

4/4 6/6



M2608.001300

Machine Learning Fundamentals & Applications

[4: Logistic Regression]

Electrical and Computer Engineering
Seoul National University

© 2018 Sungroh Yoon. This material is for educational uses only. Some contents are based on the material provided by the textbook authors and may be copyrighted by them.

Outline

Logistic Regression

Predicting Probability

Cross-Entropy Error Measure

Training via Gradient Descent

Logistic Regression Algorithm

Summary

Appendix: Entropy

Readings

- *Learning from Data* by Abu-Mostafa, Magdon-Ismail, and Lin
 - ▶ Chapter 3: The Linear Model (Sections 3.3 & 3.4)
- *Introduction to Statistical Learning* by James *et al.*
 - ▶ Figures 1 and 2 are from this book (Chapter 4)
- *Machine Learning* by Murphy
 - ▶ Gradient descent examples are from this book (Chapter 8)

Outline

Logistic Regression

Predicting Probability

Cross-Entropy Error Measure

Training via Gradient Descent

Logistic Regression Algorithm

Summary

Appendix: Entropy

Introduction

- core of linear models
 - ▶ 'signal' $s = \underline{w^T X}$: combines input variables linearly
 - ▶ we have seen two models based on this
- 1. linear regression
 - ▶ signal itself = output
 - ▶ for predicting real (unbounded) response
- 2. linear classification
 - ▶ signal is thresholded at zero to produce ± 1 output
 - ▶ for binary decisions

Motivating example

heart attack prediction:

- ▶ based on cholesterol level, blood pressure, age, weight, . . .
- cannot predict a heart attack with any certainty
 - ▶ but can predict how likely it is to occur given these factors
- a more suitable model than binary decision would be:
 - ▶ output y that varies continuously between 0 and 1
 - ▶ the closer y is to 1, the more likely heart attack will occur

Logistic regression: ‘soft’ binary classification

- we introduce the third linear model:
 - ▶ outputs probability of a binary response
e.g. heart attack or not, dead or alive
 - ▶ returns ‘soft labels’ (probability)
- our new model: called “*logistic* regression”
 - ▶ output: **real** (like regression) but bounded (like classification)
sigmoid
- comparison: linear classification vs logistic regression
 - ▶ both deal with a binary event
 - ▶ logistic regression: allowed to be “uncertain”
⇒ intermediate values between 0 and 1 reflect this uncertainty

Example: probability of default

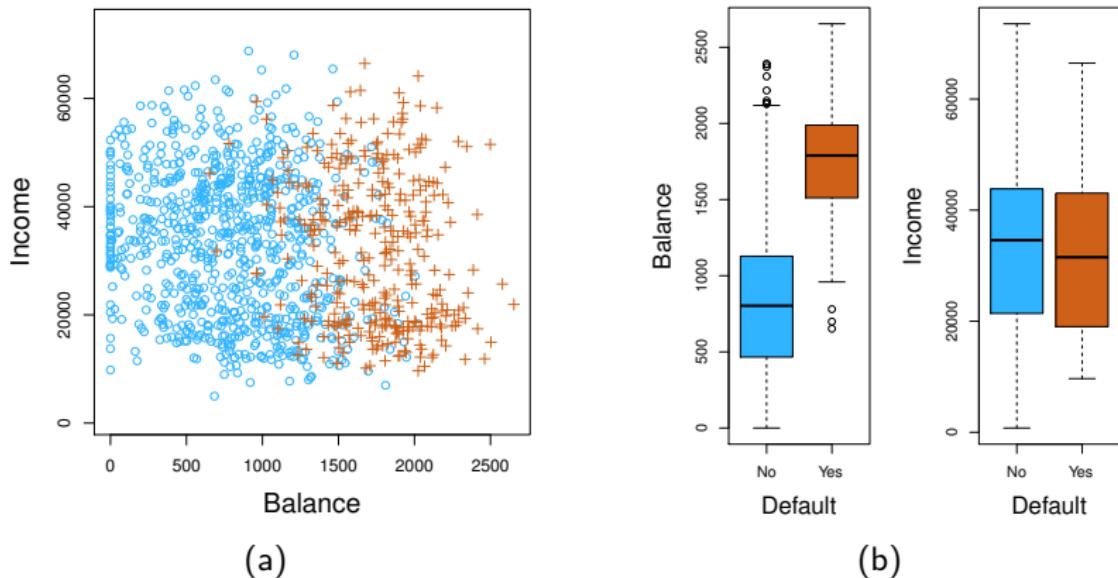
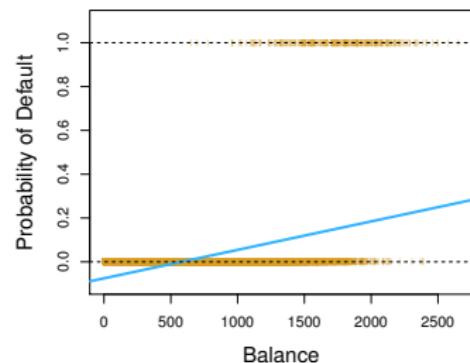
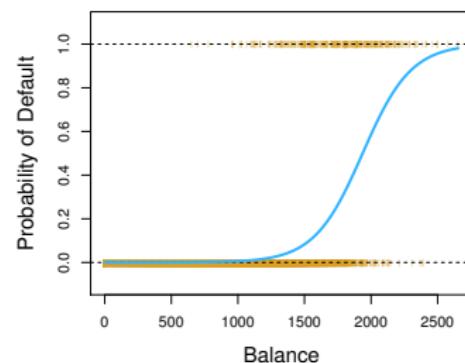


Figure 1: annual incomes and monthly credit card balances: (a) the individuals who defaulted on their credit card payments are shown in orange, and those who did not are shown in blue (b) boxplots of balance/income as a function of default status

- $\mathbb{P}[\text{default} = \text{yes} | \text{balance}]$: probability of default given balance
 - ▶ can be estimated/predicted by regression analysis
 - ▶ logistic regression: more appropriate than linear regression



(a) linear regression



(b) logistic regression

Figure 2: estimated probability of default: (a) some probabilities are negative (b) all probabilities lie between 0 and 1

Logistic regression model

- linear classification: uses a hard threshold on signal $s = \mathbf{w}^T \mathbf{x}$

$$h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x})$$

- linear regression: uses no threshold at all

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

- our new model: needs something between these two
 - ▶ smoothly restricts output to probability range [0,1]
- one choice: the logistic regression model

$$h(\mathbf{x}) = \underline{\theta(\mathbf{w}^T \mathbf{x})}$$

- ▶ θ is so-called logistic function

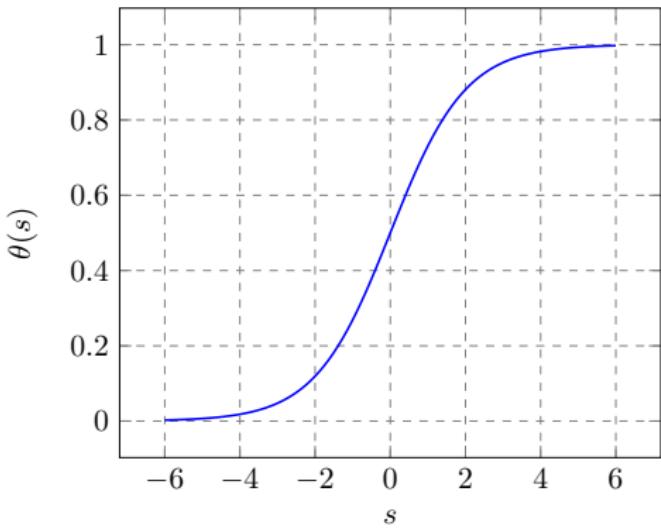
θ: logistic sigmoid

Logistic function θ

- definition

► for $-\infty < s < \infty$:

$$\theta(s) = \frac{e^s}{1 + e^s} = \frac{1}{1 + e^{-s}}$$



- output lies between 0 and 1

► can be interpreted as probability for binary events

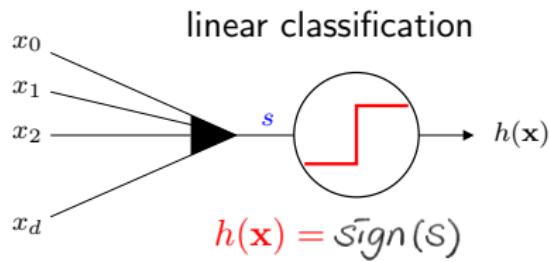
- the above formula of $\theta(s)$
 - ▶ allows us to define an error measure that has analytical and computational advantages (details will follow)
- other names of logistic function θ
 - “soft threshold” in contrast to the hard one in classification
 - “sigmoid” its shape looks like flattened-out ‘s’
- note:

$$1 - \theta(s) = 1 - \frac{1}{1 + e^{-s}} = \frac{e^{-s}}{1 + e^{-s}} = \theta(-s)$$

★ $1 - \theta(s) = \theta(-s)$

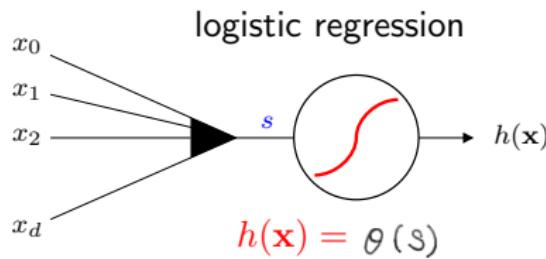
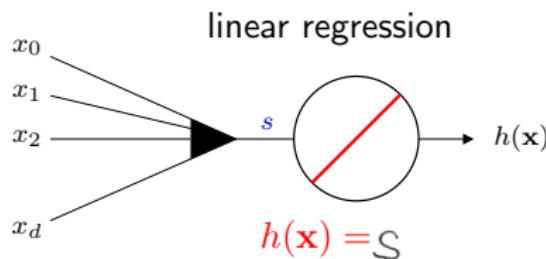
이므로 두어개 징후 모두

Linear models



- based on
“signal” s :

$$s = \sum_{i=0}^d w_i x_i$$



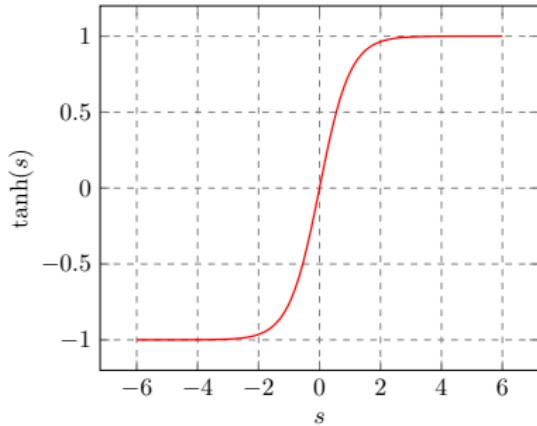
Example: heart attack

- prediction of heart attacks
 - ▶ input \mathbf{x} : cholesterol level, age, weight, etc.
 - ▶ signal $s = \mathbf{w}^T \mathbf{x}$: “risk score”
- linear classification
 - ▶ $h(\mathbf{x})$ returns ± 1 : heart attack (+1) or not (-1) for sure
- linear regression
 - ▶ $h(\mathbf{x})$ returns risk score s itself
- logistic regression
 - ▶ $h(\mathbf{x})$ returns $\theta(s)$: probability of heart attack

Another popular sigmoid function

- hyperbolic tangent:

$$\tanh(s) = \frac{e^s - e^{-s}}{e^s + e^{-s}}$$



- properties
 - ▶ $\tanh(s)$ converges to a hard threshold for large $|s|$
 - ▶ converges to no threshold for small $|s|$

Comparison

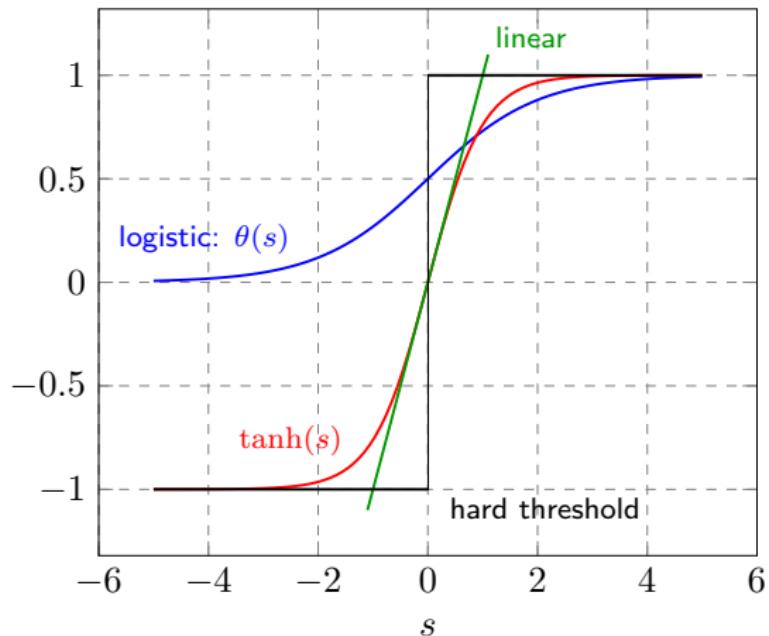


Figure 3: linear, logistic, tanh, and hard threshold

Outline

Logistic Regression

Predicting Probability

Cross-Entropy Error Measure

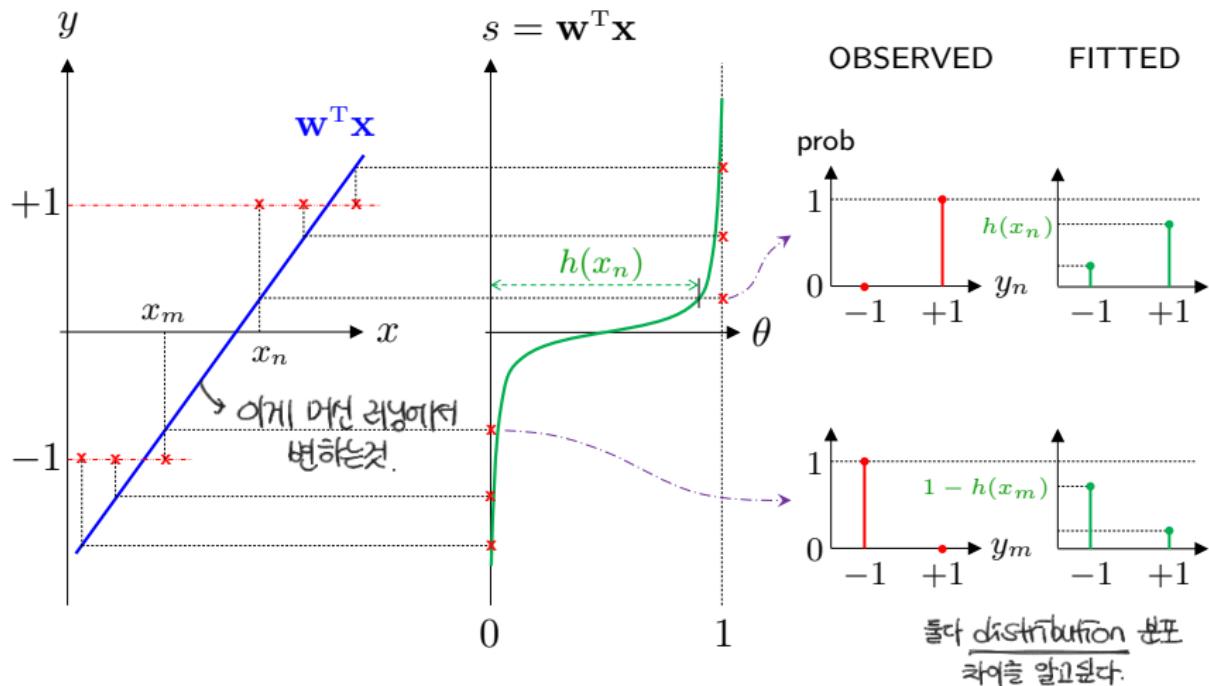
Training via Gradient Descent

Logistic Regression Algorithm

Summary

Appendix: Entropy

Big picture



$$\mathbb{P}[y = +1 | \mathbf{x}] = f(\mathbf{x}) \sim h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x})$$

⇒ 분포 사이의 거리를 재는 것

$$P(y|\mathbf{x}) = h(\mathbf{x})^{[y=+1]} (1-h(\mathbf{x}))^{[y=-1]} = \theta(y\mathbf{w}^T \mathbf{x})$$

Learning target in logistic regression

- learning target:

- ▶ probability of event $y = +1$ given input \mathbf{x}

$$f(\mathbf{x}) = \mathbb{P}[y=+1 | \mathbf{x}]$$

event에 대한 확률

- ▶ e.g. probability of a patient being at risk for heart attack given the characteristics of the patient

- training data do not give us the value of f explicitly

- ▶ rather, give us samples generated by this probability
 - ▶ e.g. patients with/without heart attacks

⇒ learning target: inherently noisy

► target is $P(y|\mathbf{x})$ rather than $y = f(\mathbf{x})$

ideal (noiseless) data

$$(\mathbf{x}_1, f_1 = 0.99 = \mathbb{P}[+1|\mathbf{x}_1])$$

$$(\mathbf{x}_2, f_2 = 0.13 = \mathbb{P}[+1|\mathbf{x}_2])$$

⋮

$$(\mathbf{x}_N, f_N = 0.51 = \mathbb{P}[+1|\mathbf{x}_N])$$

we can learn
from this.

actual (noisy) data

$$(\mathbf{x}_1, y_1 = +1)$$

$$(\mathbf{x}_2, y_2 = -1)$$

⋮

$$(\mathbf{x}_N, y_N = -1)$$

we can't learn
from this.

학습에 따라 distribution 을 learn.

- we thus view data as generated by target distribution $P(y|\mathbf{x})$

$$P(y|\mathbf{x}) = \begin{cases} f(\mathbf{x}) & \text{for } y = +1 \\ 1 - f(\mathbf{x}) & \text{for } y = -1 \end{cases} \quad (1)$$

- to learn from such data, we need a proper error measure
 - ▶ it gauges how close a given hypothesis h is to target f
 - ▶ in terms of noisy ± 1 examples

What makes an h good?

- 'fitting' data \mathcal{D} means finding a good h

► h is good if $\begin{cases} h(\mathbf{x}_n) \approx 1 & \text{whenever } y_n = +1 \\ h(\mathbf{x}_n) \approx 0 & \text{whenever } y_n = -1 \end{cases}$

- a simple error measure that captures this

$$E_{\text{in}}(h) = \frac{1}{N} \sum_{n=1}^N \left(h(\mathbf{x}_n) - \frac{1}{2}(1 + y_n) \right)^2$$

has exponential effect.

► not very convenient (hard to minimize)

- preview: we will use cross-entropy error measure

$$E_{in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln \left(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n} \right)$$

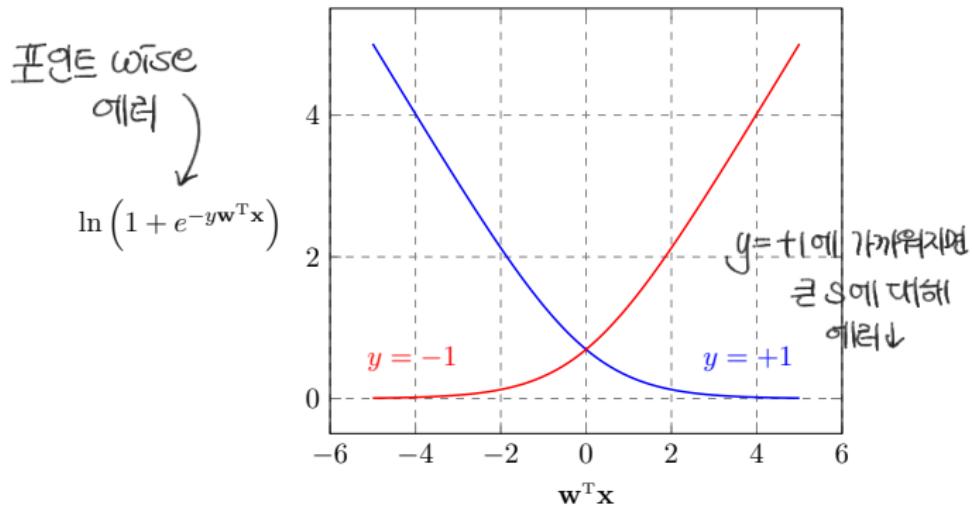
((2))

- ▶ looks complicated and ugly; however (2) is
 - based on an intuitive probabilistic interpretation (next page)
 - has nice property for gradient-based optimization
(although not as easy as linear regression)

(1. 예측의 성질 가져
2. Convex 꼴을 가진다. 가장 높은점 쉽게)



- $y_n = +1$ encourages $\mathbf{w}^T \mathbf{x}_n \gg 0$, so $\theta(\mathbf{w}^T \mathbf{x}_n) \approx 1$, and vice versa



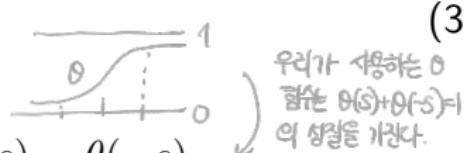
↳ 예외의 성질을 가진다.

Defining error measure

- (1) 확률적: ML. (maximum likelihood)
(2) 정보이론: cross entropy.

- standard error measure in logistic regression: based on "likelihood"
 - ▶ how 'likely' is it to get output y from input \mathbf{x} , if target distribution $P(y|\mathbf{x})$ was indeed captured by $h(\mathbf{x})$?
- based on (1), that likelihood would be

$$P(y|\mathbf{x}) = \begin{cases} \frac{h(\mathbf{x})}{1-h(\mathbf{x})} & \text{for } y = +1 \\ 1 - h(\mathbf{x}) & \text{for } y = -1 \end{cases}$$
$$= \boxed{\theta(y\mathbf{w}^T\mathbf{x})} \quad (3)$$



- recall: $h(\mathbf{x}) = \theta(\mathbf{w}^T\mathbf{x})$ and $1 - \theta(s) = \theta(-s)$
- this simplicity in (3): reason for defining $\theta(s)$ as $e^s/(1 + e^s)$

Criterion for choosing h : maximum likelihood

- assume $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ are independently generated
⇒ probability of getting all y_n 's from corresponding \mathbf{x}_n 's:

$$P(y_1|\mathbf{x}_1)P(y_2|\mathbf{x}_2)\cdots P(y_N|\mathbf{x}_N) = \prod_{n=1}^N P(y_n|\mathbf{x}_n)$$

► this is the '*likelihood*' of training data

- the method of *maximum likelihood*
► selects the hypothesis h that maximizes this probability

- we can equivalently minimize a more convenient quantity

$$-\frac{1}{N} \ln \left(\prod_{n=1}^N P(y_n | \mathbf{x}_n) \right) = \frac{1}{N} \sum_{n=1}^N \ln \frac{1}{P(y_n | \mathbf{x}_n)}$$

$\therefore -\frac{1}{N} \ln(\cdot)$ is a monotonically decreasing function
log likelihood
NLL : negative √ : max(ML) = min(NLL)

- substituting with Eq. (3), we would be minimizing (wrt \mathbf{w}):

$$\frac{1}{N} \sum_{n=1}^N \ln \frac{1}{\theta(y_n \mathbf{w}^T \mathbf{x}_n)} = \frac{1}{N} \sum_{n=1}^N \ln \begin{cases} \frac{1}{h(\mathbf{x}_n)} & \text{for } y_n = +1 \\ \frac{1}{1-h(\mathbf{x}_n)} & \text{for } y_n = -1 \end{cases} \quad (4)$$

$$= \frac{1}{N} \sum_{n=1}^N \left\{ \underbrace{[y_n = +1]}_{\substack{\text{positive의 } \rightarrow \text{학습} \\ \text{경우}}} \ln \frac{1}{h(\mathbf{x}_n)} + \underbrace{[y_n = -1]}_{\substack{\text{negative의 } \rightarrow \text{학습}}} \ln \frac{1}{1-h(\mathbf{x}_n)} \right\} \quad (5)$$

- ▶ more details on (5) will follow

정보이거나 더 나지

E_{in} for logistic regression

- the fact that we are *minimizing* quantity (4)
⇒ allows us to treat it as an '*error measure*'
이거로 생각하자.
- substituting the functional form for $\theta(y_n \mathbf{w}^T \mathbf{x}_n)$ produces
 - in-sample error measure for logistic regression

$$E_{in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln \left(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n} \right) \quad (6)$$

pointwise 예외 같은 문제.

- implied point-wise error

$$e(h(\mathbf{x}_n), y_n) = \ln \left(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n} \right)$$

Optimization involved

$$\max \prod_{n=1}^N P(y_n | \mathbf{x}_n) \quad (7)$$

$$\Leftrightarrow \max \ln \left(\prod_{n=1}^N P(y_n | \mathbf{x}_n) \right) \quad (8)$$

$$\equiv \max \sum_{n=1}^N \ln P(y_n | \mathbf{x}_n) \quad (9)$$

$$\Leftrightarrow \min -\frac{1}{N} \sum_{n=1}^N \ln P(y_n | \mathbf{x}_n) \quad (10)$$

$$\equiv \min \frac{1}{N} \sum_{n=1}^N \ln \frac{1}{P(y_n | \mathbf{x}_n)} \quad \text{설명} \quad (11)$$

$$\equiv \min \frac{1}{N} \sum_{n=1}^N \ln \frac{1}{\theta(y_n \mathbf{w}^T \mathbf{x}_n)} \quad \text{설명} \quad (12)$$

$$\equiv \min \frac{1}{N} \sum_{n=1}^N \ln \left(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n} \right) \quad \text{설명} \quad (13)$$

$$E_{in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln \left(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n} \right)$$

$\therefore \text{Max. likelihood} = \text{Min. crossentropy}$.
► (11) → (12): by logistic modeling

Remarks

- this error measure: small when $y_n \mathbf{w}^T \mathbf{x}_n$ is large and *positive*
 - ▶ which would imply that $\text{sign}(\mathbf{w}^T \mathbf{x}_n) = y_n$
- therefore, as our intuition expect
 - ▶ the error measure encourages \mathbf{w} to ‘classify’ each \mathbf{x}_n correctly
- alternative derivation of $E_{\text{in}}(\mathbf{w})$ is possible
 - ▶ based on the notion of *cross entropy*

Cross entropy

- consider two pmfs $\{p, 1 - p\}$ and $\{q, 1 - q\}$ with binary outcomes
 - ▶ e.g. $\mathbb{P}[\text{pass}] = p$ and $\mathbb{P}[\text{fail}] = 1 - p$
- cross entropy for these two pmfs: defined by

$$p \log \frac{1}{q} + (1 - p) \log \frac{1}{1 - q} \quad (14)$$

- ▶ log denotes binary logarithm (\log_2)
- ▶ see Appendix for more details
- cross entropy measures the ‘error’ for approximating
 - ▶ ‘observed’ pmf $\{p, 1 - p\}$ by ‘fitted’ pmf $\{q, 1 - q\}$
= learned

Alternative derivation based on cross entropy

- recall:

$$P(y|\mathbf{x}) = \begin{cases} h(\mathbf{x}) & \text{for } y = +1 \\ 1 - h(\mathbf{x}) & \text{for } y = -1 \end{cases} \quad (15)$$

$$= \theta(y\mathbf{w}^T \mathbf{x}) \quad (16)$$

- (15) can also be written as

$$P(y|\mathbf{x}) = h(\mathbf{x})^{\llbracket y=+1 \rrbracket} (1 - h(\mathbf{x}))^{\llbracket y=-1 \rrbracket} \quad (17)$$

| 1일 때 1, 아니면 0.

- likelihood of data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ now becomes

$$\prod_{n=1}^N P(y_n | \mathbf{x}_n) = \prod_{n=1}^N h(\mathbf{x}_n)^{[y_n=+1]} (1 - h(\mathbf{x}_n))^{[y_n=-1]} \quad (18)$$

- negative log-likelihood (NLL) is given by

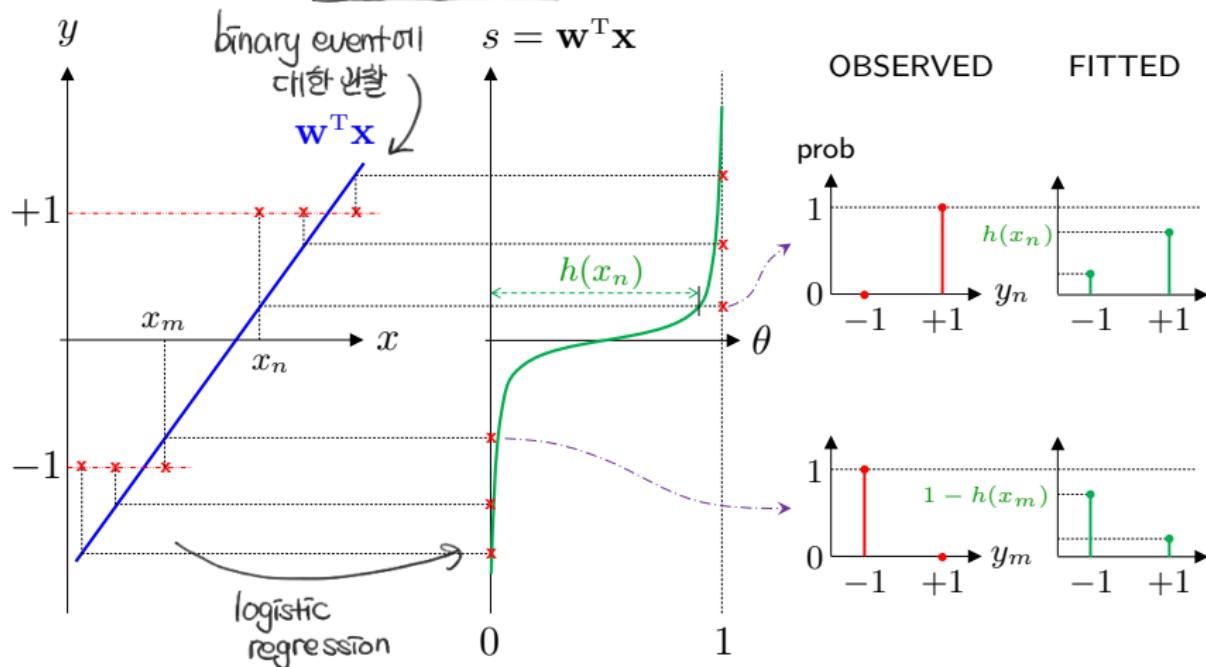
$$\begin{aligned}
 NLL(\mathbf{w}) &\triangleq -\log \left\{ \prod_{n=1}^N P(y_n | \mathbf{x}_n) \right\} \\
 &= -\log \left\{ \prod_{n=1}^N h(\mathbf{x}_n)^{\llbracket y_n = +1 \rrbracket} (1 - h(\mathbf{x}_n))^{\llbracket y_n = -1 \rrbracket} \right\} \\
 &= \sum_{n=1}^N \left\{ \frac{\llbracket y_n = +1 \rrbracket \log \frac{1}{h(\mathbf{x}_n)}}{\text{이제 } y_n \text{가 } +1 \text{인 경우 } h(\mathbf{x}_n) \text{를 } 1 \text{로 정의함}} + \llbracket y_n = -1 \rrbracket \log \frac{1}{1 - h(\mathbf{x}_n)} \right\} \quad (19)
 \end{aligned}$$

- this is also called cross-entropy error function
- related to log loss and KL divergence

- n th term in (19): cross entropy for two pmfs measured on (\mathbf{x}_n, y)
 - ▶ observed: $\{\llbracket y_n = +1 \rrbracket, 1 - \llbracket y_n = +1 \rrbracket\}$
 - ▶ fitted: $\{h(\mathbf{x}_n), 1 - h(\mathbf{x}_n)\}$
 - ▶ in (14): $p = \llbracket y_n = +1 \rrbracket$ and $q = h(\mathbf{x}_n)$
- minimizing $NLL(\mathbf{w})$ is equivalent to
 - ▶ maximizing (18)
 - ▶ minimizing (5)
 - ▶ minimizing (6) for $h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x})$

Big picture

로지스틱 회귀분석



$$\mathbb{P}[y = +1 | \mathbf{x}] = f(\mathbf{x}) \sim h(\mathbf{x}) = \theta(w^T \mathbf{x})$$

$$P(y|\mathbf{x}) = h(\mathbf{x})^{[y=+1]} (1-h(\mathbf{x}))^{[y=-1]} = \theta(yw^T \mathbf{x})$$

Outline

Logistic Regression

Predicting Probability

Cross-Entropy Error Measure

Training via Gradient Descent

Logistic Regression Algorithm

Summary

Appendix: Entropy

Training linear models

- **linear classification** (the perceptron)
 - ▶ minimize E_{in} by solving a combinatorial optimization problem
 - ▶ PLA/pocket algorithm NP hard 문제.
- **linear regression**
 - ▶ minimize E_{in} by setting $\nabla E_{in}(\mathbf{w}) = 0$
 - ▶ analytic pseudo-inverse algorithm
→ closed form solution 가능
- **logistic regression**
 - ▶ minimize E_{in} by setting $\nabla E_{in}(\mathbf{w}) = 0$
 - ▶ iterative optimization (*e.g.* gradient descent)
(경사법을 사용)

Training logistic regression 어떻게 train?

최적화 matrix 보고
이거 어떤 positive semi-definite 어떤 convex 꼴. 그러면

- take an approach similar to linear regression: set

$$\nabla E_{in}(w) = 0$$

- unfortunately, ∇E_{in} for logistic regression:

- not easy to manipulate analytically

$$\nabla E_{in}(w) = -\frac{1}{N} \sum_{n=1}^N \frac{y_n \mathbf{x}_n}{1 + e^{y_n \mathbf{w}^T \mathbf{x}_n}} \quad (20)$$

$$= \frac{1}{N} \sum_{n=1}^N -y_n \mathbf{x}_n \theta(-y_n \mathbf{w}^T \mathbf{x}_n) \quad (21)$$

- an analytic solution is not feasible!

⇒ 딱히 iteratively 문제는 푼다.

Iterative optimization

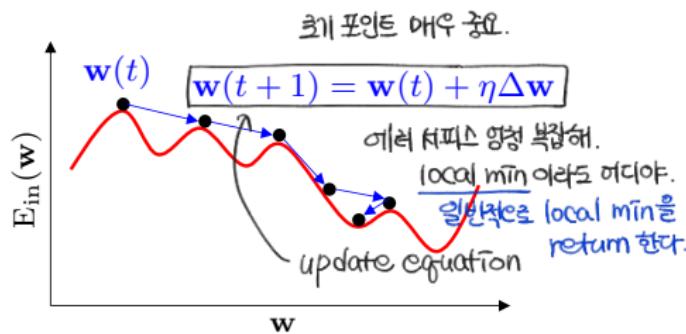
- we will iteratively set the gradient to zero
 - ▶ instead of analytically setting it to zero 이해하기
미분값=0 일때 min. max. saddle 이해하기 힘들다.
 - gradient descent: a very general algorithm
 - ▶ can train many other models with smooth error measures
 - ▶ has particularly nice properties for logistic regression
 - gradient: a fancy word for the derivative of vector functions
 - ▶ direction of greatest *increase* of a function (*i.e.* "steepest")
 - ▶ zero at (local) max/min (\because no single direction of increase)
우리는 error surface 가장 가파른 방향

Gradient descent

- logistic regression에서 종종 사용된다.
- optimization 문제에서 많이 사용된다.

- general technique for minimizing twice-differentiable function
 - ▶ e.g. $E_{in}(w)$ in logistic regression and neural networks
- physical analogy: a ball rolling down a hilly surface
 - ▶ if placed on a hill, it will roll down
 - ▶ it comes to rest at the bottom of a valley

- same basic idea underlies gradient descent
 - ▶ $E_{in}(\mathbf{w})$ is a ‘surface’ in high-dimensional space
 - ▶ in step 0, we start somewhere on this surface, at $\mathbf{w}(0)$
 - ▶ and try to roll down the surface, thereby decreasing E_{in}
 빙향이 거치는 정해진 훈련
 (RH Iteration에서 결정사항)



- caution: local optimum
 - ▶ initial location is crucial

- two things to decide:
 1. which direction?
 2. how much?

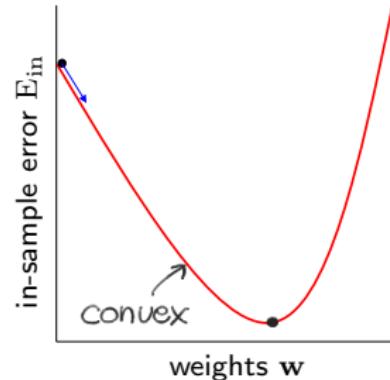
$$\Delta \mathbf{w} = -\nabla E_{in}(\mathbf{w}(t))$$

(η : learning rate
 $-\nabla$: negative gradient
 빙향으로 뛴다.)

→ Iterative 알고리즘

Advantage of using cross-entropy error

- there is only one valley!
 - ▶ no matter where we start
 - ▶ same (unique) *global minimum*
- $E_{in}(w)$ for cross-entropy error:
 - ▶ a convex function of w
(check: Hessian of E_{in} is positive semidefinite)
→ no local minima



How to ‘roll’ down E_{in} -surface

- assume:
 - ▶ you are at weights $\mathbf{w}(t)$, and
 - ▶ take a step of size η in the direction of *unit vector* $\hat{\mathbf{v}}$
unit vector

- new weights:

$$\underline{\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \cdot \hat{\mathbf{v}}}$$

- question: which is the best direction of $\hat{\mathbf{v}}$?

- ▶ pick $\hat{\mathbf{v}}$ s.t. $E_{in}(\mathbf{w}(t+1))$ becomes as small as possible
- ↔ maximize $E_{in}(\mathbf{w}(t)) - E_{in}(\mathbf{w}(t+1))$
- ↔ minimize $E_{in}(\mathbf{w}(t+1)) - E_{in}(\mathbf{w}(t))$

Recall

- Taylor expansion of vector function $f(\mathbf{x})$ to "1st order"

$$f(\mathbf{x} + \boldsymbol{\delta}) - f(\mathbf{x}) \approx \underbrace{\nabla f(\mathbf{x})^T \boldsymbol{\delta}}_{\text{기울기 } \times \text{ 밀변}}$$

- for arbitrary vector \mathbf{x} and unit vector \mathbf{y}

$$\underbrace{\mathbf{x}^T \mathbf{y}}_{\geq -||\mathbf{x}||} \stackrel{\text{($|\mathbf{x}| \cos \theta$)}}{\sim}$$

▶ equality holds if and only if

$\mathbf{x}^T \mathbf{y} \neq 0$ 인 때
 $\mathbf{y} = ?$

$$\mathbf{y} = -\frac{\mathbf{x}}{||\mathbf{x}||}$$

(*i.e.*, the angle between \mathbf{x} and \mathbf{y} is 180°)

Negative gradient gives the steepest descent

- approximating the change in E_{in}

$$\begin{aligned}\Delta E_{in} &= E_{in}(\mathbf{w}(t+1)) - E_{in}(\mathbf{w}(t)) \\ &= E_{in}(\mathbf{w}(t) + \eta \hat{\mathbf{v}}) - E_{in}(\mathbf{w}(t)) \\ &\approx \eta \nabla E_{in}(\mathbf{w}(t))^T \hat{\mathbf{v}} \\ &\geq -\eta \|\nabla E_{in}(\mathbf{w}(t))\|\end{aligned}$$

- the best (steepest) direction to move is the negative gradient:

$$\boxed{\hat{\mathbf{v}} = -\frac{\nabla E_{in}(\mathbf{w}(t))}{\|\nabla E_{in}(\mathbf{w}(t))\|}} \quad (22)$$

가장 steepest

Fixed learning-rate gradient descent

- gradient descent algorithm (with fixed η):

1: initialize weights at time step $t = 0$ to $\mathbf{w}(0)$

2: **for** $t = 0, 1, 2, \dots$ **do**

3: compute gradient $\nabla E_{in}(\mathbf{w}(t))$

4: set direction to move: $\hat{\mathbf{v}}_t = -\frac{\nabla E_{in}(\mathbf{w}(t))}{\|\nabla E_{in}(\mathbf{w}(t))\|}$

5: update weights: $\mathbf{w}(t + 1) = \mathbf{w}(t) + \eta \hat{\mathbf{v}}_t$

$$\boxed{\mathbf{w}(t + 1) = \mathbf{w}(t) - \eta \frac{\nabla E_{in}(\mathbf{w}(t))}{\|\nabla E_{in}(\mathbf{w}(t))\|}} \quad (23)$$

6: iterate to next step until it is time to stop

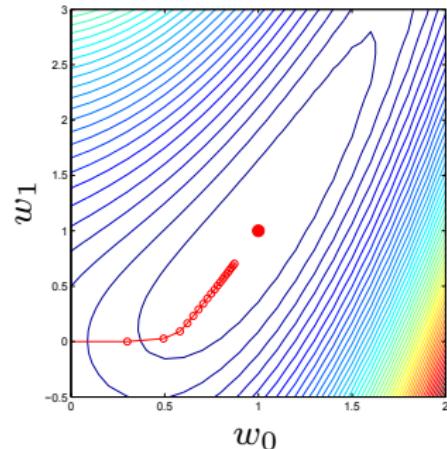
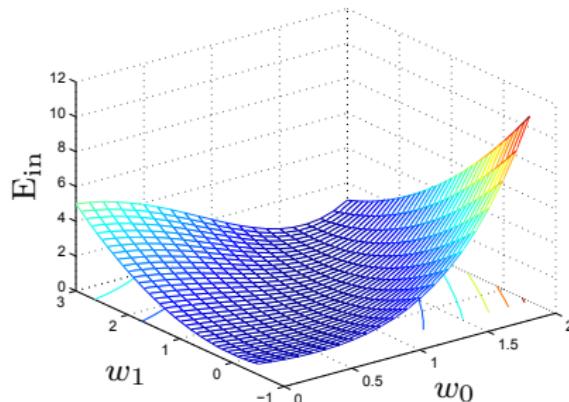
7: return final weights

Remarks

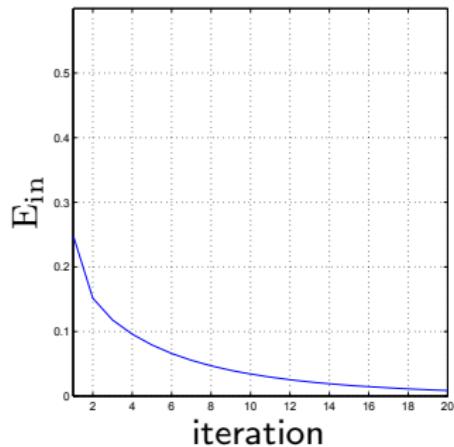
- to use gradient descent, one must compute gradient
 - ▶ this can be done explicitly for logistic regression: Eq. (21)
- we can apply gradient descent to logistic regression
 - ▶ to approximately minimize $E_{in}(\mathbf{w})$ in Eq. (6)
→ convex 어떤 정학적가
- parameter η has to be specified
 - ▶ step size $\underline{\eta}$ is called learning rate (가장 중요한 하이퍼 파라미터)
 - ▶ typically good choice: around ~~0.1~~ (empirical)
case by case, 각을 봄
 - ▶ how does the value of η affect gradient descent?

Example

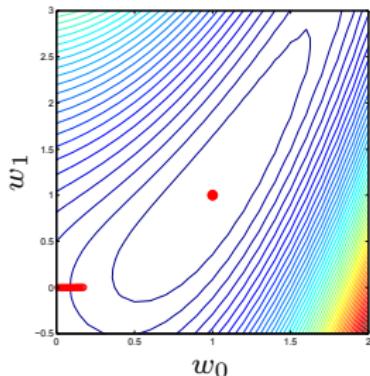
- assume a simple function for E_{in}
 - global minimum 0 at $(1,1)$



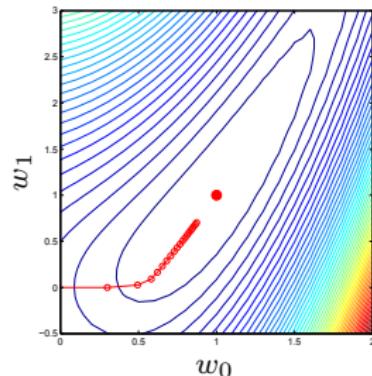
- gradient descent starting at $(0,0)$
 - # iterations (steps) = 20
 - step size (learning rate) $\eta = 0.3$



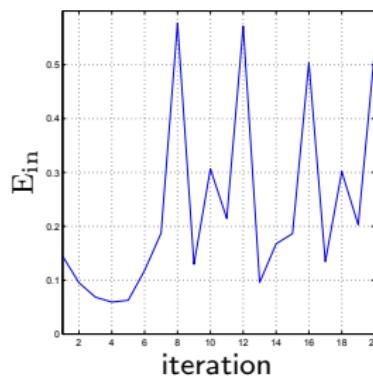
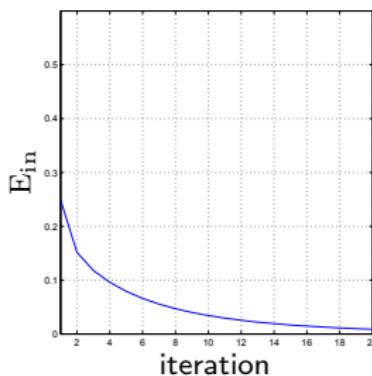
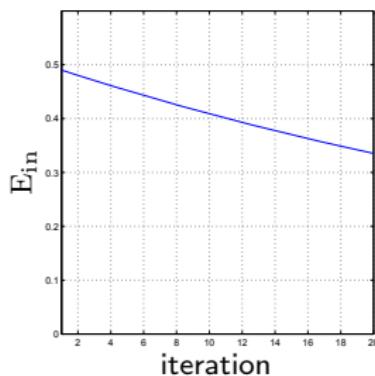
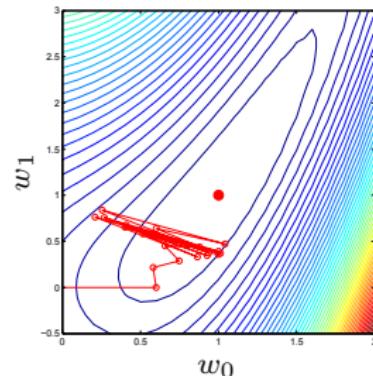
$$\eta = 0.01$$



$$\eta = 0.3$$

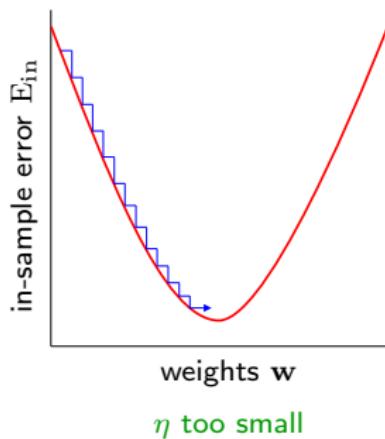


$$\eta = 0.6$$

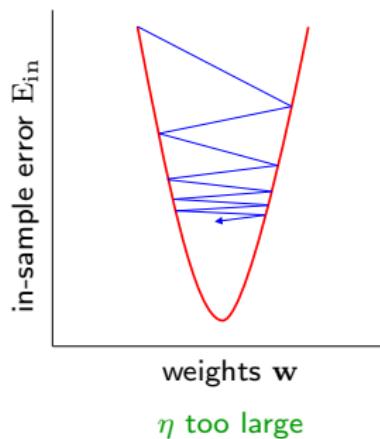


Effect of step size η

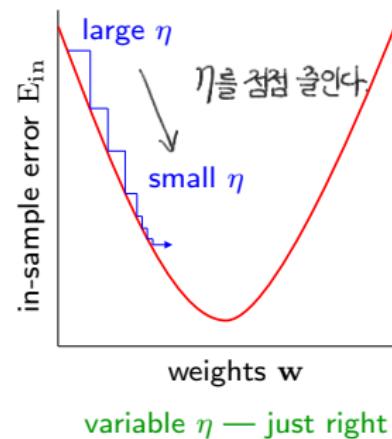
- η affects the convergence behavior of gradient descent:



η too small



η too large



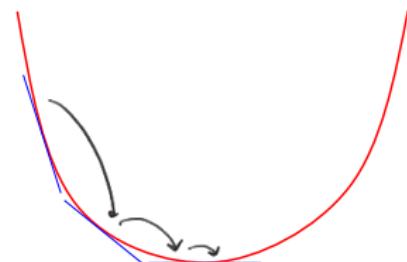
variable η — just right

How to choose η ?

- a fixed step size (if too small)
 - ▶ insufficient when you are far from the local minimum
- too large a step size when you are close to the minimum
 - ▶ leads to bouncing around, possibly even increasing E_{in}
- ideal strategy:
 - ▶ take large steps when far from the minimum
(this is to get in the right ballpark quickly)
 - ▶ and then small (more careful) steps when close to the minimum
- a simple heuristic can accomplish this

Achieving variable step size η

- observe that the norm of gradient
 - ▶ typically large when far from the minimum
 - ▶ small when close to the minimum
- thus, we could set $\eta_t = \frac{\eta}{\|\nabla E_{in}(\mathbf{w}(t))\|}$
 - ▶ to obtain the desired behavior for the variable step size
 - ▶ also conveniently cancels normalizing term of \hat{v} in Eq. (22)



Gradient descent (revised)

- gradient descent algorithm (with redefined learning rate):

- 1: initialize weights at time step $t = 0$ to $\mathbf{w}(0)$
- 2: **for** $t = 0, 1, 2, \dots$ **do**
- 3: compute gradient $\nabla E_{in}(\mathbf{w}(t))$
- 4: set the direction to move: $\hat{\mathbf{v}}_t = -\frac{\nabla E_{in}(\mathbf{w}(t))}{\|\nabla E_{in}(\mathbf{w}(t))\|}$
- 5: set learning rate: $\eta_t = \eta \|\nabla E_{in}(\mathbf{w}(t))\|$
- 6: update weights: $\mathbf{w}(t + 1) = \mathbf{w}(t) + \eta_t \hat{\mathbf{v}}_t$

$$\boxed{\mathbf{w}(t + 1) = \mathbf{w}(t) - \eta \nabla E_{in}(\mathbf{w}(t))}$$



(24)

- 7: iterate to next step until it is time to stop
- 8: return final weights

Remarks

- update rule (24)

$$\boxed{\mathbf{w}(t+1) = \mathbf{w}(t) - \eta \nabla E_{in}(\mathbf{w}(t))}$$

- ▶ needs to examine all examples at each iteration: $O(N)$ time
 - ▶ can be approximated in $O(1)$ by stochastic gradient descent
SGD (모든 exemplo를 필요 없게 비껴남)
- ※ least mean square (LMS) rule (Widrow-Hoff learning rule)
- ▶ when the cost function is mean squared error

$$\begin{aligned}\mathbf{w} &:= \mathbf{w} - \eta \nabla E_{in}(\mathbf{w}) \\ &= \mathbf{w} + \eta \underbrace{(\mathbf{y} - \mathbf{w}^T \mathbf{x})}_{\text{error}} \cdot \underbrace{\mathbf{x}}_{\text{input}}\end{aligned}$$

Outline

Logistic Regression

Predicting Probability

Cross-Entropy Error Measure

Training via Gradient Descent

Logistic Regression Algorithm

Summary

Appendix: Entropy

Logistic regression algorithm

- 1: initialize weights at time step $t = 0$ to $\mathbf{w}(0)$
- 2: **for** $t = 0, 1, 2, \dots$ **do**
- 3: compute the gradient

$$\nabla E_{in}(\mathbf{w}(t)) = -\frac{1}{N} \sum_{n=1}^N \frac{y_n \mathbf{x}_n}{1 + e^{y_n \mathbf{w}^T(t) \mathbf{x}_n}}$$

- 4: set the direction to move: $\mathbf{v}_t = -\nabla E_{in}(\mathbf{w}(t))$
- 5: update weights: $\mathbf{w}(t+1) = \mathbf{w}(t) + \eta \mathbf{v}_t$
- 6: iterate to next step until it is time to stop
- 7: return final weights \mathbf{w}

Remarks

- vector v_t is no longer restricted to unit length
 - ▶ due to "cancellation" of $\|\nabla E_{in}(\mathbf{w}(t))\|$ term
- we have two more loose ends to tie
 1. how to choose initial weights $\mathbf{w}(0)$
 2. how to set the criterion for stopping gradient descent

Initialization

- in some cases (*e.g.* logistic regression) \rightarrow convex
 - ▶ setting $\mathbf{w}(0) = \mathbf{0}$ works well
- however, in general
 - ▶ it is safer to initialize weights randomly
 - ▶ to avoid getting stuck on a perfectly symmetric hilltop
- in practice
 - ▶ choose each weight independently from a normal distribution with zero mean and small variance

Two basic criteria for termination

종료기

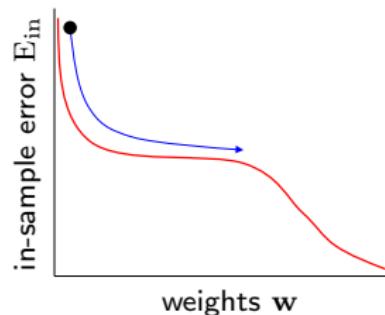
- non-trivial topic in optimization
- one simple approach (as in pocket algorithm)
 - ▶ set an upper bound on the number of iterations
 - ▶ problem: no guarantee on quality of final weights
- another natural termination criterion: $\|\nabla E_{in}\|$ 이 적은 ϵ 까지 stop.
 - ▶ stop once $\|\nabla E_{in}(\mathbf{w}(t))\|$ drops below a threshold
 - ▶ eventually this must happen, but we do not know when
- for logistic regression, combination works well:
 1. large upper bound for the number of Iterations
 2. small lower bound for the size of gradient

Third criterion for termination

- problem with relying solely on the size of gradient
 - ▶ you might stop prematurely

- when relatively flat region is reached
(more common than we suspect)
 - ▶ algorithm will prematurely stop when we may want to continue

- one solution is to require both:
 - (i) change in error is small, and (ii) error itself is small



Termination criteria in practice



언제 끝내야 하는가?

- ultimately, a combination of criteria works reasonably well:
 1. maximum number of iterations
 2. marginal error improvement (기울기 ↓)
 3. small value for the error itself (에러 자체가 ↓)

Variants of gradient decent

1st order
method

- **batch** gradient descent

- ▶ use (all N examples) in each iteration

전체만 훈련할 경우.

- **stochastic** gradient descent

- ▶ use (1 example) in each iteration

- **mini-batch** gradient descent

- ▶ use b examples in each iteration
 - ▶ b : mini-batch size (typically 2 ~ 100)

Outline

Logistic Regression

Predicting Probability

Cross-Entropy Error Measure

Training via Gradient Descent

Logistic Regression Algorithm

Summary

Appendix: Entropy

Three linear models

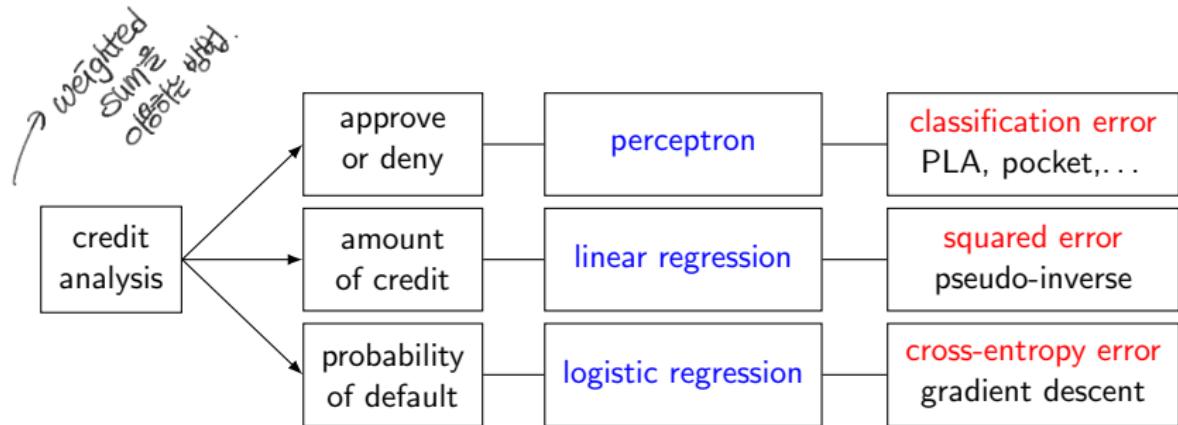


Figure 4: comparison of the three linear models considered

- have their respective goals, error measures, and algorithms
 - ▶ nonetheless, share similar sets of linear hypotheses
- you should "first try" a linear model when learning from data!

Comparison

	linear classification	linear regression	logistic regression
\mathcal{Y}	$\{-1, +1\}$	\mathbb{R}	$\{-1, +1\} \cup \{\infty\}$
$\hat{y} = h(\mathbf{x})$	$\text{sign}(\mathbf{w}^T \mathbf{x})$	$\mathbf{w}^T \mathbf{x}$	$\theta^*(\mathbf{w}^T \mathbf{x})$
$e(\hat{y}, y)$	0-1 loss $\llbracket \hat{y} \neq y \rrbracket$	squared error $(\hat{y} - y)^2$	cross-entropy error $\llbracket y=+1 \rrbracket \ln \frac{1}{\hat{y}} + \llbracket y=-1 \rrbracket \ln \frac{1}{1-\hat{y}}$
$E_{\text{in}}(h)$	$\frac{1}{N} \sum_{n=1}^N \llbracket h(\mathbf{x}_n) \neq y_n \rrbracket$	$\frac{1}{N} \sum_{n=1}^N (h(\mathbf{x}_n) - y_n)^2$	$\frac{1}{N} \sum_{n=1}^N \ln \left(1 + e^{-y_n \mathbf{w}^T \mathbf{x}_n} \right)$
opt.	combinatorial optimization (NP-hard)	set $\nabla E_{\text{in}}(\mathbf{w}) = 0$ (closed-form solution exists)	set $\nabla E_{\text{in}}(\mathbf{w}) = 0$ iterative optimization (e.g. gradient descent)

* logistic sigmoid $\theta(s) = 1/(1 + e^{-s})$

Summary

- logistic regression: estimates probabilities of binary events
 - ▶ approximates $\mathbb{P}[y = +1|\mathbf{x}] = f(\mathbf{x})$ by $h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x})$
 - ▶ likelihood: $P(y|\mathbf{x}) = h(\mathbf{x})^{\mathbb{I}[y=+1]}(1 - h(\mathbf{x}))^{\mathbb{I}[y=-1]} = \theta(y\mathbf{w}^T \mathbf{x})$ ~~prob.~~
 - ▶ training: minimize cross-entropy error iteratively
- gradient descent: iteratively solves unconstrained optimization
 - ▶ general technique for minimizing twice-differentiable functions
 - ▶ batch gradient descent: $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla E_{in}(\mathbf{w})$ [over all data]
 - ▶ stochastic gradient descent: $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla e_n(\mathbf{w})$ [pointwise]
 - ▶ mini-batch gradient descent, batch normalization
- linear models should be a first try: simple/robust/work & generalize well

Outline

Logistic Regression

Predicting Probability

Cross-Entropy Error Measure

Training via Gradient Descent

Logistic Regression Algorithm

Summary

Appendix: Entropy

(Shannon) entropy and relative entropy

- entropy
 - ▶ entropy $H(X)$ of discrete rv X with alphabet \mathcal{X} and pmf p :

$$H(X) = H(p) = - \sum_{x \in \mathcal{X}} p(x) \log p(x)$$

- relative entropy (KL divergence)
 - ▶ relative entropy $D(p||q)$ of pmf p wrt pmf q :
- $$D(p||q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$$
- ▶ measures the information lost when q is used to approximate p

Cross entropy

- cross entropy
 - ▶ cross entropy for two pmfs p and q :

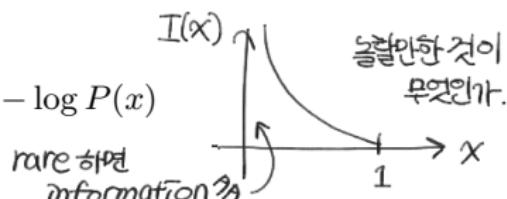
$$\begin{aligned} H(p, q) &= H(p) + D(p||q) \\ &= \sum_x p(x) (-\log p(x) + \log p(x) - \log q(x)) \\ &= - \sum_x p(x) \log q(x) \end{aligned}$$

- ▶ measures the average number of bits needed to identify an event from a set of possibilities, if a coding scheme is used based on a given probability distribution q , rather than the “true” distribution p

Comparison

- self information:

$$I(x) = -\log P(x)$$



눌렀을 때
무엇인가.

- Shannon entropy:

- ▶ min # of bits per msg needed (on average) to encode events from P (比特)

$$H(P) = \mathbb{E}_{x \sim P}[I(x)] = -\mathbb{E}_{x \sim P}[\log P(x)]$$

random valuable

- KL divergence: (두 분포사이 거리)
 (P 는 x 가 따르는 분포)
 Q 는 x 가 따르는 분포.)

- ▶ expected # of extra bits per msg to encode events from true P if using an optimal code for Q (rather than P)

$Q || P$ 이 둘이 다르다.

$$D_{KL}(P||Q) = \mathbb{E}_{x \sim P} \left[\log \frac{P(x)}{Q(x)} \right] = -H(P) - \mathbb{E}_{x \sim P}[\log Q(x)]$$

- cross entropy:

- ▶ expected # of total bits per msg to encode events from true P if using an optimal code for Q (rather than P)

(P : 그분이 아는 data 발생방법.
 Q : 우리가 알고 있는 것.)

$$H(P, Q) = \underbrace{H(P)}_{\text{total bit}} + \underbrace{D_{KL}(P||Q)}_{\text{extra bit}} = -\mathbb{E}_{x \sim P}[\log Q(x)]$$