

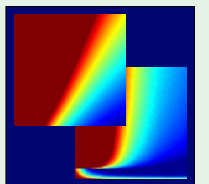
Learning From Data

Yaser S. Abu-Mostafa
California Institute of Technology

Lecture 9: **The Linear Model II**



Sponsored by Caltech's Provost Office, E&AS Division, and IST • Tuesday, May 1, 2012



Where we are

- Linear classification ✓
- Linear regression ✓
- Logistic regression
- Nonlinear transforms ✗

Logistic regression - Outline

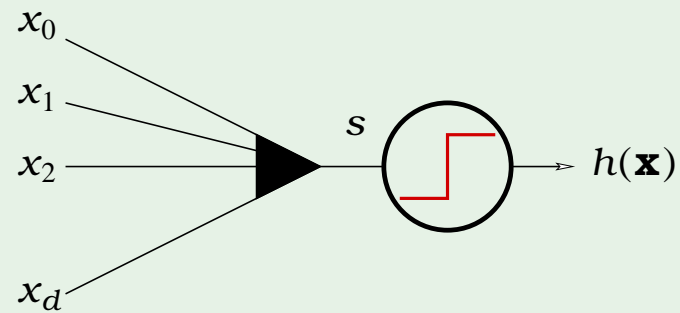
- The model
- Error measure
- Learning algorithm

A third linear model

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

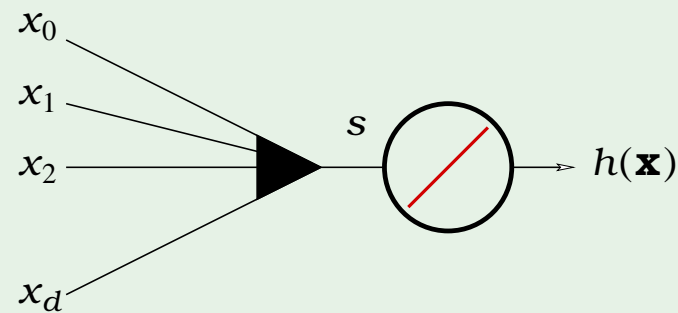
linear classification

$$h(\mathbf{x}) = \text{sign}(s)$$



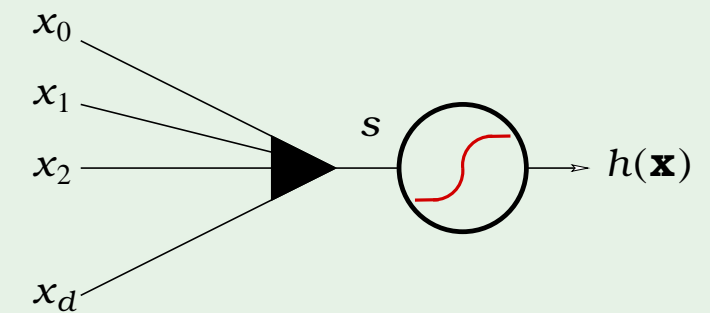
linear regression

$$h(\mathbf{x}) = s$$



logistic regression

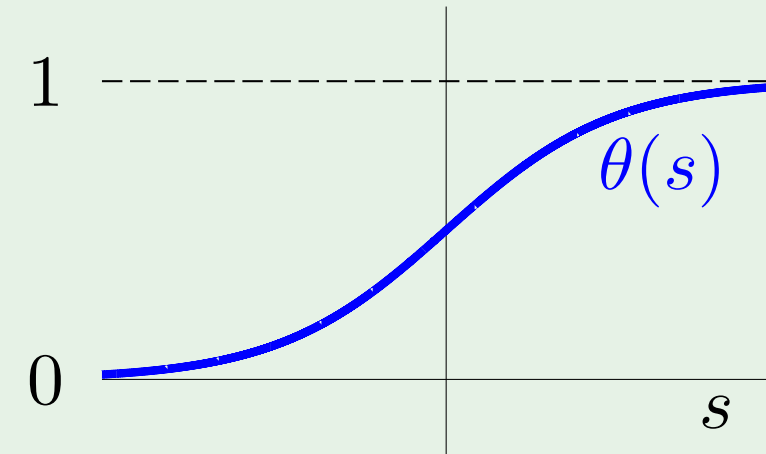
$$h(\mathbf{x}) = \theta(s)$$



The logistic function θ

The formula:

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



soft threshold: uncertainty

sigmoid: flattened out 's'

Probability interpretation

$h(\mathbf{x}) = \theta(s)$ is interpreted as a probability

Example. Prediction of heart attacks

Input \mathbf{x} : cholesterol level, age, weight, etc.

$\theta(s)$: probability of a heart attack

The signal $s = \mathbf{w}^T \mathbf{x}$ “risk score”

Genuine probability

Data (\mathbf{x}, y) with **binary** y , generated by a noisy target:

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

The target $f : \mathbb{R}^d \rightarrow [0, 1]$ is the probability

$$\text{Learn } g(\mathbf{x}) = \theta(\mathbf{w}^\top \mathbf{x}) \approx f(\mathbf{x})$$

Error measure

For each (\mathbf{x}, y) , y is generated by probability $f(\mathbf{x})$

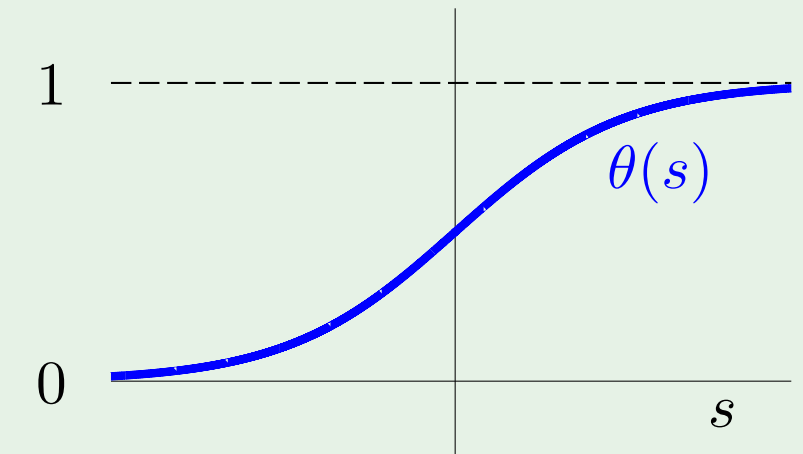
Plausible error measure based on **likelihood**:

If $h = f$, how likely to get y from \mathbf{x} ?

$$P(y \mid \mathbf{x}) = \begin{cases} h(\mathbf{x}) & \text{for } y = +1; \\ 1 - h(\mathbf{x}) & \text{for } y = -1. \end{cases}$$

Formula for likelihood

$$P(y \mid \mathbf{x}) = \begin{cases} h(\mathbf{x}) & \text{for } y = +1; \\ 1 - h(\mathbf{x}) & \text{for } y = -1. \end{cases}$$



Substitute $h(\mathbf{x}) = \theta(\mathbf{w}^\top \mathbf{x})$, noting $\theta(-s) = 1 - \theta(s)$

$$P(y \mid \mathbf{x}) = \theta(y \mathbf{w}^\top \mathbf{x})$$

Likelihood of $\mathcal{D} = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ is

$$\prod_{n=1}^N P(y_n \mid \mathbf{x}_n) = \prod_{n=1}^N \theta(y_n \mathbf{w}^\top \mathbf{x}_n)$$

$$\begin{aligned} \theta(s) &= e^s / (1 + e^s) \\ \text{so, } \theta(-s) &= e^{-s} / (1 + e^{-s}) \\ &= 1/e^s \times (1/(1 + e^{-s})) \\ &= 1/(e^s + 1) \\ &= 1 - e^s / (1 + e^s) = 1 - \theta(s) \\ \text{hence, } \theta(-s) &= 1 - \theta(s) \end{aligned}$$

Maximizing the likelihood

Minimize

$$-\frac{1}{N} \ln \left(\prod_{n=1}^N \theta(y_n \mathbf{w}^\top \mathbf{x}_n) \right)$$

Apply rule: $\ln(M \times N) = \ln M + \ln N$
then, $\ln(M^{-1}) = -\ln(M)$

$$= \frac{1}{N} \sum_{n=1}^N \ln \left(\frac{1}{\theta(y_n \mathbf{w}^\top \mathbf{x}_n)} \right)$$

Divide numerator and denominator of $\theta(s)$ by e^s

$$\left[\theta(s) = \frac{1}{1 + e^{-s}} \right]$$

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \underbrace{\ln \left(1 + e^{-y_n \mathbf{w}^\top \mathbf{x}_n} \right)}_{e(h(\mathbf{x}_n), y_n)}$$

Intuitively, $\ln(1 + \exp(-y_n \mathbf{w}^\top \mathbf{x}_n))$ should work.
if $\mathbf{w}^\top \mathbf{x}_n$ is very positive and y_n is 1 then it gives $\exp(\text{high negative}) \cong 0$ and then $\ln(1) \cong \text{approx } 0$ which is expected.

“cross-entropy” error

Logistic regression - Outline

- The model
- Error measure
- Learning algorithm

How to minimize E_{in}

For logistic regression,

Cross Entropy Error

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

Compare to linear regression:

its derivation is easy: $e^{-y\mathbf{w}^T \mathbf{x}} / (1 + e^{-y\mathbf{w}^T \mathbf{x}}) * (-y\mathbf{x})$ since: $\text{derv}(\ln x) = 1/x$
 $= (1/e^{y\mathbf{w}^T \mathbf{x}}) * (1 + 1/e^{y\mathbf{w}^T \mathbf{x}}) * (-y\mathbf{x})$
 $= (1/e^{y\mathbf{w}^T \mathbf{x}}) * ((e^{y\mathbf{w}^T \mathbf{x}} + 1)/e^{y\mathbf{w}^T \mathbf{x}}) * (-y\mathbf{x})$
 $= -y\mathbf{x} / (1 + e^{y\mathbf{w}^T \mathbf{x}})$

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - y_n)^2 \quad \begin{array}{l} \text{squared error} \\ \leftarrow \text{closed-form solution} \end{array}$$

Iterative method: gradient descent

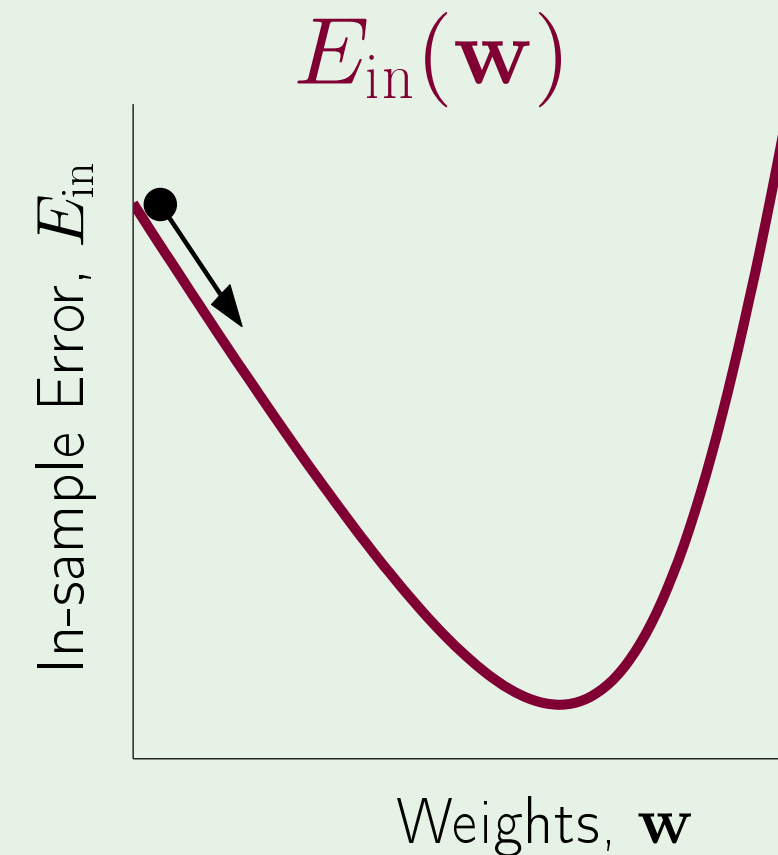
General method for nonlinear optimization

Start at $\mathbf{w}(0)$; take a step along steepest slope

Fixed step size: $\mathbf{w}(1) = \mathbf{w}(0) + \eta \hat{\mathbf{v}}$

What is the direction $\hat{\mathbf{v}}$?

For multi-dimension V has infinite directions. Here only 2D is shown.



Formula for the direction $\hat{\mathbf{v}}$

Taylor Series:
Approximate $f(x)$

$$f(x) \cong f(a) + (x-a)f'(a)/1! + (x-a)^2f''(a)/2! + \dots$$

or

$$f(x+a) \cong f(a) + (x)f'(a)/1! + (x)^2f''(a)/2! + \dots$$

for Error E_{in} or E :

$$E(\eta\mathbf{v}+\mathbf{w}_0) \cong E(\mathbf{w}_0) + (\eta\mathbf{v})E'(\mathbf{w}_0)/1! + (\eta\mathbf{v})^2E''(\mathbf{w}_0)/2! + \dots$$

$$E(\eta\mathbf{v}+\mathbf{w}_0) - E(\mathbf{w}_0) \cong E(\mathbf{w}_0) - E(\mathbf{w}_0) + (\eta\mathbf{v})E'(\mathbf{w}_0)/1! + (\eta\mathbf{v})^2E''(\mathbf{w}_0)/2! + \dots$$

$$\cong (\eta\mathbf{v})E'(\mathbf{w}_0) + (\eta\mathbf{v})^2E''(\mathbf{w}_0)/2! + \dots$$

$$\therefore \Delta E \cong (\eta\mathbf{v})\nabla E(\mathbf{w}_0) + O(\eta^2)$$

$$\Delta E_{in} = E_{in}(\mathbf{w}(0) + \eta\hat{\mathbf{v}}) - E_{in}(\mathbf{w}(0))$$

$$= \eta \nabla E_{in}(\mathbf{w}(0))^T \hat{\mathbf{v}} + O(\eta^2)$$

$$\geq \eta(\mathbf{v})\nabla E_{in}$$

$$\geq \eta(\nabla E_{in}/\|\nabla E_{in}\|) \nabla E_{in} \geq \eta\|\nabla E_{in}\| \quad \text{Nooo... it is always } > 0.$$

$$\geq -\eta\|\nabla E_{in}(\mathbf{w}(0))\|$$

We need ΔE_{in} to be less than 0.

$$\boxed{A \cdot A^T = \|A\| \|A\| \cos 0 = \|A\|^2}$$

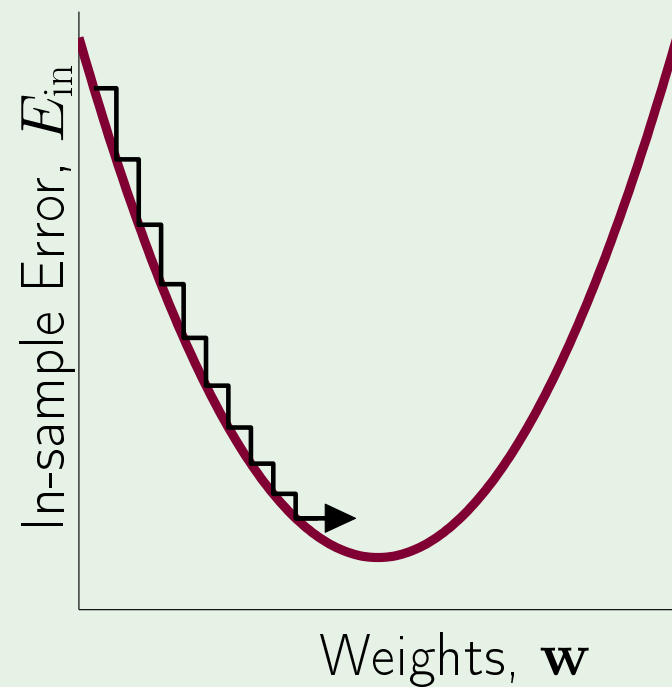
$$\therefore A^* A^T / \|A\| = \|A\|$$

Since $\hat{\mathbf{v}}$ is a unit vector,

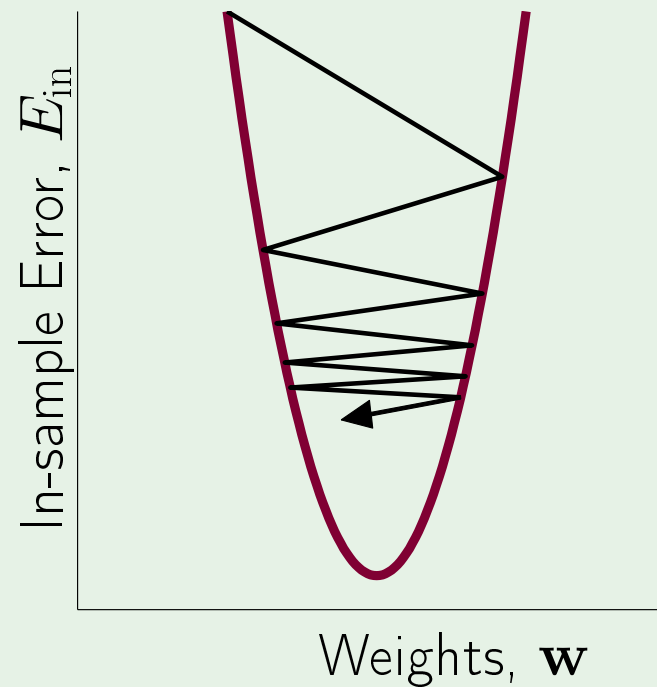
$$\hat{\mathbf{v}} = - \frac{\nabla E_{in}(\mathbf{w}(0))}{\|\nabla E_{in}(\mathbf{w}(0))\|}$$

Fixed-size step?

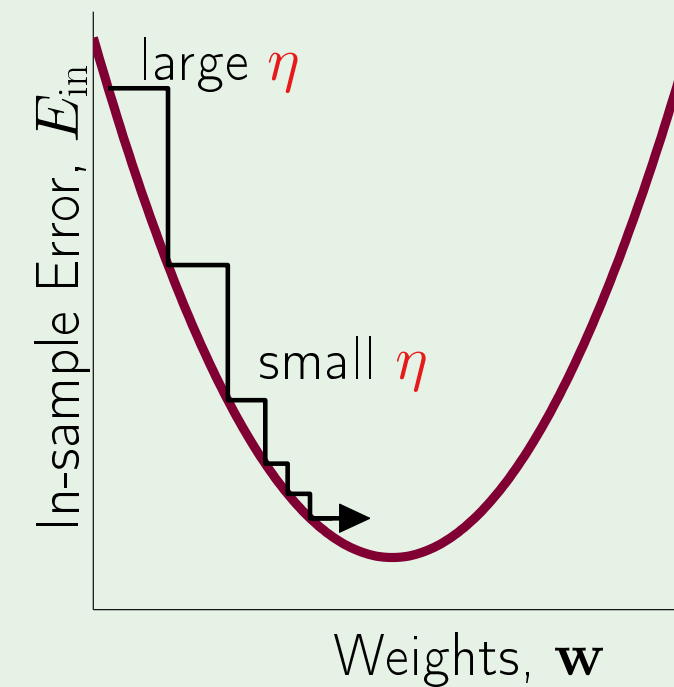
How η affects the algorithm:



η too small



η too large



variable η – just right

η should increase with the slope

Easy implementation

Instead of

$$\begin{aligned}\Delta \mathbf{w} &= \eta \hat{\mathbf{v}} \\ &= -\eta \frac{\nabla E_{\text{in}}(\mathbf{w}(0))}{\|\nabla E_{\text{in}}(\mathbf{w}(0))\|}\end{aligned}$$

since eta is prop size of ∇E
 $\eta = \text{const} * \|\nabla E\|$
 $\therefore \Delta W = -\eta \nabla E / \|\nabla E\|$
 $= -\text{const} * \|\nabla E\| * \nabla E / \|\nabla E\|$
 $= -\text{const} \nabla E$
lets call this constant another eta
 $\therefore \Delta W = -\eta \nabla E$

Have

$$\Delta \mathbf{w} = -\eta \nabla E_{\text{in}}(\mathbf{w}(0))$$

Fixed learning rate η

Logistic regression algorithm

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

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

- 4: Update the weights: $\mathbf{w}(t + 1) = \mathbf{w}(t) - \eta \nabla E_{\text{in}}$
- 5: Iterate to the next step until it is time to stop
- 6: Return the final weights \mathbf{w}

Summary of Linear Models

