



1920



linear regression : $\hat{y} = b_0 + b_1 x_1 + \dots + b_n x_n$

logistic regression : $\hat{y} = b_0 + b_1 x_1 + \dots + b_n x_n = \log\left(\frac{P}{1-p}\right)$
log odds

$$\hookrightarrow f(x) = \frac{1}{1+e^{-x}}$$

nilai y nya
0 ... 1
(pertamaan klasifikasi)



Masayu Leylia Khodra
(masayu@informatika.org)

KK IF – Teknik Informatika – STEI ITB

$$p = p(\text{kelas : 1} | x)$$

$$\begin{aligned} &= \frac{1}{1 + e^{-(b_0 + b_1 x_1 + \dots + b_n x_n)}} \\ &= \frac{1}{1 + e^{-\text{net}}} \end{aligned}$$

Logistic Regression

What & Why

Pembelajaran Mesin
(Machine Learning)

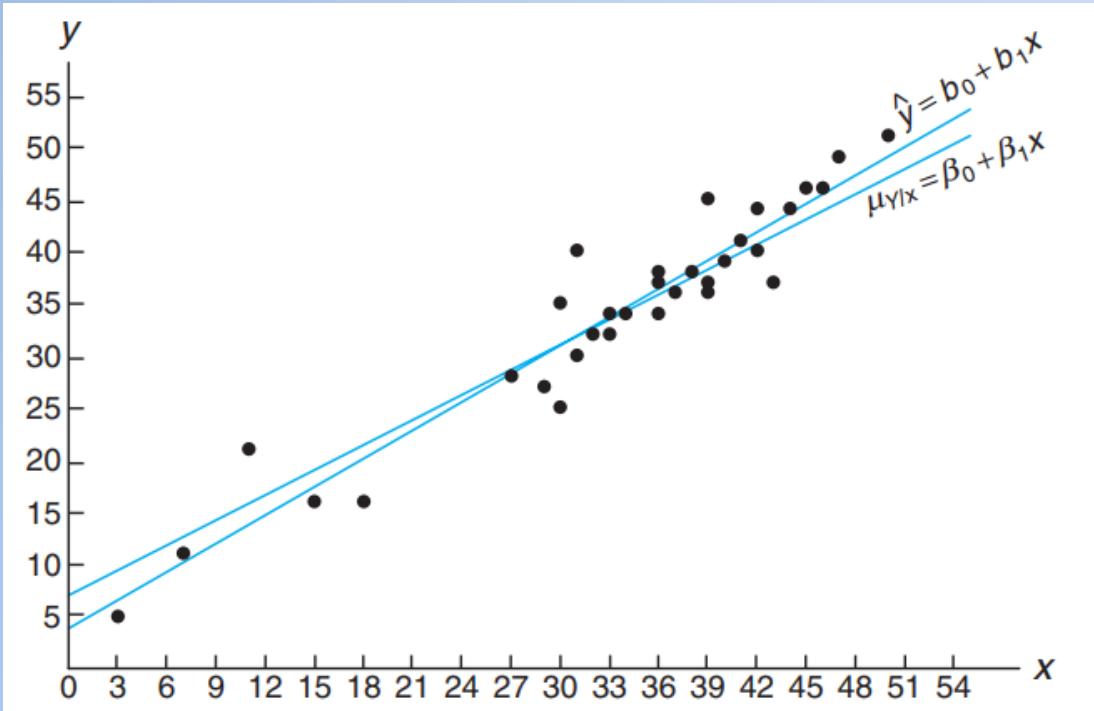


logistic regression.

dataset
label → inductive
learning → $\begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_n \end{bmatrix}$

unseen → inference →
soft class
hard class

REGRESSION ANALYSIS



Finding the best relationship between Y and x: not deterministic, random error

Quantifying the strength of that relationship. Random error with $E(\text{error})=0$ and homoscedasticity. Least Squared Estimation (LSE)

Predicting of the response value y given values of the regressor x.

harus
dipenuhi

Walpole, R. E., Myers, R. H., Myers, S. L., & Ye, K. (2012). Probability and Statistics for engineering and sciences. Pearson Education, 430-435.

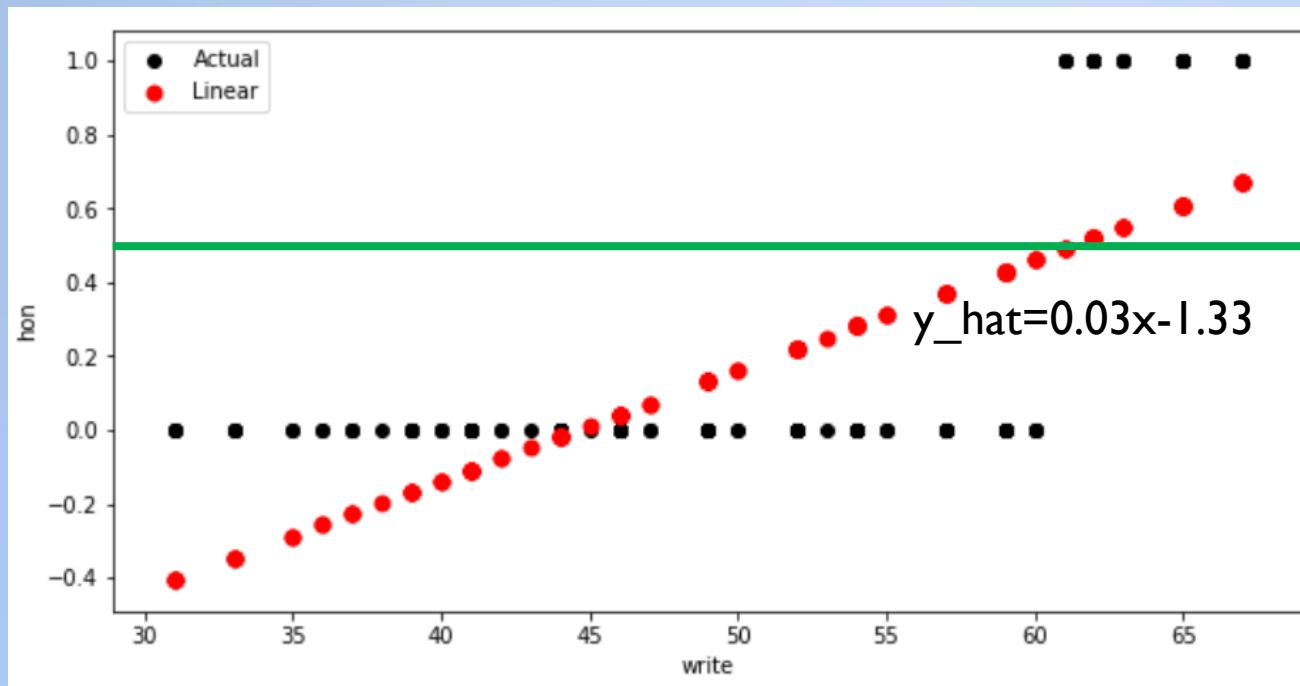




CLASSIFICATION USING LINEAR REGRESSION

LINEAR REGRESSION + THRESHOLD

Example: Dataset has 200 observations (5 attributes), and the target attribute is **hon**, indicating if a student is in an honors class or not.



We are able to predict the value along the Y-axis. If Y is greater than 0.5 (above the green line), predict that the student is in honor class otherwise not in honor class.

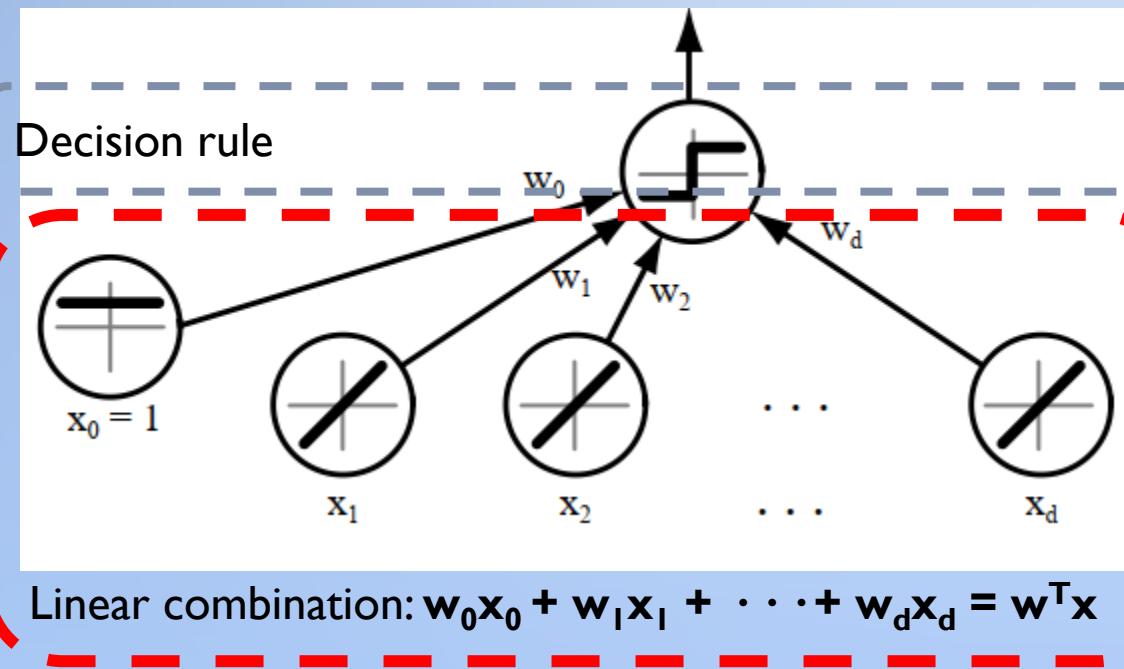
Linear Model or
Linear Discriminant Function
or Linear Classifier





LINEAR DISCRIMINANT FUNCTION

LINEAR COMBINATION + DECISION RULE

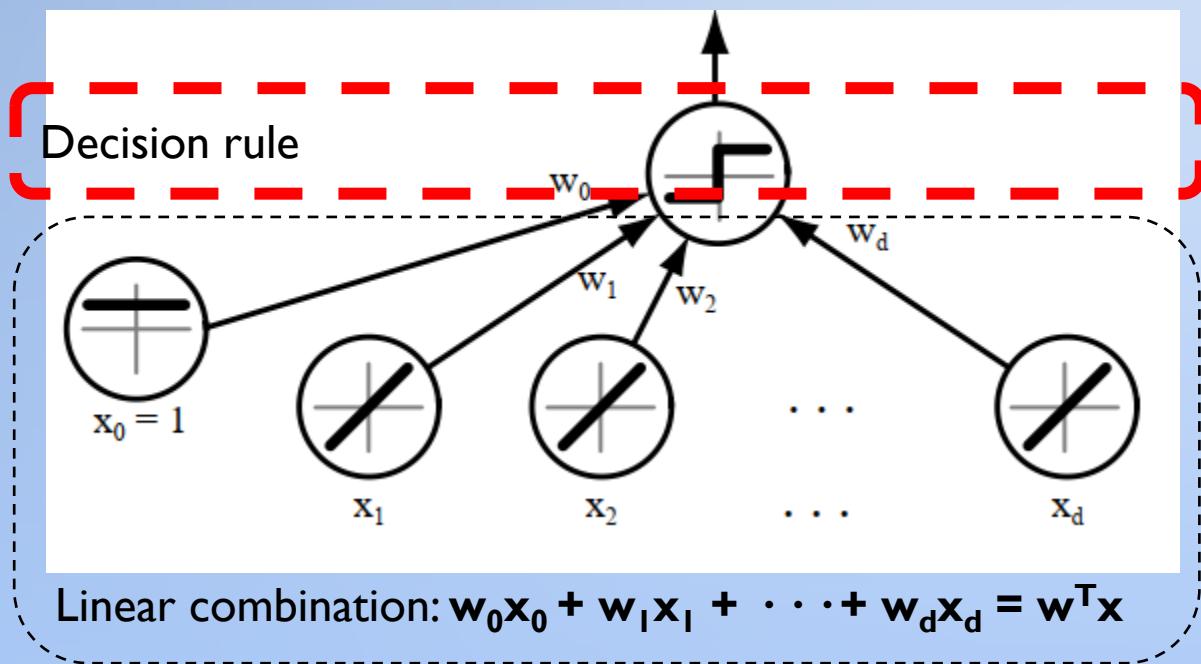


Discriminant function:
a linear combination of the
components of x

$$g(x) = w_0x_0 + w_1x_1 + \dots + w_dx_d = w^T x$$



HYPERPLANE DECISION SURFACE



Discriminant function:

$$g(x) = w^T x$$

Decision rule:
 Decide ω_1 if $g(x) > 0$ and ω_2 if $g(x) < 0$
 $g(x) = 0$: decision surface
 ω_1 and ω_2 are target classes



HYPERPLANE DECISION SURFACE

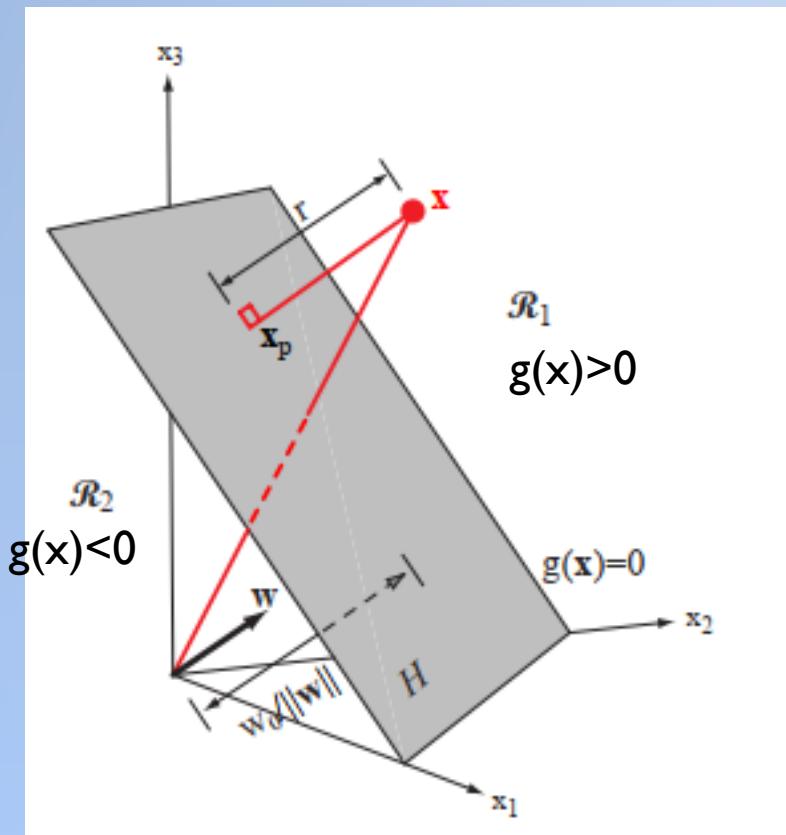


Figure 5.2: The linear decision boundary H , where $g(\mathbf{x}) = \mathbf{w}^t \mathbf{x} + w_0 = 0$, separates the feature space into two half-spaces \mathcal{R}_1 (where $g(\mathbf{x}) > 0$) and \mathcal{R}_2 (where $g(\mathbf{x}) < 0$).

Discriminant function:

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

Decision rule:
 Decide ω_1 if $g(\mathbf{x}) > 0$ and ω_2 if $g(\mathbf{x}) < 0$
 $g(\mathbf{x}) = 0$: decision surface
 ω_1 and ω_2 are target classes



Linear Regression

Least squares

Constant variance

Response variable
is normal

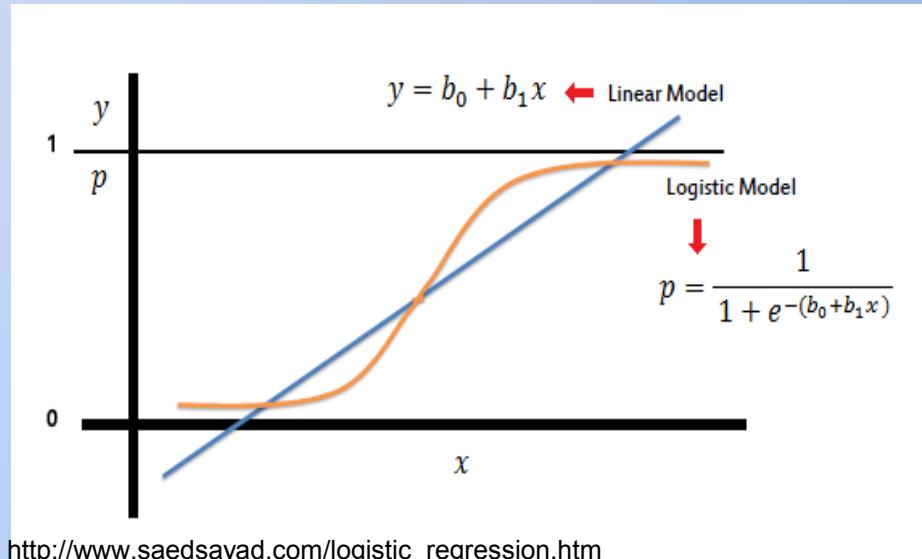
Logistic Regression

Estimate the probability of classes

Maximum likelihood

Non-constant variance

Response variable is binary



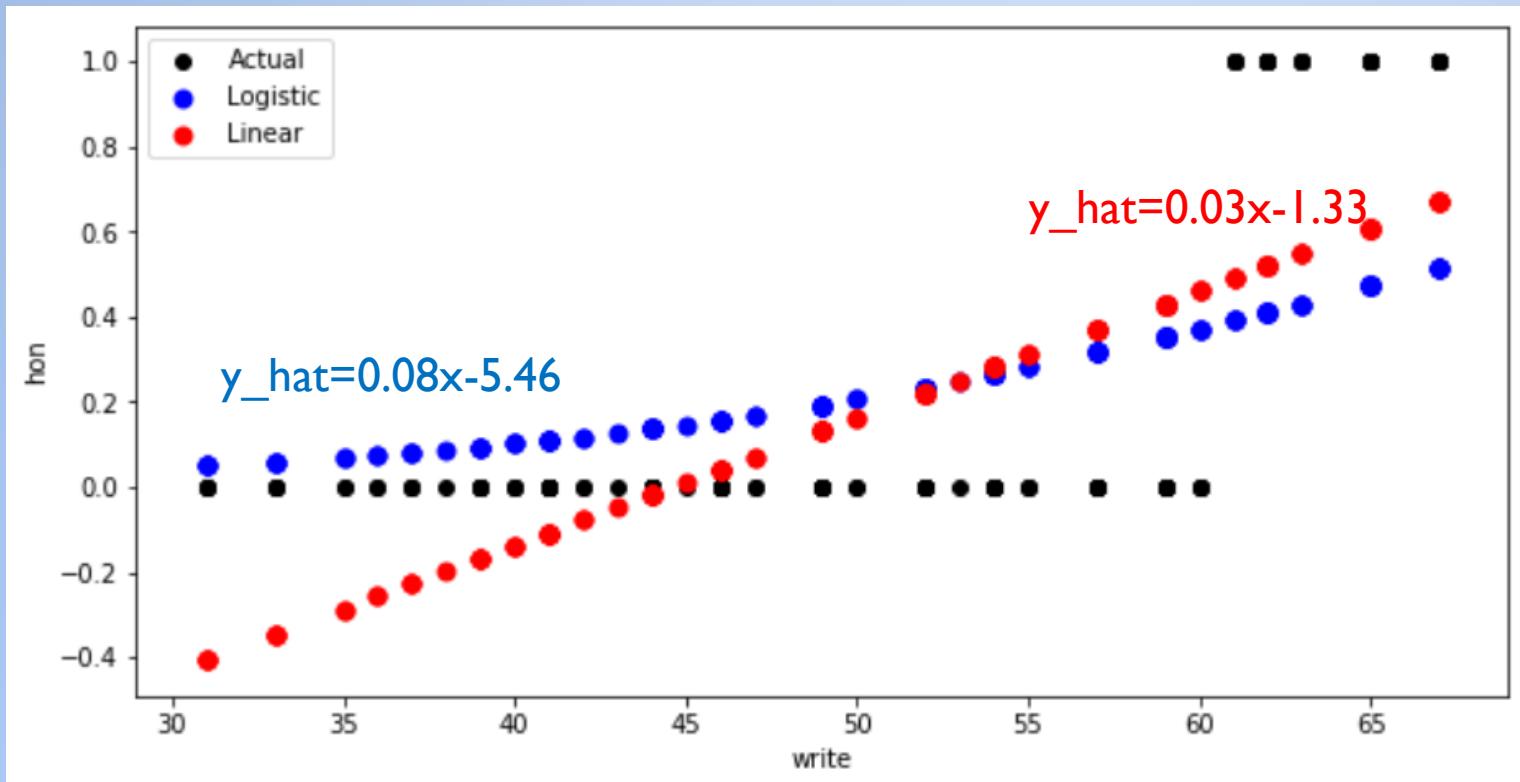
Walpole, R. E., Myers, R. H., Myers, S. L., & Ye, K. (2012). Probability and Statistics for engineering and sciences. Pearson Education, 430-435.





LOGISTIC REGRESSION

$$\hat{y} = \log\left(\frac{p}{1-p}\right) = \mathbf{b}^T \mathbf{x} = b_0 \cdot 1 + b_1 x_1 + \dots + b_d x_d$$



Walpole et al. (2012):
 Odds of success = $p/(1-p)$
 Logistic regression estimates the probability of classes:

$$p = P(y = 1 | \mathbf{x}, \mathbf{b}) = \frac{1}{1 + e^{-\mathbf{b}^T \mathbf{x}}}$$

Coefficients in logistic regression are in terms of the log odds, that is, the coefficient 0.08 implies that a one unit change in "write" results in a 0.08 unit change in the log of the odds.

<https://stats.idre.ucla.edu/stata/faq/how-do-i-interpret-odds-ratios-in-logistic-regression/>



SUMMARY: LOGISTIC REGRESSION



Linear model

$$\hat{y} = \log\left(\frac{p}{1-p}\right) = \mathbf{b}^T \mathbf{x} = b_0 \cdot 1 + b_1 x_1 + \cdots + b_d x_d$$

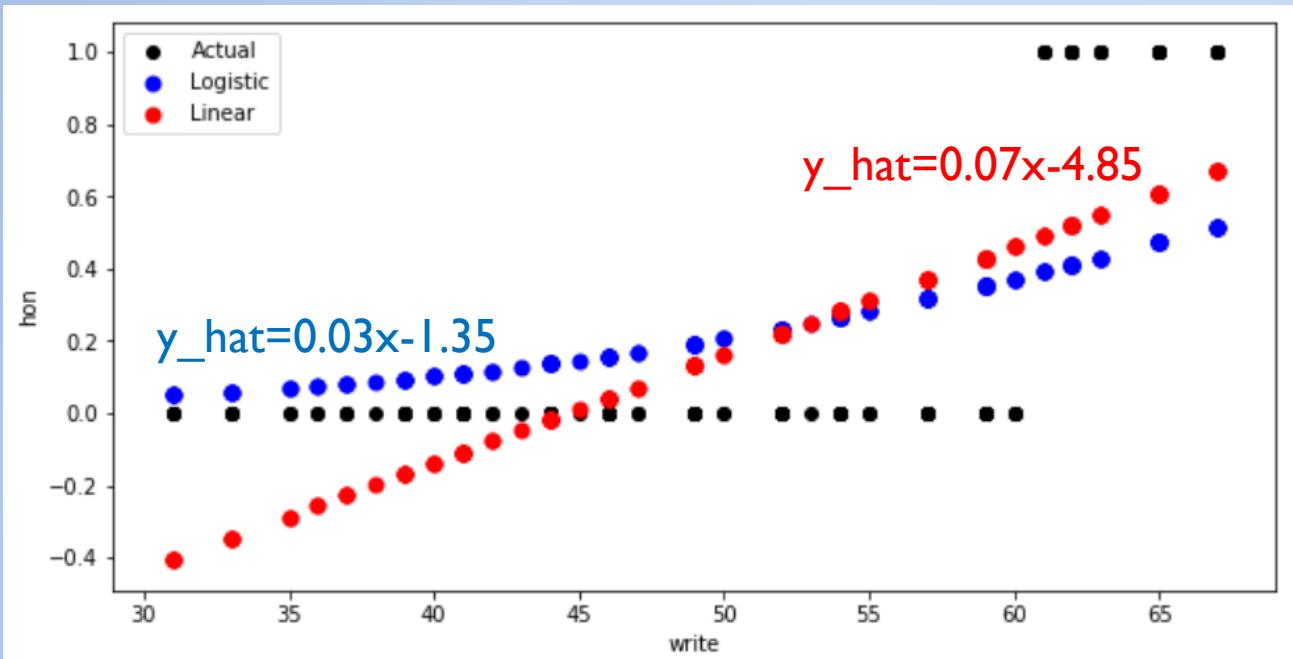
Binary classification

$$p = P(y = 1 | \mathbf{x}, \mathbf{b}) = \frac{1}{1 + e^{-\mathbf{b}^T \mathbf{x}}}$$

Soft classification



EXERCISE



Given these linear regression and logistic regression models, determine whether a student that has write score 65 is in honors class.

Linear regression

$$y = 0.07x - 4.85 \\ x = 65 \rightarrow y = (0.07 \cdot 65) - 4.85 = 0.3 \rightarrow y = 0.3 < 0.5 \rightarrow \text{kelas 0}$$

Logistic regression

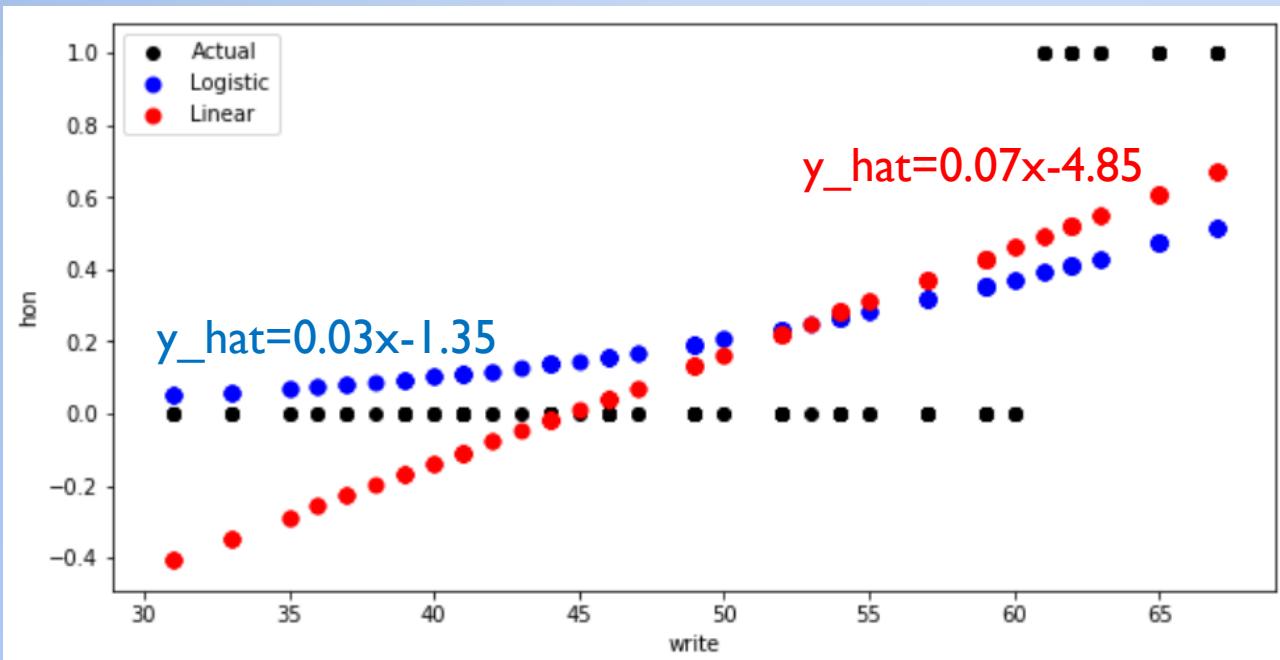
$$\hat{y} = 0.03x - 1.35 \\ x = 65 \rightarrow \hat{y} = (0.03 \cdot 65) - 1.35 = 0.6 \rightarrow p = P(\text{kelas } 1 | x) \\ = \frac{1}{1 + e^{-0.6}} \\ = 0.645 > 0.5 \rightarrow \text{kelas 1.}$$



EXERCISE: SOLUTION



Given these linear regression and logistic regression models, determine whether a student that has write score 65 is in honors class.



Linear regression:

$$y_{\text{hat}}(x) = 0.07x - 4.85$$

$$y_{\text{hat}}(65) = 0.07 * 65 - 4.85 = -0.3$$

Class: 0 ($y_{\text{hat}} = -0.3 < 0.5$)

Logistic regression:

$$y_{\text{hat}}(x) = 0.03x - 1.35$$

$$y_{\text{hat}}(65) = 0.03 * 65 - 1.35 = 0.6$$

$$p = 1 / (1 + e^{-0.6}) = 0.65$$

Class: I ($y_{\text{hat}} = 0.65 > 0.5$)



I3 REFERENCES

- Walpole, R. E., Myers, R. H., Myers, S. L., & Ye, K. (2012). Probability and Statistics for engineering and sciences. Pearson Education, 430-435. Chapter 11 & 12.12, 9.14
- RO Duda, PE Hart, and DG Stork, Pattern Classification, 2nd edition, John Wiley & Sons, 2001. Chapter 5
- Charles Elkan (2014). Maximum Likelihood, Logistic Regression, and Stochastic Gradient Training. <https://cseweb.ucsd.edu/~elkan/250B/logreg.pdf>
- Russell, S., & Norvig, P. (2010). Artificial intelligence: a modern approach. 3rd edition. Chapter 18.6.4.



Masayu Leylia Khodra
(masayu@informatika.org)

KK IF – Teknik Informatika – STEI ITB

Logistic Regression

Stochastic Gradient Ascent

Pembelajaran Mesin
(*Machine Learning*)



15 LOGISTIC REGRESSION



$$\hat{y} = \log\left(\frac{p}{1-p}\right) = \mathbf{b}^T \mathbf{x} = b_0 \cdot 1 + b_1 x_1 + \cdots + b_d x_d$$

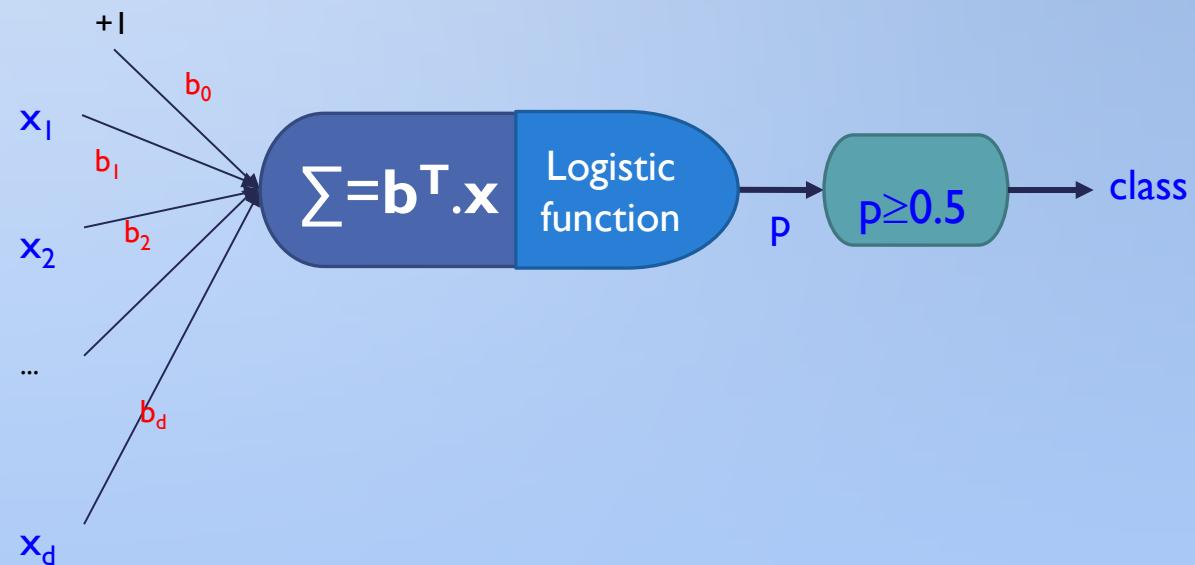
$$p = P(y = 1 | \mathbf{x}, \mathbf{b}) = \frac{1}{1 + e^{-\mathbf{b}^T \mathbf{x}}}$$

Input $\mathbf{x} = (1, x_1, x_2, \dots, x_d)$

Model $\mathbf{b} = (b_0, b_1, b_2, \dots, b_d)$

$$\sum = \mathbf{b}^T \mathbf{x} = b_0 \cdot 1 + b_1 x_1 + \cdots + b_d x_d$$

$$\text{Output} = \sigma(\sum)$$





16 MAXIMUM LIKELIHOOD ESTIMATOR FOR LOGISTIC REGRESSION

ESTIMATOR THAT RESULTS IN A MAXIMUM VALUE FOR ITS JOINT PROBABILITY OR MAXIMIZES THE LIKELIHOOD OF THE SAMPLE

Formal definition

(Elkan, 2014):

Given the training set $\{<\mathbf{x}_1, y_1>..<\mathbf{x}_n, y_n>\}$, learn logistic regression classifier by maximizing the log joint conditional likelihood, that is the sum of log conditional likelihood (LCL) for each training example. x_{ij} is the value of the jth feature of the ith training example.

$$LCL = \sum_{i=1}^n \log L(\theta; y_i | \mathbf{x}_i) = \sum_{i=1; y_i=1}^n \log p_i + \sum_{i=1; y_i=0}^n \log(1 - p_i)$$

$$\rightarrow \frac{\partial LCL}{\partial b_j} = \sum_i (y_i - p_i)x_{ij}$$

Stochastic gradient ascent is optimization method that changes the coefficient values (as random approximation to true derivative) to **increase** the log likelihood based on a randomly chosen example at a time.

Stochastic gradient update of b_j is:

η : learning rate

$$b_j = b_j + \eta(y_i - p_i)x_{ij}$$

<https://cseweb.ucsd.edu/~elkan/250B/logreg.pdf>

learning rate ($0 < \eta < 1$)
makin kecil makin bagus (?)

EDUNEX ITB





STOCHASTIC GRADIENT ASCENT FOR LOGISTIC REGRESSION

INPUT: TRAINING DATA $D=\{<X_1, Y_1>, \dots, <X_N, Y_N>\}$; MAX-ITER T; LEARNING RATE η

Initialize \mathbf{b}

For $t=1, \dots, T$:

For each example $<x_i, y_i>$: #randomly chosen example

$p_i = \text{prediction for } x_i \text{ using the current coefficients } \mathbf{b}$

For each non-zero feature of x_{ij} : $b_j = b_j + \eta(y_i - p_i)x_{ij}$

Return \mathbf{b}

One iteration = one epoch

$D: \{<[52, 41], 0>, <[62, 58], 1>\}; T=1; \eta=0.1$

$\mathbf{b}=[0,0,0] \ #b_0=0; b_1=0; b_2=0 \rightarrow \text{random, biar mudah} = D$

$t=1:$

$$<[62, 58], 1>: \text{dari } b_0=0 + b_1=0 + b_2=0$$

$$p_i = 1 / (1 + e^{-0}) = 0.5$$

$$b_0 = 0 + 0.1(1 - 0.5) * 1 = 0.05$$

$$b_1 = 0 + 0.1(1 - 0.5) * 62 = 3.1$$

$$b_2 = 0 + 0.1(1 - 0.5) * 58 = 2.9$$

$$<[52, 41], 0>:$$

$$p_i = 1 / (1 + e^{-(0.05 + 3.1 * 52 + 2.9 * 41)})$$

$$= 1 / (1 + e^{-280.15}) = 1$$

$$b_0 = 0.05 + 0.1(0 - 1) * 1 = -0.05$$

$$b_1 = 3.1 + 0.1(0 - 1) * 52 = -2.1$$

$$b_2 = 2.9 + 0.1(0 - 1) * 41 = -1.2$$



18 PREDICTION

- $b_0 = -0.05; b_1 = -2.1; b_2 = -1.2$
- $D: \{<[52,41],0>, <[62,58],1>\}$
 - $x_1 = [52, 41]: p_1 = 1/(1+e^{(-0.05-2.1*52-1.2*41)}) = 1/(1+e^{158.45}) = 1.53 \cdot 10^{-69} \rightarrow \text{class}=0 (p_1 < 0.5)$
 - $x_2 = [62, 58]: p_2 = 1/(1+e^{(-0.05-2.1*62-1.2*58)}) = 1/(1+e^{199.85}) = 1.61 \cdot 10^{-87} \rightarrow \text{class}=0 (p_2 < 0.5)$
- Akurasi training = $\frac{1}{2} = 50\%$

x_1	x_2	y	p	Kelas pred	
52	41	0	$1.53 \cdot 10^{-69}$	0	TN
62	58	1	$1.61 \cdot 10^{-87}$	0	FN

confusion matrix

		predicted		Akurasi
		1	0	
label / actual	1	TP : 0	FN : 1	$= \frac{0+1}{0+1+0+1} \cdot \frac{1}{2}$
	0	FP : 0	TN : 1	$= 50\%$

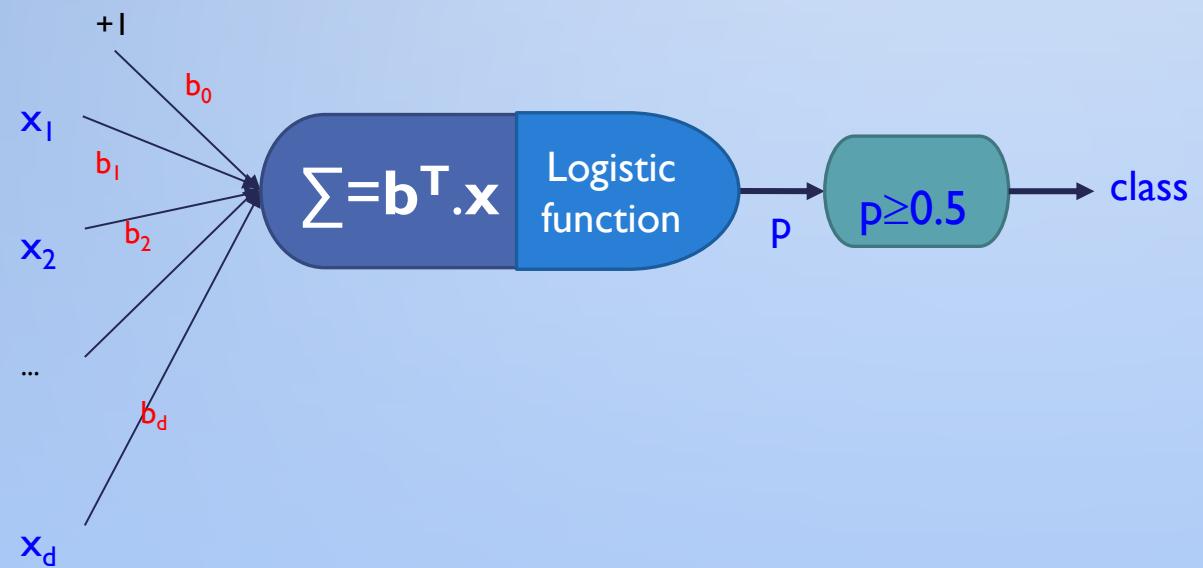
$$\text{Recall}_+ = 0$$

$$\text{Precision}_+ = \frac{1}{2} : 0.5$$

$$\text{Recall}_- = \frac{1}{1} : 1$$



I9 SUMMARY: LOGISTIC REGRESSION



Model: $\mathbf{b} \in \Re^{d+1}$

Maximum Likelihood estimator

Stochastic gradient ascent



20 REFERENCES

- Walpole, R. E., Myers, R. H., Myers, S. L., & Ye, K. (2012). Probability and Statistics for engineering and sciences. Pearson Education, 430-435. Chapter 11 & 12.12, 9.14
- RO Duda, PE Hart, and DG Stork, Pattern Classification, 2nd edition, John Wiley & Sons, 2001. Chapter 5
- Charles Elkan (2014). Maximum Likelihood, Logistic Regression, and Stochastic Gradient Training. <https://cseweb.ucsd.edu/~elkan/250B/logreg.pdf>
- Russell, S., & Norvig, P. (2010). Artificial intelligence: a modern approach. 3rd edition. Chapter 18.6.4.



1920

