ 이내로 유지하면서 마진을 최대화하게 되는데 이를 수리계획모형으로 나타내면 다음과 같다[8].

$$\min w^T w + C \sum_i (\xi_i^+ + \xi_i^-)$$

$$\begin{aligned} (RP) : \quad & \text{s.t. } (w^T x^i + b) - y_i \leq \epsilon + \xi_i^+, \forall i \\ & y_i - (w^T x^i + b) \leq \epsilon + \xi_i^-, \forall i \\ & \xi_i^+, \xi_i^- \geq 0, \forall i \end{aligned}$$

$$\max \sum_i y_i (\alpha_i^- - \alpha_i^+) - \epsilon \sum_i (\alpha_i^- + \alpha_i^+)$$

$$\begin{aligned} (RD) : \quad & -\frac{1}{2} \sum_{i,j} (\alpha_i^- - \alpha_i^+) (\alpha_j^- - \alpha_j^+) Q_{i,j} \\ & \text{s.t. } \sum_i (\alpha_i^- - \alpha_i^+) = 0, \\ & 0 \leq \alpha_i^+, \alpha_i^- \leq C, \forall i \end{aligned}$$

SVR도 SVM에서와 마찬가지로 커널함수를 써서 학습데이터를 특징공간의 점으로 변화시킨 다음 특징공간에서 학습을 수행하게 되는데 이때 특징공간에서의 (RD)에 나타나는 $Q_{i,j} = \phi(x^i)^T \phi(x^j)$ 값은 SVM과 동일하게 커널함수 $k(x^i, x^j)$ 를 사용하여 계산하게 된다.

3. 국내외 소프트웨어 개발비 예측모형

SVR을 이용하여 소프트웨어 개발비를 예측하기 위해서는 먼저 소프트웨어 개발비용에 영향을 주는 요인들을 식별해야 한다. 이를 위해 기존에 개발된 여러 가지 알고리즘 기반의 예측모형을 분석하여 소프트웨어 개발의 비용요소를 추출해 본다.

일반적으로 알고리즘 기반의 예측모형은 먼저 개발하고자 하는 소프트웨어의 규모를 산정하고, 시스템의 특성이나 프로젝트 환경 등을 반영하는 보정요소를 고려하여 소프트웨어 개발비를 예측하게 된다. 소프트웨어 규모는 LOC(Line of Code) 또는 기능점수(function point) 방식을 가장 널리 사용하고 있는데, 예측의 편이성, 측정결과의 정확성 및 일관성 측면에서 기능점수 방식이 우수하여 본 연구에서는 기능점수 방식의 소프트웨어 규모 측정만을 고려하기로 한다. 기능점수 방식에 관한 상세한 설명은 [11]를 참조하기 바란다.

〈표 1〉 여러 가지 모형의 보정요소

모형	개발자	주요 보정요소
COCOMO	B. W. Boehm	소프트웨어 신뢰성, DB크기, 제품복잡도, 실행시간제약, 주기억장치제약, 가상기계가변성, 반응시간, 분석능력, 어플리케이션경험, 프로그래머능력, 가상기계경험, 개발언어경험, 최신프로그래밍기술사용 정도, 소프트웨어 도구사용 정도, 개발일정
IFPUG	IFPUG	데이터 통신, 분산처리, 성능, 하드웨어 제약, 처리율, 온라인데이터입력, 최종사용자 용이성, 온라인 갱신, 처리복잡도, 재사용성, 설치용이성, 운영용이성, 복수사이트, 변경용이성
ISBSG	ISBSG	플랫폼(메인프레임, 중간, PC), 개발언어(3GL, 4GL, ApG),
PRICE-S	PRICE Systems	응용도메인(상용, 군용) 통합수준, 개발언어, 생산성요소, 복잡도
Knowledge PLAN	Caper Jones	플랫폼, 응용도메인, 복잡도, 개발언어, 인력구성, 사용기술, 소프트웨어 도구, 개발환경, 보안/성능 요구수준, 문서화수준, 테스트 복잡도

1965년 미국 SDC(Systems Development Corporation)사에서 169개의 소프트웨어 프로젝트를 분석하여 104개의 비용보정요소 중 통계적으로 의미가 있는 14개의 보정요소를 갖는 소프트웨어 비용예측모형을 개발한 후 현재까지 많은 비용예측모형이 개발·개선되어 왔다. 그 중에서 가장 널리 사용되고 있는 공개모형인 COCOMO/COCOMO II 모형[6], IFPUG 모형[11], ISBSG 모형[12]과 상용모형인 PRICE-S[17], KnowledgePLAN[21]에 포함된 보정요소를 간략하게 정리하면 <표 1>과 같다. <표 1>에서 IFPUG 모형은 원래는 시스템의 특성을 반영하여 기능점수를 산정하는 방법이지만, 시스템 일반특성에서 고려하는 14 요소가 다른 모형의 보정요소와 같다고 보았다.

<표 1>에 포함된 모형 중에서 공개된 모형은 COCOMO와 ISBSG 모형이다. IFPUG 모형은 공개되어 있긴 하지만 시스템 특성을 반영한 기능점수 규모 산정 과정까지만 제시할 뿐 비용을 추정하는 모형은 제시하지 않고 있다. COCOMO 모형의 비용예측 공식은 다음과 같다.

$$PM = a \times S^b \times \prod_{i=1}^{15} f_i \quad (3.1)$$

PM는 소요인력(person-month)을 나타내고, S는 소프트웨어 규모이고, a, b는 프로젝트 특성과 개발유형에 따라 정해지는 상수이다. 또한, f_i 는

노력승수(effort multiplier)라고 하는데, 프로젝트의 특성과 여러 가지 보정요소를 고려하여 결정되는 값이다. COCOMO에서 사용하는 보정요소에는 15가지가 있으며 각 보정요소의 복잡도를 매우 낮음(very low), 낮음(low), 보통(nominal), 높음(high), 매우 높음(very high), 극히 높음(extra high), 6단계로 구분하여 가중치를 주고 있다. 예를 들어, 요구되는 소프트웨어 신뢰성이 매우 낮은 경우 0.75에서부터 매우 높은 경우 1.40까지의 값을 갖고 있다.

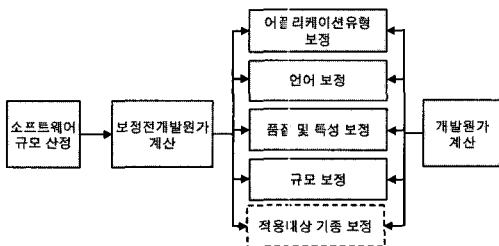
ISBSG 모형은 ISBSG에 수집한 약 731개의 소프트웨어 비용자료를 분석하여 제시한 회귀모형으로서 프로젝트 인도율(delivery rates), 소요노력(effort), 프로젝트기간(duration) 등을 예측할 수 있다. 소요노력을 예측하는 식은 다음과 같다[12].

$$PWE = C \times S^{E_1} \times M \times Team^{E_2} \quad (3.2)$$

$$PWE = C \times S^E \quad (3.3)$$

식 (3.2)는 소프트웨어의 기능점수 규모와 개발팀의 최대 크기를 알고 있을 경우 개발비를 예측하는 식이고 식 (3.3)은 기능점수 규모만 알고 있을 경우 개발비를 예측하는 식이다. 식 (3.2)와 식 (3.3)에 포함된 상수 C, E, E_1 , E_2 등의 값은 보정요소에 따라 달라진다. 예를 들어 식 (3.3)에서 플랫폼이 PC이고 4GL을 사용하는 경우 $C=3.38$, $E=0.974$ 를 사용하여 개발비를 예측하게 된다.

다양한 비용예측 모형들이 개발·활용되고 있는 해외와는 달리 국내에서는 정보통신부에서 제정·고시한 소프트웨어사업대가기준을 통일적으로 사용해 왔다. 소프트웨어사업대가기준은 1989년 제정된 이래 공공기관에 발주하는 소프트웨어 사업의 대가를 산정하는 기준으로 활용되어 왔으며, 소프트웨어사업에 관한 전문지식 없이도 소프트웨어 사업의 예산을 산정할 수 있는 통일적 기준을 제시함으로써 소프트웨어사업의 발주와 관리를 용이하게 하는데 기여하였다. 그러나 산정결과의 정확성 면에서 꾸준한 문제점이 제기되고 있어 지속적인 보완이 필요한 실정이다. 소프트웨어사업대가기준에 따른 소프트웨어 개발비 예측 절차는 <그림 4>와 같다.



<그림 4> 소프트웨어사업대가기준의 개발비 산정 절차

기능점수방식으로 개발 대상 소프트웨어의 규모를 산정한 다음 기능점수 단가를 곱하여 보정전 개발원가를 계산하고, 이를 보정하여 최종적으로 개발원가를 계산한다. <그림 4>에서 알 수 있듯이 소프트웨어사업대가기준에서는 여러 가지 보정요소를 포함하고 있다. 어플리케이션 유형, 언어, 적용대상 기종, 규모 보정은 2004년 이전의 소프트웨어사업대가기준에 포함된 보정요소이고, 2004년에 개정된 사업대가기준에는 적용대상 기종 보정이 삭제되고 품질 및 특성 보정요소로 통합되었다. 식 (3.1)~식 (3.3)의 예측모형은 규모가 증가함에 따라 소요노력이 지수적으로 증가하지만, 소프트웨어사업대가기준의 기능점수 단가는 규모에 상관없이 일정하다. 따라서 이러한 문제점을 해결할 수 있도

록 포함된 보정요소가 규모보정이다.

4. 비용요소 선정

3장에서 알아보았듯이 소프트웨어 개발비 예측을 위해서는 먼저 고려할 비용요소를 정해야 한다. 비용요소 선정 시에는 개발비와의 관련성뿐만 아니라 데이터의 획득가능성, 다른 비용요소와의 상관관계 등이 종합적으로 고려되어야 한다. 본 연구에서는 해외 소프트웨어 개발사업 비용자료로 ISBSG 데이터와 국내 소프트웨어 개발사업 비용자료로 한국전산원에서 구축한 데이터(이하 '국내 비용자료'라 부른다)를 활용한다. 먼저 데이터의 획득가능성을 확인하기 위해 ISBSG 데이터와 국내 비용자료의 주요 수집항목을 정리하면 <표 2>와 같다. 비용과 전혀 관련이 없는 수집자료는 <표 2>에서 제외하였다.

<표 2> ISBSG 및 국내 비용자료의 수집항목

구 분	주요 수집항목
ISBSG	프로젝트 ID, 규모산정 방법, 기능점수, 기능점수보정인자, 기능점수측정방식, 소요노력(총괄), 소요노력(측정수준), 데이터품질, 최대팀 크기, 개발유형, 개발플랫폼, 언어유형, 주프로그래밍언어, 사용 DBMS, CASE 사용여부, 방법론사용여부, 개발방법, 사업기간, 사업중단기간, 구현일자, 불량개수, 사용자수(사업, 위치, 동시사용자 기준), 조직 유형, 업무영역 유형, 어플리케이션 유형, 패키지사용여부, 사업범위
국내 비용자료	사업번호, 사업연도, 사업비(예산, 계약, 개발비, 장비도입비 등), 사업기간, 사업범위, 어플리케이션 유형, 프로그램본수, 기능점수, 기능점수보정인자, 개발언어, 입력화면수, 보고서수, 출력화면수, 적용대상 기종, 소요인력

<표 2>에서 볼 수 있듯이 ISBSG는 비용에 영향을 미칠 가능성이 있는 다양한 자료들을 수집하고 있지만 회원들의 자발적인 참여로 수집된 자료이기 때문에 많은 결측치(missing value)를 포함하고 있다. 반면에 국내 비용자료는 사업결과물로

부터 획득 가능한 제한된 항목들만을 포함하고 있지만, 대부분의 사업에 대해 결측치 없이 완전한 비용자료를 갖고 있다. ISBSG 데이터와 국내 비용자료에 공통으로 수집되는 항목은 <표 2>에 굵게 밑줄 처진 항목들이다. 공통항목들에 대해 살펴보고 이들을 어떻게 예측모형에 반영할 수 있는지 알아본다.

4.1 소프트웨어 규모

먼저 소프트웨어 규모는 ISBSG와 국내 비용자료 모두 기능점수 방식에 의해 계산되었다. 국내비용자료는 모두 IFPUG 측정규칙에 따라 기능점수 규모를 측정하였다. 반면, ISBSG는 IFPUG 측정규칙뿐만 아니라 NESA, MARK II 등 여러 가지 다른 기능점수 규칙에 따라 측정된 데이터들도 포함되어 있다. 국내 비용자료와의 호환을 위해 IFPUG 규칙 이외의 다른 기능점수 규칙에 의해 소프트웨어 규모를 측정한 데이터는 제외하기로 한다. 또한, 국내비용자료에는 보정전(unadjusted) 기능점수와 기능점수 보정인자(adjustment factor)가 모두 조사되어 있으나, ISBSG에는 기능점수 보정인자가 없이 단지 보정후 기능점수만 포함된 자료들이 많다. 따라서 보다 많은 데이터를 분석에 포함할 수 있도록 보정후 기능점수를 소프트웨어 규모로 사용하기로 한다.

4.2 소요인력

ISBSG의 소요노력(work effort)은 투입된 인력의 총시간으로 표시되어 있으나, 국내비용자료의 소요인력은 인-월(man-month)로 표시되어 있다. 따라서 소요노력의 측정단위를 시간으로 통일하기 위해 1개월을 4주로 보고 주 44시간 근무를 기준으로 인-월 단위의 소요인력을 시간으로 변환하였다.

4.3 개발플랫폼

ISBSG의 개발플랫폼(development platform)은

대상사업의 하드웨어 유형을 PC, 중형(mid-range), 대형(mainframe) 등 3 가지로 분류한다. 반면에 소프트웨어사업대기기준에 따라 국내비용자료의 적용대상 기종은 PC, 워크스테이션(client-server, UNIX 시스템 포함), 중대형 등 3가지로 분류되어 ISBSG의 분류기준과는 약간 다르다. 따라서 ISBSG의 개발플랫폼과 국내비용자료의 적용대상 기종을 종합하여 개발플랫폼을 PC와 중대형 2가지로 분류하기로 한다. 중대형은 대부분 UNIX 계열의 클라이언트-서버(client-server) 구조를 포함하며 비용에 미치는 영향이 PC 기반 플랫폼과 다를 수 있음을 고려하였다.

4.4 언 어

ISBSG의 언어유형은 2GL, 3GL, 4GL, AP(Application Generator) 4가지로 분류되고, 주프로그램밍언어 항목은 프로젝트에 주로 사용된 언어를 의미한다. 그러나 ISBSG 데이터에 포함된 대부분의 사업에서 주프로그램밍언어 항목이 빠진 채 언어유형만 표시되어 있다. 또 국내비용자료의 개발언어에는 사업에 사용된 언어가 모두 표기되어 있으나, 전체 소프트웨어 중 어느 정도 비율로 각 언어가 사용되었는지는 알 수 없는 실정이다. 따라서 사용언어를 ISBSG와 같이 2GL, 3GL, 4GL, AP로 구분하고 국내사업도 이와 같이 재분류하였다.

4.5 사업기간과 사업시작연도

사업기간은 ISBSG와 국내비용자료 모두 개월수로 표시되어 있다. 사업범위는 ISBSG와 국내비용자료 모두 소프트웨어개발 기간의 어느 단계에 해당하는 사업인가를 나타내는데 요구사항분석에서 시스템 개발, 시험, 설치까지 전단계를 포함하는 사업만을 분석대상으로 한다. 또한 ISBSG의 구현일자와 국내비용자료의 사업기간으로부터 소프트웨어 개발연도를 유추해 낼 수 있다. 시간에 지남에 따라 소프트웨어 생산성이 변화할 수 있고 이로 인해 소요인력의 변화 여부를 검증하기 위해 사

업수행연도를 예측모형에 포함하기로 한다. 1980년을 0으로 하여 1년 단위로 사업수행연도를 표시하여 분석한다.

4.6 어플리케이션 유형

ISBSG에서는 어플리케이션 유형을 MIS, 업무/생산시스템, 사무정보시스템, 의사결정지원, 공정 제어 등 10가지 이상 다양한 형태로 분류하고 있다. 반면 국내 비용자료는 소프트웨어사업대가기준에 따라 업무처리용, 과학기술용, 멀티미디어용, 지능정보용, 시스템용, 통신제어용, 공정제어용, 지휘통제용 8가지로 분류하고 있다. ISBSG의 분류와 국내 비용자료의 분류가 일치하는 부분은 그대로 사용하고 서로 다르게 분류된 부분은 재분류하여 다음 <표 3>과 같이 어플리케이션 유형을 구분하기로 한다. <표 3>에서 지휘통제용, 과학기술용 등 ISBSG 및 국내비용자료에 나타나지 않은 어플리케이션 유형은 제외하였다.

4.7 조직 유형

발주기관에 따른 소프트웨어 개발비 영향여부를 파악하기 위해 조직유형에 관한 비용요소를 도입한다. ISBSG에는 발주기관 조직유형을 보험, 은행, 금융업, 정부 및 공공기관, 제조업, 도소매업, 컴퓨터, 운송업, 통신업, 전기/가스/수도 등 여러 가지 유형으로 분류하고 있다. 이를 유사한 몇 개의 그룹으로 재분류할 필요가 있는데, 본 연구에서는 보험/은행/금융업, 공공기관/공공재분야, 제조업, 컴퓨터/통신업, 기타 등 5가지 분야로 조직유형을 분류하여 조직유형간 개발비 차이 유무를 분석한다. 국내 비용자료는 모두 공공기관으로 조직유형을 분류하였다.

지금까지 논의된 비용관련 항목을 5장의 SVR 모형의 독립변수와 종속변수로 구분하여 정리해 보면 <표 4>과 같다.

마지막으로 ISBSG의 개발유형(development type)은 프로젝트가 신규개발, 기능개선, 재개발

<표 3> 어플리케이션 유형 분류

ISBSG 분류	국내비용자료 분류	본 연구의 분류
MIS, 업무처리/생산시스템, 사무정보시스템	업무처리용	업무처리 및 MIS용
네트워크 관리, 실시간, 공정 제어, 통신제어	시스템용, 통신제어용, 공정제어용	시스템관리/공정제어용
의사결정지원시스템	지능정보용	의사결정지원용
멀티미디어	멀티미디어용	멀티미디어

<표 4> 비용예측모형의 독립변수와 종속변수

독립 변수		종속 변수	
변수명	단위 또는 값	변수명	단위 또는 값
SW규모	기능점수(Function Point)	소요인력	시간(man-hour)
개발플랫폼	PC기반, 중대형		
언어	2GL, 3GL, 4GL, AP		
사업기간	개월(month)		
사업시작연도	연(year)		
어플리케이션 유형	업무처리 및 MIS용, 시스템관리/공정제어용, 의사결정지원용, 멀티미디어용		
조직유형	보험/은행/금융업, 공공기관/공공재분야, 제조업, 컴퓨터/통신업, 기타		

중 어디에 해당하는가를 의미하는데 국내 비용자료는 모두 신규개발 사업을 대상으로 수집된 자료이다. 따라서 ISBSG의 데이터 중 개발유형이 기능개선이거나 재개발인 사업은 모두 제외하고 신규개발인 자료만을 분석대상으로 한다.

5. SVR 적용결과

본 절에서는 4장에서 선정된 비용요소를 독립변수로 하여 소프트웨어 개발비를 예측한 결과를 제시한다. 먼저 불완전한 데이터 및 이상치(outlier)를 갖는 데이터 제거 과정을 설명하고 최종적으로 분석대상으로 남은 데이터에 대한 분포를 설명한다. 또한 SVR 적용 결과를 제시하기 이전에 SVR에 의한 예측 정확성을 비교해 볼 수 있는 기준으로 선형회귀모형에 의한 예측결과를 제시하고, 마지막으로 SVR 적용결과를 제시한다.

5.1 적용데이터

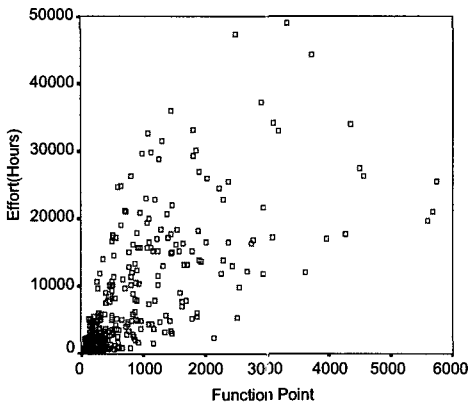
4장에서 비용요소로 선정된 속성에 관한 측정값이 없는 자료, 소프트웨어 개발 프로젝트로 보기 어려운 사업 등 분석 대상에 포함되기 부적합한 자료를 제거하였는데, ISBSG와 국내비용자료에서 제거된 데이터들과 제거 기준을 요약하면 <표 5>

과 같다. 총 1,238개의 ISBSG 데이터 중 908개가 제거되었고 남은 데이터는 332개이고, 총 76개의 국내비용자료에서 10개가 제거되고 남은 데이터는 66개이다. 결과적으로 <표 5>과 같이 916개의 자료가 제거되고 남은 자료는 총 398개이다.

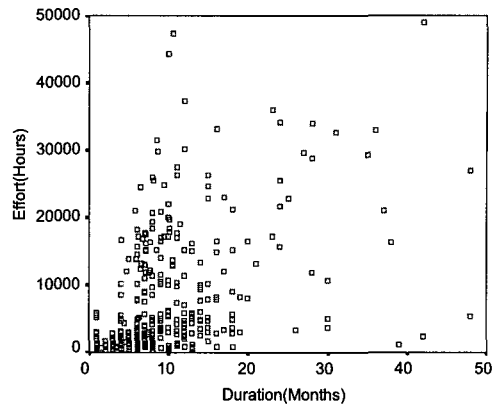
다음으로 이상치에 의해 예측 결과가 왜곡되는 것을 방지하기 위해 다른 데이터에 비해 매우 작거나 큰 속성값을 갖는 자료를 제거하기로 한다. 첫 번째로 기능점수 규모가 매우 작거나 또는 매우 큰 사업들을 제외한다. <표 5>의 결과로 남은 총 398개의 데이터의 기능점수 규모 분포는 11FP~26,968FP이다. 기능점수 규모가 다른 사업에 비해 매우 작거나 큰 사업은 제외하는데, 평균으로부터 2.58σ를 벗어나는 6,000FP 이상인 사업과 50FP 이하인 사업 13개를 제외한다. 두 번째로 기능점수 당 소요노력 비율이 1Hour/FP로 매우 작은 사업과 45Hours/FP로 평균으로부터 2.58σ보다 큰 사업 22개를 제외하였다. 마지막으로 소요노력이 60,000Hours로 평균으로부터 2.58σ를 벗어나 다른 데이터에 비해 현격히 큰 사업 6개를 제외하였다. 따라서 ISBSG 데이터와 국내비용자료 중에서 불완전한 자료와 이상치를 제외하고 분석 대상으로 선정된 데이터는 총 347개이다. 이들 데이터의 기능점수 규모와 소요노력, 사업기간과 소요노력의 분포를 그래프로 나타내면 <그림 5>, <그림 6>과 같다.

<표 5> 제거 기준 및 제거 사업 개수

ISBSG 데이터		국내비용자료	
제거 기준	제거 사업수	제거 기준	제거 사업수
IFPUG 이외의 규모산정	120개	DB구축 인력 비율이 높은 사업	3개
기능점수 규모가 없는 자료	7개		
기능개선, 재개발, 패키지 구매	647개	소요노력 자료가 없는 사업	6개
개발 플랫폼, 언어 유형 결측 자료	38개		
사업기간이 없는 자료	61개		
어플리케이션 유형, 조직 유형, 사업분야 유형이 없는 자료	32개	동일 사업이나 기능점수 측정범위에 따라 2개로 구분된 사업	1개
언어유형이 2GL인 자료	1개		
사업시작연도가 없는 자료	10개		
합 계	916개	합 계	10개



〈그림 5〉 기능점수와 소요노력



〈그림 6〉 사업기간과 소요노력

5.2 선형회귀모형에 의한 비용예측

4장에서 선정된 비용요소들을 고려하여 식 (5.1)과 같이 회귀모형을 구성하였다. 예측할 종속변수 E 는 소요노력이고, S 는 소프트웨어의 기능점수 규모이다. T 는 사업기간, Y 는 사업시작연도, P 는 개발플랫폼, L 는 개발언어, A 는 어플리케이션 유형, O 는 조직유형, K 는 ISBSG 데이터와 국내비용자료간의 생산성 차이를 반영하기 위해 도입한 변수이다. 소문자로 표시된 a, b, c, d 는 회귀모형에 의해 도출되는 계수들이다.

$$E = a \times S^b \times T^c \times Y^d \times P \times L \times A \times O \times K \quad (5.1)$$

식 (5.1)은 네 가지 가정을 토대로 도출되었다. 첫 번째 가정은 소요노력은 소프트웨어 규모가 증가함에 따라 지수적으로 증가한다는 것이다. 이 가정은 대부분의 기존 소프트웨어 개발비용 예측모형에서 사용되고 있다. 두 번째 가정은 소프트웨어 개발기간을 단축시킴에 따라 필요한 소요노력은 지수적으로 증가한다는 것이다. 이 가정 또한 SLIM, ISBSG 등 대부분의 예측모형에 통용되고 있다. 세 번째 가정은 소프트웨어 개발 생산성은 개발연도가 최근에 가까울수록 지수적으로 증가한다는 것이다. 이는 소프트웨어 개발 생산성이 매년 균일한 비율로 개선됨을 의미하는데, 소프트웨어

개발 생산성이 시간이 지남에 따라 개선되기는 하지만 균일한 비율로 개선되지는 않기 때문에 현실적으로 성립되기 어려운 가정이다. 그러나 소프트웨어 개발 생산성의 개선비율을 매년 다르게 하여 회귀모형을 구성하는 것은 이론적으로나 계산복잡도면에서 간단치 않다. 따라서 개선비율의 기하평균을 선형회귀모형에 반영시키기 위해서는 세 번째 가정이 필요하다. 네 번째 가정은 개발플랫폼, 언어유형, 어플리케이션유형, 조직유형 등이 달라짐에 따라 일정 비율만큼 소요노력에 영향을 준다는 가정이다. 이는 COCOMO 모형, 소프트웨어사업대가기준 등에서 일반적으로 사용되고 있는 가정이다. 마지막으로 비용자료의 국내의 구분에 따라 일정 비율만큼 소요노력이 달라진다는 가정이다. 이는 국내 소프트웨어사업과 해외 소프트웨어사업 간의 생산성 차이를 반영하기 위한 가정이다.

식 (5.1)을 선형회귀식으로 변환하기 위해 양변에 상용로그(log)를 취하고, 이후 개발플랫폼, 언어유형, 어플리케이션유형, 조직유형, 국내의 구분 등의 변수는 0 또는 1의 값을 갖는 더미변수(dummy variable)를 도입하여 표현하였다. 예를 들면, 개발언어의 경우 3GL, 4GL, AP 등 3개의 언어 유형을 고려하므로 더미변수 $LANG_{4GL}$ 와 $LANG_{AP}$ 를 도입하여 언어 유형이 4GL이면 $LANG_{GL} = 1$, $LANG_{AP} = 0$ 으로 나타내고, 언어유형이 AP이면

$LANG_GL = 0$, $LANG_AP = 1$ 로 나타낸다. 언어유형이 3GL이면 두 더미변수 모두 0을 갖도록 표현하였다.

회귀모형의 계수를 구하기 이전에 독립변수들간의 상관관계를 검토한 결과 조직유형과 국내외 구분간의 상관성, 기능점수와 사업기간 사이의 상관성이 매우 높은 것으로 나타났다.

<그림 7>에서와 같이 국내비용자료가 전부 공공기관을 상대로 수집되어서 총 분석대상 데이터 347개 중에서 90개의 자료가 공공기관의 프로젝트이고 그중 58개가 국내비용자료이기 때문이다. 상관분석에서도 국내외 구분과 조직유형간의 상관관계가 유의한 것으로 나타났다. <그림 8>에서 나타나듯이 사업기간과 기능점수 간에 상관관계가 통계적으로 유의한 것으로 나타났다. 또한 회귀식의 계수를 구해 본 결과 조직유형과 사업기간은 소요 노력에 통계적으로 영향을 미친다고 볼 수 없는 것으로 나타나 조직유형과 사업기간은 회귀모형에서 제외하기로 한다. 조직유형을 제거한 후 회귀식의 계수를 다시 구해본 결과 어플리케이션 유형과 사업수행연도의 계수가 모두 통계적으로 유의하다고 볼 수 없어 어플리케이션 유형과 사업수행연도도 회귀모형에서 제거하기로 한다. 결과적으로 다음과 같은 선형회귀모형과 회귀계수를 얻을 수 있었다. 식 (5.2)의 회귀모형에서 $R^2 = 0.674$ 이다.

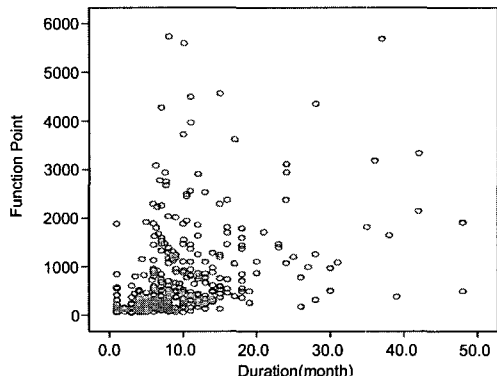
$$\log E = a + b \log S + c \log P + d \log L + e \log K \quad (5.2)$$

식 (5.2)에 표현된 회귀모형의 기본가정이 성립되는가를 살펴보기 위해 공선성(multi-collinearity), 이분산성(heteroscedasticity), 자기상관성(autocorrelation) 존재 여부를 확인해야 하는데, 분석대상 자료가 시계열 자료가 아니므로 그중에 자기상관성은 제외하기로 한다. 공선성 존재여부는 공차한계(tolerance)를 통해 살펴볼 수 있는데, 기능점수 규모, 개발플랫폼, 개발언어, 데이터의 국내외 여부 등이 회귀식에 각각 추가될 때 공차한계의 최소값은 0.688이므로 공선성이 존재한다고 볼 수 없다. 또한 이분산성 존재 여부를 확인하기 위해 회귀식의 잔차를 히스토그램과 정규 P-P도표로 그려보면 <그림 9>, <그림 10>과 같다. <그림 9>에서 보면 잔차가 정규분포에서 크게 벗어나지 않는 것을 알 수 있고, <그림 10>을 보면 관측누적확률과 기대누적확률이 거의 일치하는 것을 볼 수 있다. 따라서 이분산성도 존재한다고 볼 수 없으므로 식 (5.2)는 회귀모형의 기본가정을 모두 충족하고 있다.

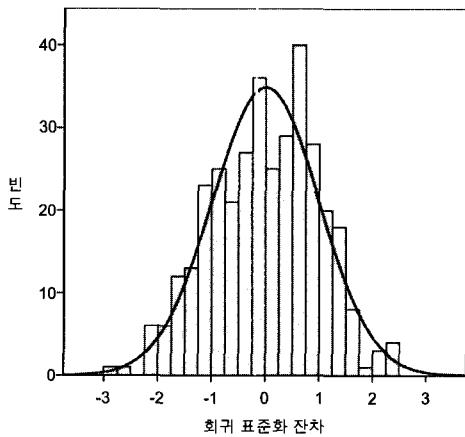
다음으로 <표 6>에 나타난 회귀식의 계수를 살펴봄으로써 식 (5.2)의 선형회귀모형이 의미하는 바를 고찰해보자. 먼저 회귀식 계수 b 는 기능점수 규모의 증가에 따라 개발에 필요한 소요노력이 어느 정도로 증가하는 가를 나타낸다. 일반적으로 $b > 1$ 이면 규모의 비경제효과가 나타남을 의미하고 $b < 1$ 이면 규모의 경제 효과가 나타남을 의미하는데 본 연구에서는 $b = 0.827$ 를 얻어 분석대상 사업

		국내외 구분		합 계
		국내	국외	
조직유형	보험/은행 금융업 분야	0(0%)	110(31.7%)	110(31.7%)
	공공기관/ 공공재 분야	58(16.7%)	32(9.2%)	90(25.9%)
	제조업	0(0%)	42(12.1%)	42(12.1%)
	컴퓨터/ 통신업 분야	0(0%)	20(5.8%)	20(5.8%)
	기타	0(0%)	85(24.5%)	85(24.5%)
합 계		58(16.75%)	289(83.3%)	347(100%)

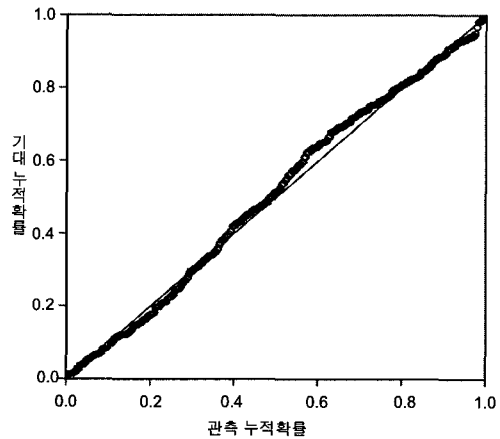
<그림 7> 조직유형의 분포



<그림 8> 사업기간에 따른 기능점수 분포



〈그림 9〉 잔차의 히스토그램



〈그림 10〉 잔차의 정규P-P도표

〈표 6〉 회귀식의 계수

회귀식 계수		계수값	유의확률
a		1.237	0.000
b		0.827	0.000
c	PC일 경우	0	-
	중대형일 경우	0.269	0.000
d	3GL일 경우	0	-
	4GL일 경우	-0.158	0.000
	AP일 경우	-0.249	0.002
e	해외사업의 경우	0	-
	국내사업의 경우	0.219	0.000

에서는 규모의 경제 효과가 발생함을 알 수 있다. 이러한 현상이 발생하는 원인은 크게 2가지로 볼 수 있다. 첫째, ISBSG 자체 분석에서 밝힌 바와 같이 ISBSG 데이터는 우수 사례인 경우가 많다. 따라서 일반적인 소프트웨어개발사업에 비해 생산성이 매우 높을 수 있고 이러한 경우에는 규모의 경제 효과가 발생한다. 유사하게 COCOMO II의 경우에도 SW-CMM(Software Capability Maturity Model) 수준이 5인 경우 오히려 규모의 경제 효과가 나타난다[6]. 두 번째로 규모의 비경제를 초래할 만큼 충분히 큰 기능점수 규모를 갖는 사업이 분석데이터에 적게 포함되었을 수 있다. 전체 분석 데이터 347개 중 2,000FP 이상의 규모를 갖는 사업

은 33개로 9.5%에 지나지 않는다. 소프트웨어개발 사업도 다른 사업과 유사하게 규모가 증가함에 따라 규모의 경제효과를 보이다가 어느 규모 이상으로 커지면 규모의 비경제 효과가 나타나게 되는데, 본 연구의 분석에 사용된 데이터는 규모의 비경제 효과가 나타나기 이전의 기능점수 규모를 갖는 사업이 많았기 때문에 b 가 1보다 작게 나왔다고 볼 수 있다.

다음으로 개발플랫폼이 소요시간에 미치는 영향을 살펴보자. <표 6>에서 개발플랫폼이 PC인 경우에는 회귀식의 계수가 0이고 중대형일 경우 회귀식의 계수가 0.269이다. log를 취하기 이전의 소요노력으로 환원하면 중대형일 경우에는 PC인 경우보다 $10^{0.269} \approx 1.857$ 배 더 많은 소요노력을 필요로 한다. 다음으로 개발언어에 따른 소요노력 차이를 보자. <표 6>에서 개발언어가 3GL, 4GL, AP일 경우에 회귀식의 계수는 각각 -0.158, -0.249이다. 즉 개발언어가 4GL, AP인 경우에는 3GL일 때보다 소요노력이 각각 $10^{-0.158} \approx 1.439$, $10^{-0.249} \approx 1.774$ 배 더 적게 필요함을 알 수 있다.

마지막으로 ISBSG 데이터와 국내 비용자료 사이의 생산성 차이를 살펴보자. <표 6>에서 ISBSG 데이터인 경우에는 회귀식의 계수가 0이고 국내 비용자료인 경우에는 회귀식의 계수가 0.219이다. 이는 국내사업이 ISBSG에 수집된 사업에 비해

$10^{0.219} \approx 1.656$ 배 더 많은 소요노력이 필요함을 의미한다. 회귀모형에서 이렇게 국내사업과 ISBSG 사업 간에 큰 생산성 차이가 발생한 이유는 국내 소프트웨어산업과 선진국의 소프트웨어산업 간의 생산성 차이와 앞서 언급한 바와 같이 ISBSG 데이터가 모범 사업 위주로 데이터가 수집되었다는 사실이 복합적으로 작용했기 때문으로 판단된다.

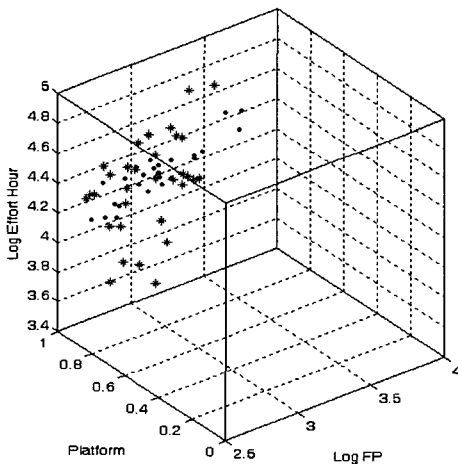
5.3 SVR 적용 결과

식 (5.1)에 쓰인 독립변수를 모두 사용하고 소요노력을 종속변수로 하여 SVR로 학습을 수행하였다. SVR 학습프로그램으로는 LIBSVM[7]을 사용하였으며, 커널함수로는 RBF를 사용하였고 RBF 커널의 파라미터 σ 는 $\sigma^2 = 10$ 으로 설정하였다. 또한, $C=2$ 로 두었고 ϵ -무감도 손실함수(ϵ -insensitive loss function) $L_\epsilon(x, y, f)$ 의 ϵ 값은 $\epsilon = 0.1$ 로 설정하였다. 2장에서 살펴본 바와 같이 (RP)와 (RD)의 최적해를 (α^-, α^+, b) 라 하고 RBF 커널을 사용하는 경우 SVR에 의한 예측치 \hat{y} 는

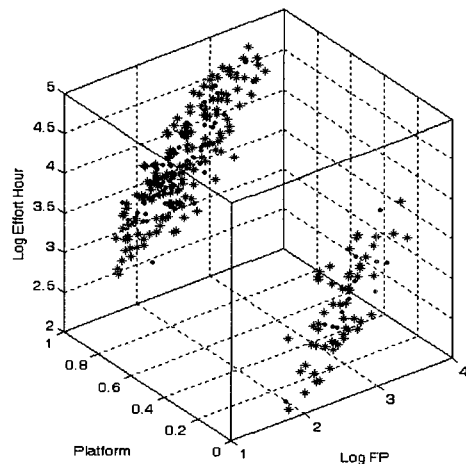
$$\begin{aligned} \hat{y} = f_\phi(x) &= \sum_i (\alpha_i^- - \alpha_i^+) y_i (\phi(x^i))^T \phi(x) + b \\ &= \sum_i (\alpha_i^- - \alpha_i^+) y_i e^{-\frac{\|x^i - x\|^2}{\sigma^2}} + b \end{aligned} \quad (5.3)$$

와 같이 구해진다. 최적조건에 의해 $\alpha_i^- \cdot \alpha_i^+ = 0$ 이므로 α_i^- 또는 α_i^+ 가 0이 아닌 학습데이터 x^i 가 지지벡터가 되는데, 347개의 학습데이터 중 249개가 지지벡터로 선정되었다. 지지벡터로 포함된 데이터와 그렇지 않은 데이터를 소요시간과 기능점수 규모 그래프에 표현하면 <그림 11>과 같다. ‘*’로 표시된 점이 지지벡터이고 ‘•’로 표시된 점은 비지지벡터(non-support vector)를 나타낸다. 회귀모형과 유사하게 중앙에 비지지벡터가 위치하고 중앙에서 멀리 떨어진 자료가 지지벡터가 됨을 볼 수 있다.

SVR을 통한 예측결과와 식 (5.2)의 선형회귀식을 통한 예측결과를 비교해보면 <표 7>와 같이 정리할 수 있다. <표 7>는 크게 학습과정에서의 오차와 검증(validation)과정의 오차에 관한 결과로 구분된다. 회귀모형의 학습과정에서의 오차란 모든 학습데이터를 대상으로 식 (5.2)의 회귀모형계수를 구하고 구해진 회귀식에 의한 예측결과와 실제 소요노력 간의 차이를 비교한 것이다. 마찬가지로 SVR의 학습과정에서의 오차는 모든 학습데이터를 대상으로 (RP), (RD)를 풀어 식 (5.3)의 예측식을 구하고 구해진 예측식에 의한 예측결과와 실제 소요노력 간의 차이를 비교한 것이다. <표 7>의 비



(a) 국내비용자료



(b) ISBSG 비용자료

<그림 11> SVR 학습결과에서 지지벡터

교항목을 살펴보면, i 번째 학습데이터의 소요노력에 로그를 취한 값을 y_i 라 하고 y_i 에 대한 예측치를 \hat{y}_i 라 하자. 학습데이터의 개수를 l 이라 하면, 예측오차 범위는 다음과 같이 구해진다.

$$\min_i \frac{\hat{y}_i - y_i}{y_i} \times 100 \sim \max_i \frac{\hat{y}_i - y_i}{y_i} \times 100$$

또한 MAPE(Mean Absolute Percent Error)는

$$MAPE = \frac{1}{l} \sum_{i=1}^l \left| \frac{\hat{y}_i - y_i}{y_i} \right| \times 100$$

이고, MSE(Mean Squared Error)는

$$MSE = \frac{1}{l} \sum_{i=1}^l (\hat{y}_i - y_i)^2$$

이다. Pred(0.05), Pred(0.10), Pred(0.20)는 예측오차가 각각 5%, 10%, 20% 이내인 사업의 개수를 말한다.

다음으로 검증과정의 오차는 LOO(Leave-One-Out) 예측치와 실제 소요시간 간의 차이를 나타낸 것이다. 학습데이터 (x^i, y_i) 의 LOO 오차는 (x^i, y_i) 를 제외한 나머지 학습데이터를 이용하여 예측모형을 구한 다음 예측모형에 의한 예측값 $\hat{y}_{loo,i}$ 과 y_i 간의 차이를 말한다. LOO 오차는 예측모형을 실제 적용할 때 발생하는 오차에 대한 불편추정치로서 [15], 예측모형의 적합도를 검증해보거나 예측모형의 파라미터를 설정하는데 자주 사용된다. <표 9>에서는 모든 학습데이터에 대해 식 (5.2)의 회귀모형과 SVR의 LOO 오차를 각각 구하여 LOO 오차의 범위와 MAPE, MSE, Pred(·)를 나타내었다.

우선 학습과정에서 발생하는 오차를 비교해 보면 R^2 는 SVR이 선형회귀모형보다 훨씬 높는데 이는 SVR 모형이 선형회귀모형보다 설명력이 높음을 의미한다. 예측오차의 범위 면에서는 SVR이 약간 우수한 것으로 보이나, MAPE, MSE를 살펴보면 SVR이 선형회귀모형에 비해 훨씬 좋은 결과를 나타내고 있다. PRED(·)면에서도 예측오차가 5%, 10% 이내인 사업 개수가 SVR에서는 선형회

귀모형보다 훨씬 많음을 알 수 있다. 검증과정의 오차를 비교해보더라도 학습과정과 유사하게 모든 비교항목에서 SVR이 선형회귀모형보다 예측 정확성이 높은 것을 볼 수 있다. 따라서 소프트웨어 개발비용을 예측하는 데 있어 SVR 모형이 유용함을 확인할 수 있다.

<표 7> 선형회귀모형과 SVR의 예측 정확성 비교

비교항목		선형회귀모형	SVR	
학습	R^2	0.674	0.809	
	예측오차 범위	-37.3%~ +21.3%	-22.8%~ +30.5%	
	MAPE	7.56%	5.30%	
	MSE	0.319	0.058	
	Pred(·)	Pred(0.05)	140	219
		Pred(0.10)	251	302
		Pred(0.20)	336	337
검증	LOO 오차 범위	-21.6%~ +37.7%	-23.2%~ +30.7%	
	MAPE	7.69%	6.49%	
	MSE	0.103	0.082	
	Pred(·)	Pred(0.05)	140	173
		Pred(0.10)	250	281
		Pred(0.20)	331	331

6. 결론 및 추후과제

본 연구는 해외의 소프트웨어 개발 사업 비용자료와 국내 비용자료를 종합적으로 분석하여 소프트웨어 개발비용에 영향을 미치는 요소를 정의하고 그에 따라 소프트웨어 개발비용을 예측하는 모형을 제시하였다. 본 연구에 새롭게 사용한 SVR 모형은 기존의 선형회귀모형을 이용한 예측 방법보다 훨씬 정확한 예측 결과를 보였다.

본 연구는 데이터 획득의 한계로 인해 소프트웨어 개발비용에 영향을 미치는 다양한 요소 중에 일부밖에 고려하지 못하였는데, 추후 다양한 보정요소를 추가적으로 포함하는 예측 모형에 관한 연구가 필요할 것으로 보인다. 또한 소프트웨어 개발비

용 예측에 있어서 선형회귀모형뿐만 아니라 신경망 기법 등 다른 기계학습기법들과 SVR 방법과의 성능비교도 필요할 것이다. 마지막으로 본 연구에서 제시한 예측모형을 국내 소프트웨어 개발 사업에 실제 적용하고 그 결과를 분석하는 것도 흥미로운 과제가 될 것이다.

참고 문헌

- [1] 김우제, 박찬규, 신수정, "AHP를 이용한 소프트웨어개발비 보정계수 산정", 『IE Interfaces』, 제17권, 제5호(2005), pp.1-11.
- [2] 김현수, "기능점수를 이용한 소프트웨어 규모 및 비용산정 방안에 관한 연구", 『경영과학』, 제14권, 제1호(1997), pp.131-149.
- [3] 박찬규, 구자환, 김성희, 신수정, 송병선, "공공부문 정보화사업의 소프트웨어 개발비용 예측에 관한 연구", 『경영과학』, 제19권, 제2호(2002), pp.191-204.
- [4] 박찬규, 신수정, 이현욱, "국내 소프트웨어개발사업에 적합한 기능점수규모 예측방법에 관한 연구", 『경영과학』, 제20권, 제2호(2003), pp.179-196.
- [5] 이양규, "기능점수모형을 이용한 소프트웨어 개발비용 산정", 『경영연구』, 제6권(1997), pp. 241-261.
- [6] Boehm, B.-W. et al., *Software cost estimation with COCOMO II*, Prentice Hall PTR, 2000.
- [7] Chang, C.-C. and C.-J. Lin, *LIBSVM : a library for support vector machines*, <http://www.csie.ntu.edu.tw/~cjlin/libsvm>, 2001.
- [8] Cristianini, N. and J. Shawe-Taylor, *An introduction to support vector machines and other kernel-based learning methods*, Cambridge University Press, 2000.
- [9] Finnie, G.R., G.E. Witting, and J-M. Desharnais, "A comparison of software effort estimation techniques : Using function points with neural networks, case-based reasoning and regression models," *Journal of Systems Software*, Vol.39(1997), pp.281-289
- [10] Pfleeger, S.L., *Software Engineering Theory and Practice 2nd ed.*, Prentice Hall, 2001.
- [11] IFPUG, *Function Point Counting Practices Manual(Release 4.2)*, International Function Point Users Group, 2000.
- [12] ISBSG, *Practical project estimation : A tool kit for estimating software development effort and duration*, International Software Benchmarking Standards Group, 2001.
- [13] Jones, T.C., *Estimating Software Costs*, McGraw-Hill, 1998.
- [14] Kitchenham, B.A. and S.L. Pfleeger, "Software quality : The elusive target," *IEEE Software*, Vol.13(1996), pp.12-21.
- [15] Luntz, A. and V. Brailovsky, "On estimation of characters obtained in statistical procedure of recognition," *Technicheskaya Kibernetica*, Vol.3(1969).
- [16] Porter, A. and R. Selby, "Empirically-guided software development using metric-based classification trees," *IEEE Software*, Vol.7(1990), pp.46-54.
- [17] Price Systems, <http://www.pricystems.com>.
- [18] Schoelkopf, B. and A.J. Smola, *Learning with Kernels : Support Vector Machines, Regularization, Optimization, and Beyond*, MIT Press, 2002.
- [19] Shepperd, M., "Effort and size estimation : An appraisal," *Software Reliability and Metrics Club Newsletter*, January : 6-8 (1997), London, Centre for Software Reliability.

- [20] Sommerville, I., *Software Engineering 6th ed.*, Addison-Wesley, 2001.
- [21] Software Productivity Research, <http://www.spr.com>.
- [22] Srinivasan, K. and D. Fisher, "Machine learning approaches to estimating development effort," *IEEE Transactions on Software Engineering*, Vol.21(1995), pp.126-137.
- [23] Vapnik, V., *Statistical Learning Theory*, Wiley, 1998.
- [24] Vidger, M.R. and A.W. Kark, *Software cost estimation and control*, Institute for Information Technology, National Research Council Canada, 1994.