

## Determination of Protein and Oil Contents in Soybean Seed by Near Infrared Reflectance Spectroscopy

Myoung Gun Choung\*<sup>†</sup>, In Youl Baek\*, Sung Taeg Kang\*, Won Young Han\*,  
Doo Chull Shin\*, Huhn Pal Moon\* and Kwang Hee Kang\*\*

\*National Yeongnam Agricultural Experiment Station, RDA, Milyang 627-803, Korea

\*\*College of Natural Resources, Yeungnam University, Kyongsan 712-749, Korea

**ABSTRACT:** The applicability of near infrared reflectance spectroscopy (NIRS) was tested to determine the protein and oil contents in ground soybean [*Glycine max* (L.) Merr.] seeds. A total of 189 soybean calibration samples and 103 validation samples were used for NIRS equation development and validation, respectively. In the NIRS equation of protein, the most accurate equation was obtained at 2, 8, 6, 1 (2nd derivative, 8 nm gap, 6 points smoothing and 1 point second smoothing) math treatment condition with SNV-D (Standard Normal Variate and Detrend) scatter correction method and entire spectrum by using MPLS (Modified Partial Least Squares) regression. In the case of oil, the best equation was obtained at 1, 4, 4, 1 condition with SNV-D scatter correction method and near infrared (1100 ~2500 nm) region by using MPLS regression. Validation of these NIRS equations showed very low bias (protein : -0.016%, oil : -0.011%) and standard error of prediction (SEP, protein : 0.437%, oil : 0.377%) and very high coefficient of determination ( $R^2$ , protein : 0.985, oil : 0.965). Therefore, these NIRS equation seems reliable for determining the protein and oil content, and NIRS method could be used as a mass screening method of soybean seed.

**Keywords :** soybean [*Glycine max* (L.) Merr.], NIRS, protein, oil content

Soybean [*Glycine max* (L.) Merr.] is an increasingly important human food source and animal feed, mainly due to its high protein and oil contents (Pazdernik *et al.*, 1997). Unripe seeds are eaten as vegetable and dried seeds are eaten whole, split or sprouted. Processed soybeans give soy milk, a valuable protein supplement in infant formula which also provides curds and cheese. Soy source made from the mature fermented beans and soybean used an ingredient in other sauces. In Asia, the highly nutritious sprouts are used a material of food-stuff (James, 1983).

Soybeans seeds contain more protein than any other cultivated commercial crop. Approximately 40% of the dry

weight of the soybean seed is storage protein, and 20% is oil. Soybean protein used for manufacture of food and in manufacture of synthetic fiber, adhesives, textile sizing, waterproofing, fire-fighting foam and many other uses. In the case of soybean oil, it used as salad oil and for manufacture of margarine and shortening. Also, soybean oil used for industrial purpose such as paints, linoleum, oilcloth, printing inks, soap, insecticides and disinfectants (James, 1983).

Recently, more emphasis has been placed on breeding for improving the soybean protein and oil content, but the lack of a fast, efficient mass screening method to determine protein and oil content has slowed breeding progress.

For measuring protein and oil content in soybean seed, the kjeldahl and soxhlet method are widely utilized. However, these methods and related methodologies are relatively complicated and time consuming and involved corrosive chemicals and required elaborate laboratory facilities, which has deterred use in many breeding programs (Williams *et al.*, 1984; Pazdernik *et al.*, 1997). Due to these difficulties, rapid and less hazardous methods, such as the use of near infrared reflectance spectroscopy (NIRS), needed to be developed to estimate protein and oil contents in soybean seeds.

The NIRS is a multi-trait technique that fulfills most of the requirements for rapid, accurate, and cost-effective mass screening for several seed quality traits in many crops (Velasco *et al.*, 1997; Pazdernik *et al.*, 1997; Perez-Vich *et al.*, 1998; Oh *et al.*, 2000). NIRS was first used to measure moisture content in soybean and NIRS has been used to measure moisture, protein, oil and starch contents in many cereals, legumes, forages and other food commodities over the past 20 years (Halgerson *et al.*, 1995; Hatty *et al.*, 1994; Roy *et al.*, 1993). Already the soybean protein and oil contents have been accurately estimated using NIRS at foreign state (Rinne *et al.*, 1975; Hilliard and Daynard, 1976). However, NIRS applications and studies of soybean protein and oil were insufficient in Korea. Therefore, the objectives of this study were to develop the accurate NIRS equation to estimate protein and oil content in soybean and to provide the mass screening technique for high quality soybean breeding.

<sup>†</sup>Corresponding author: (Phone) +82-55-350-1223 (E-mail) cmg7004@rda.go.kr

<Received April 6, 2001>

## MATERIALS AND METHODS

### Soybean samples

The 300 soybean germplasms were used in this study. The soybeans were grown at the experimental field of National Yeongnam Agricultural Experiment Station, Milyang, Korea in 1998 and 1999.

The soybean seed samples were ground with a ball mill and sieved with a 1.0 mm screen. The ground samples were well-mixed and used for scan of NIRS spectral data and analysis of protein and oil contents by standard methods.

### Measuring protein, oil and moisture content

The protein contents of soybean were determined by auto-kjeldahl system. 0.2 g of ground sample was digested by Buchi B-435 digestion system and Buchi B-412 scrubber with 20 ml of sulfuric acid and 3 g of catalyst ( $\text{CuSO}_4 : \text{K}_2\text{SO}_4 = 1 : 9$ ). Percent nitrogen was calculated by Buchi B-339 auto-kjeldahl system and then converted to percent protein by multiplication by 6.25.

The oil contents were determined by auto-soxhlet method with Buchi B-811 extracted system. 2 g of ground samples was extracted by hexane for 2 hours, preheated for ten minutes, and then dried 2 hours at 105°C. This conditions was confirmed in preconditioning experiment (data not shown).

The moisture contents were analyzed by oven-dry method with 105°C for 2 hours, and then all protein and oil contents were converted to dry matter base.

### Scanning and pretreatment of NIRS spectra

The spectra in the visible-near infrared region were measured on a NIRSystem Model 6500 (Silver Springs, MD) monochromator near infrared reflectance spectrophotometer by using a standard cell cup. The NIRS spectral data were recorded between 400 nm and 2500 nm at 2-nm intervals and stored as the reciprocal logarithm ( $\log 1/R$ ) of the reflected energy. NIRS instrument control as well as all graphics and NIRS specific calculations were all performed with the software package WinISI (version 1.02a) by Infra-soft International (Port Matilda, PA). In WinISI software package, two programs, *Center* and *Select*, were used to screen samples for spectral outliers and to choose samples that represented the 300 soybean samples. The *Center* program defined spectral boundaries that eliminated outliers, defined as having a maximum standardized Mahalanobis distance ( $H$ -distance) of 3.0 from the samples mean, and the

*Select* program eliminated samples with similar spectra (minimum standardized  $H$ -distance of 0.6 from their nearest neighboring samples) (Shenk and Westerhaus, 1991a). The results of pretreatment of NIRS spectra, one calibration set (189 samples) and one validation set (103 samples) were randomly selected from the  $\log 1/R$  spectra of 300 soybean germplasms.

### Calibration

The NIRS calibration equations for protein and oil contents analysis were developed for ground seed soybean calibration sample set using the WinISI program *Calibrate* with the MPLS (Modified Partial Least Squares) regression of 3 different derivative math treatment ( $\log 1/R$ ,  $D^1\log 1/R$  and  $D^2\log 1/R$ ). Three additional regression methods, such as PLS (Partial Least Squares), PCR (Principle Component Regression) and MLR (Multiple Linear Regression) were tested on the calibration sample set. The "SNV-D" (Standard Normal Variate and Detrend) and "None" transformations were implemented for scatter correction (Shenk and Westerhaus, 1991b). And the wavelengths at every 2 nm across the entire visible (408~1092 nm) plus near infrared (1108~2492 nm) spectrum were used for calibration. A trimmed spectrum including only the near infrared range was tested against the entire spectrum. The SEC (Standard Error of Calibration),  $R^2$  (Coefficient of determination), SECV (Standard Error Cross-Validation) and/or 1-VR (One minus the ratio of unexplained variance to total variance) statistics were used to select the best calibration equation (Windham *et al.*, 1989).

### Validation

The protein and oil NIRS equations of ground soybean seeds were monitored with the WinISI program *Monitor*, using the validation set of 103 samples. The SEP (Standard Error of Prediction),  $R^2$ , bias, standard deviation of residual and SEP/Mean (Standard Error of Prediction per Mean) statistics were analyzed to determine the accuracy of prediction (Windham *et al.*, 1989). The 103 validation samples had the standardized  $H$ -distance of 3.0 or less from the mean of the calibration sample set.

## RESULTS AND DISCUSSION

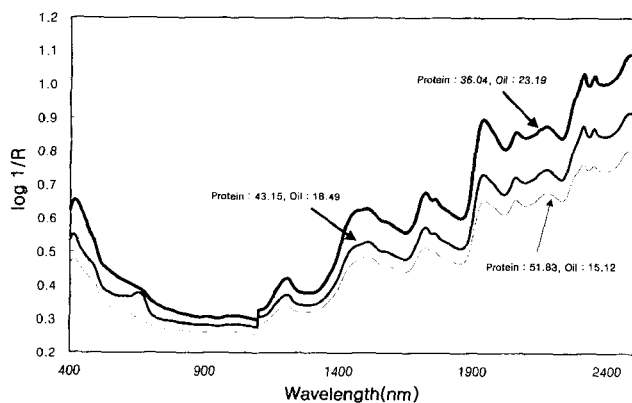
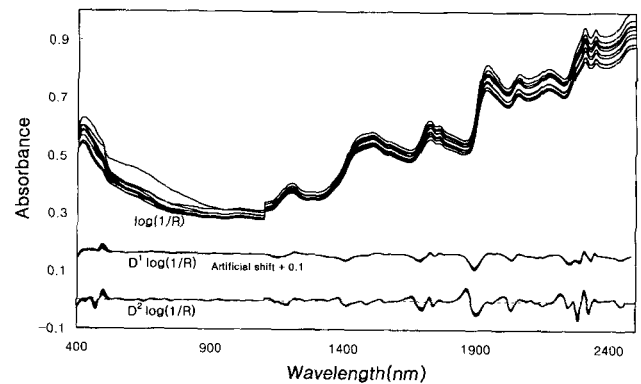
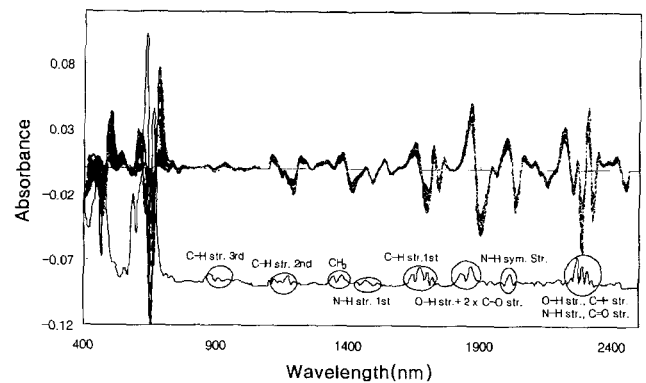
### Protein and oil contents and NIRS spectra

The mean protein and oil contents of the calibration sample sets were 43.23% (range : 36.04% to 51.83%) and

**Table 1.** Laboratory reference value statistics for protein and oil content based on ground soybean seed samples.

Sample set		n	Mean	Range	SD
Calibration	Protein	189	43.23	36.0451.83	3.24
	Oil	189	19.15	14.5724.05	2.05
Validation	Protein	103	42.13	36.1749.83	3.47
	Oil	103	20.06	14.8823.38	2.03

19.15% (range : 14.57% to 24.05%) with a standard deviation of 3.24% and 2.05% as determined by auto-kjeldahl and auto-soxhlet system, respectively (Table 1). There were significant differences among the 189 soybean samples for protein and oil contents based on their standard deviations. These results suggest that sufficient protein and oil variation exist among the samples to develop useful NIRS equations. And the means, ranges and standard deviations of protein and oil contents in validation sample sets were similar to calibration sample set (Table 1). The means and ranges of protein and oil contents were similar to previously reported values (Rinne *et al.*, 1975; Hilliard and Daynard, 1976). The  $\log(1/R)$  spectra of the ground soybean seeds with the high and low contents of protein and oil are shown in Fig. 1. The spectral difference for math treatment effect are shown in Fig. 2. Fig. 3 shows the  $D^2 \log(1/R)$  spectra and mean standard deviation spectrum of calibration samples that are obtained by using the entire wavelength range of 400~2500 nm. Here, except for visible region, several high standard deviation peaks (about 900, 1150, 1370, 1500, 1670, 1870, 2030 and/or 2250 nm) are closely connected with the functional groups (C-H,  $\text{CH}_3$ , O-H, C-O, N-H and C=O), those peaks act to the NIRS calibration of protein and oil (Osborne and Fearn, 1988).

**Fig. 1.** Raw spectra of NIRS with different protein and oil concentration in ground soybean seed samples.**Fig. 2.** Raw, first derivative and second derivative spectra of ground soybean seed samples.**Fig. 3.** Second derivative and mean standard deviation spectra of ground soybean seed calibration sample set.

### Calibration and validation analysis for protein content

The NIRS equations using MPLS method for ground soybean seeds protein are shown in Table 2. The difference of scatter correction method, wavelength and math treatment effect did not highly improve the MPLS model performance, but the optimal equation condition was obtained at 2, 8, 6, 1 (2nd derivative, 8 nm gap, 6 points smoothing and 1 point second smoothing) math treatment condition with SNV-D scatter correction method and entire spectrum (Table 2 and Fig. 4, left panel). In the difference of regression method with same 2, 8, 6, 1 math treatment and SNV-D scatter correction, the equation of MPLS method was showed the lowest SEC and the highest  $R^2$  among other methods (Table 3).

One important criterion for evaluating NIRS equations involves the test of prediction accuracy with unknown samples. Validation sample set allows NIRS equation to be validated for prediction accuracy based on random samples not used in calibration sample set (Pazdernik *et al.*, 1997). Based on the SEP,  $R^2$ , bias, residual of standard deviation

**Table 2.** Comparison on the statistics for protein calibration and validation results with different MPLS conditions.

Math <sup>†</sup> condition	Wavelength (nm)	Scatter <sup>‡</sup>	Calibration <sup>§</sup>				Validation <sup>¶</sup>			
			Term	SEC	R <sup>2</sup>	SEP	R <sup>2</sup>	Bias	R. SD(%)	SEP/M(%)
0. 0. 1. 1	400~2500	None	11	0.614	0.964	0.676	0.963	-0.023	0.68	1.60
	1100~2500	None	11	0.583	0.968	0.652	0.966	0.003	0.66	1.55
	400~2500	SNV-D	11	0.501	0.976	0.567	0.973	0.014	0.57	1.35
	1100~2500	SNV-D	9	0.517	0.975	0.580	0.972	0.021	0.58	1.38
1. 4. 4. 1	400~2500	None	10	0.509	0.975	0.602	0.971	0.001	0.61	1.43
	1100~2500	None	10	0.478	0.978	0.576	0.973	0.019	0.58	1.37
	400~2500	SNV-D	10	0.433	0.982	0.493	0.980	0.009	0.58	1.17
	1100~2500	SNV-D	8	0.436	0.982	0.506	0.979	0.017	0.51	1.20
2. 8. 6. 1	400~2500	None	9	0.492	0.977	0.544	0.976	0.010	0.55	1.29
	1100~2500	None	9	0.404	0.985	0.447	0.983	-0.003	0.45	1.06
	400~2500	SNV-D	10	0.394	0.985	0.437	0.984	-0.016	0.44	1.04
	1100~2500	SNV-D	8	0.517	0.975	0.581	0.973	0.007	0.58	1.38

<sup>†</sup>derivative, gap, smoothing, second smoothing

<sup>‡</sup>scatter correction; SNV-D (standard normal variate and detrend)

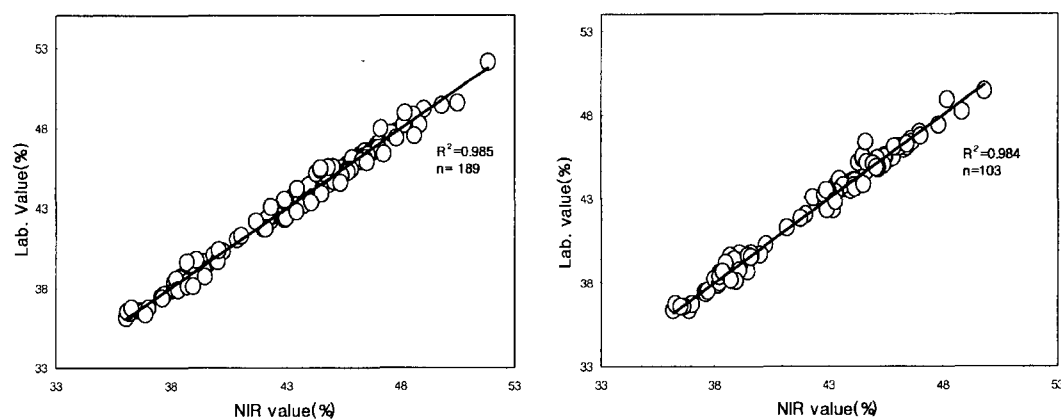
<sup>§,¶</sup>SEC; standard error of calibration, R<sup>2</sup>; coefficient of determination, SEP; standard error of prediction, Bias; difference between reference method and predicted mean, R.SD; residual of standard deviation, SEP/M; standard error of prediction/predicted mean.

**Table 3.** Comparison on the statistics for protein calibration and validation results with different regression methods.

Regression	Calibration <sup>†</sup>				Validation			
	Terms	SEC	R <sup>2</sup>	SEP	R <sup>2</sup>	Bias	R. SD (%)	SEP/M (%)
MPLS <sup>‡</sup>	10	0.394	0.985	0.437	0.984	-0.016	0.44	1.04
PLS	10	0.637	0.961	0.672	0.963	-0.036	0.67	1.59
PCR	10	1.099	0.885	1.080	0.903	0.092	0.92	2.57
MLR	9	0.484	0.978	0.492	0.980	-0.011	0.49	1.17

<sup>†</sup>Calibration condition : (2, 8, 6, 1 math treatment, SNV-D scatter correction, 400~2500 mm)

<sup>‡</sup>MPLS; Modified Partial Least Squares, PLS; Partial Least Squares, PCR; Principle Component Regression, MLR; Multiple Linear Regression

**Fig. 4.** Scatter plots of protein concentration by kjeldahl versus protein concentration by NIRS for the calibration (left) and validation (right) sample set.

and SEP/Mean, the optimal equation using MPLS method equation (2, 8, 6, 1; SNV-D; 400~2500 nm) was accurately predicting the protein contents of validation sample set

(Table 2). The right panel of Fig. 4 demonstrates the accuracy of the ground soybean seeds equation for protein on the basis of the relationship between the actual protein value

calculated from auto-kjeldahl and the predicted protein values from the NIRS. This result indicate that the NIRS analysis can be used as a mass screening method to quickly evaluate a large number of soybean breeding lines for high protein.

### Calibration and validation analysis of oil content

Table 4 and 5 show the NIRS equation statistics for oil

analysis, including SEC,  $R^2$ , SEP, bias, residual of standard deviation and SEP/Mean of the equations obtained from different method conditions. The best equation condition was obtained at 1, 4, 4, 1 (1st derivative, 4 nm gap, 4 points smoothing and 1 point second smoothing) math treatment condition with SNV-D scatter correction method and near infrared (1100~2500 nm) region (Table 4 and Fig. 5, left panel). Based on several prediction statistics, the best NIRS equation of oil analysis using MPLS method equation (1, 4,

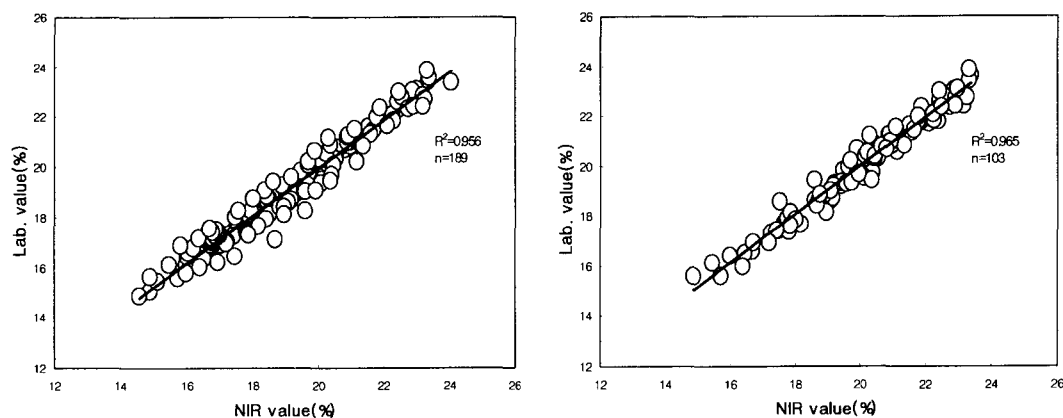
**Table 4.** Comparison on the statistics for oil calibration and validation results with different MPLS conditions.

Math condition	Wavelength (nm)	Scatter	Calibration				Validation			
			Term	SEC	$R^2$	SEP	$R^2$	Bias	R. SD (%)	SEP/M (%)
0. 0. 1. 1	400~2500	None	11	0.730	0.873	0.523	0.933	-0.014	0.53	2.60
	1100~2500	None	11	0.701	0.883	0.526	0.932	-0.027	0.53	2.62
	400~2500	SNV-D	11	0.681	0.890	0.508	0.937	0.006	0.51	2.53
	1100~2500	SNV-D	10	0.637	0.904	0.485	0.942	0.009	0.49	2.42
1. 4. 4. 1	400~2500	None	10	0.538	0.931	0.481	0.944	-0.041	0.48	2.39
	1100~2500	None	10	0.437	0.955	0.410	0.960	-0.034	0.41	2.04
	400~2500	SNV-D	10	0.489	0.943	0.431	0.954	-0.003	0.43	2.15
	1100~2500	SNV-D	9	0.432	0.956	0.377	0.965	-0.011	0.38	1.88
2. 8. 6. 1	400~2500	None	9	0.503	0.940	0.493	0.941	-0.021	0.50	2.45
	1100~2500	None	9	0.497	0.941	0.464	0.947	0.007	0.47	2.31
	400~2500	SNV-D	10	0.497	0.941	0.472	0.946	0.013	0.47	2.35
	1100~2500	SNV-D	9	0.500	0.941	0.501	0.940	-0.029	0.50	2.49

**Table 5.** Comparison on the statistics for oil calibration and validation results with different regression methods.

Regression	Calibration <sup>†</sup>				Validation			
	Terms	SEC	$R^2$	SEP	$R^2$	Bias	R. SD (%)	SEP/M (%)
MPLS	9	0.432	0.956	0.377	0.965	-0.011	0.38	1.88
PLS	11	0.520	0.936	0.453	0.950	-0.008	0.45	2.26
PCR	10	0.788	0.853	0.597	0.913	0.021	0.60	2.98
MLR	4	0.774	0.858	0.605	0.912	0.023	0.61	3.02

<sup>†</sup>Calibration condition : (1, 4, 4, 1 math treatment, SNV-D scatter correction, 1100~2500 nm)



**Fig. 5.** Scatter plots of oil concentration by soxhlet versus oil concentration by NIRS for the calibration (left) and validation (right) sample set.

4, 1; SNV-D; 1100~2500 nm) was well predicting the oil contents of validation sample set, and the SEP value and  $R^2$  of prediction were 0.377% and 0.965, respectively (Table 4 and Fig. 5, right panel). This result indicate that the NIRS analysis can be effective method for measuring soybean oil contents.

In conclusion, the main purpose of creating these NIRS equations was to develop a rapid screening method for protein and oil analysis of soybean seeds. The results of this study show that NIRS can be used as a mass screening technique to quickly evaluate a large number of soybean lines and breeding populations for protein and oil contents. Following the NIRS screening process, the standard analysis method (kjelahl and soxhlet) can be used with greater precision to further identify the very best lines within a smaller, elite group of lines initially selected by NIRS. The main advantage of NIRS is that it reduces the need to analyze the majority of the samples by kjeldahl or soxhlet method.

Future research should aim at developing equations for the non-destructive whole-seed and one-seed and improving the accuracy and sample range. These NIRS equations were developed from seed of soybeans plants grown in only one place Milyang, Kyungnam, Korea, which limits their overall utility. The future usefulness of these NIRS equations resets with the addition of appropriate future samples collected from different locations, which would be used to expand the equations. The addition of a large content range of samples also should improve the predictive accuracy of these NIRS equations.

## REFERENCES

- Halgerson, J. M., C. C. Sheaffer, O. B. Hesterman, T. S. Griffin, M. D. Stern, and G. W. Randall. 1995. Prediction of ruminal protein degradability of forages using near infrared reflectance spectroscopy. *Agron. J.* 87 : 1227-1231.
- Hatty, J. A., W. E. Sabbe, G. D. Basten and A. B. Blakeney. 1994. Nitrogen and starch analysis of cotton leaves using near infrared reflectance spectroscopy (NIRS). *Commun. Soil Sci. Plant Anal.* 25 : 1855-1863.
- Hilliard, J. H. and T. B. Daynard. 1976. Measurement of protein and oil in grains and soybean with reflected near infrared light. *Canadian Institute of Food Science and Technology Journal* 9 : 11-14.
- James, A. Duke. 1983. *Glycine max(L.) Merr.* Handbook of energy crops. unpublished. [http://www.hort.purdue.edu/newcrop/duke\\_energy/glycine\\_max.html](http://www.hort.purdue.edu/newcrop/duke_energy/glycine_max.html)
- Oh, K. W., M. G. Choung, S. B. Pae, C. S. Jung, B. J. Kim, Y. C. Kwan, J. T. Kim and Y. H. Kwack. 2000. Determination of seed lipid and protein contents in perilla and peanut by near-infrared reflectance spectroscopy. *Korean J. Crop Sci.* 45(5) : 339-342.
- Osborne, B. G. and T. Fearn. 1988. Near infrared spectroscopy in food analysis. *Longman Scientific & Technical. John Wiley & Sons, Inc.*
- Pazdernik, D. L., A. S. Killam, and J. H. Orf. 1997. Analysis of amino acid and fatty acid composition in soybean seed, using near infrared reflectance spectroscopy. *Agron. J.* 89 : 679-685.
- Perez-Vich, B., L. Velasco and J. M. Fernandez-Martinez. 1998. Determination of seed oil content and fatty acid composition in sunflower through the analysis of intact seeds, husked seeds, meal and oil by near-infrared reflectance spectroscopy. *J. Amer. Oil Chem. Soc.* 75(5) : 547-555.
- Probst, A. H. and R. W. Judd. 1973. Origin, US history and development and world distribution. In: ed. Calrwell, B. E. Soybeans: Improvement, production and uses. *Agron. monogr.* 16 1st ed. ASA, CSSA and SSSA, Madison, WI.
- Rinne, R. W., S. Gibbons, J. Bradley, R. Sief, and C. A. Brim. 1975. Soybean protein and oil percentage determined by infrared analysis. *Agric. Res. Pub. ARS-NC-26 USDA*: Washington DC.
- Roy, S., R. C. Anantheswaran, J. S. Shenk, M. O. Westerhaus and R. B. Beelman. 1993. Determination of moisture content of mushrooms by vis-NIR spectroscopy. *J. Sci. Food Agric.* 63 : 355-360.
- Shenk, J. S. and M. O. Westerhaus. 1991a. Population definition, sample selection and calibration procedures for near infrared reflectance spectroscopy. *Crop Sci.* 31 : 469-474.
- Shenk, J. S. and M. O. Westerhaus. 1991b. Population structuring of near infrared spectra and modified partial least squares regression. *Crop Sci.* 31 : 1548-1555.
- Velasco, L., J. M. Fernandez-Martinez and A. De Haro. 1997. Determination of the fatty acid composition of the oil in intact seed mustard by near-infrared reflectance spectroscopy. *J. Amer. Oil Chem. Soc.* 74(12) : 1595-1602.
- Williams, P. C., K. R. Preston, K. H. Norris and P. M. Starkey. 1984. Determination of amino acids in wheat and barley by near-infrared reflectance spectroscopy. *J. Food Sci.* 49 : 17-20.
- Windham, W. R., D. R. Mertens and F. E. Barton. 1989. Protocol for NIRS calibration: Sample selection and equation development and validation. In G. C. Marten et al.(ed.) Near infrared reflectance spectroscopy(NIRS): Analysis of forage quality. *Agric. Handb.* 643. USDA-ARS, Washington, DC.