

### Announcements

- Homework 3 has been posted.
- 

### Today's Lecture

- Logistic regression
-

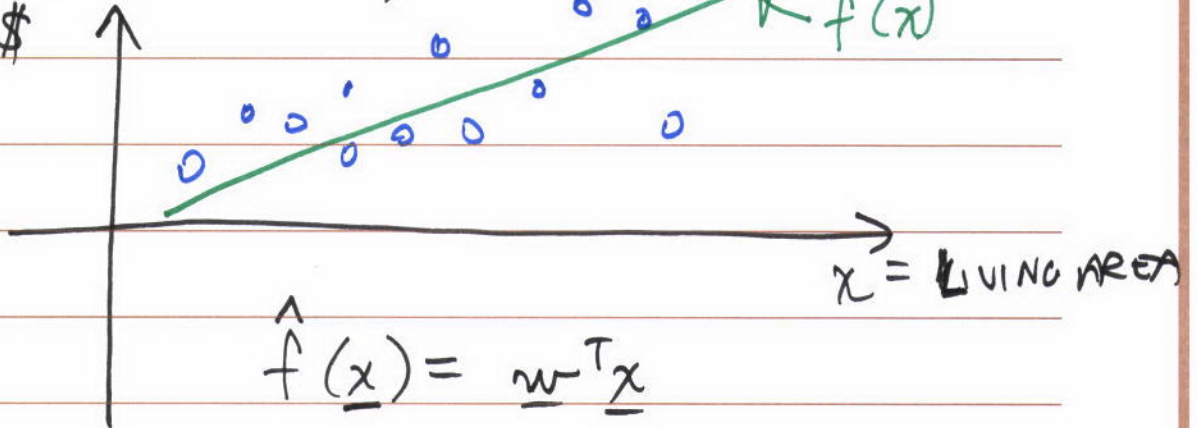
# LOGISTIC REGRESSION

## INTRO. (AUGMENTED NOTATION)

### REGRESSION

(1D, LINEAR)

$y = \$$



### CLASSIFICATION (LINEAR)

setosa +1

$y$

$\hat{f}(x)$

$x$

$y$

$x$

$x$

$x$

$x$

$x$

$x$

$x$

$x$

$x$

$x$

$x$

$x$

$\hat{f}(x)$

$g(x)$

$x$

$y$

$x$

$y$

$x$

$y$

$x$

$y$

$x$

$x = \text{PETAL LENGTH}$

virginica -1

$y$

$x$

$y$

$x$

$y$

$x$

$y$

$x$

$y$

$x$

$y$

$x$

$y$

$x$

$y$

$x$

$y$

$x$

$y$

$x$

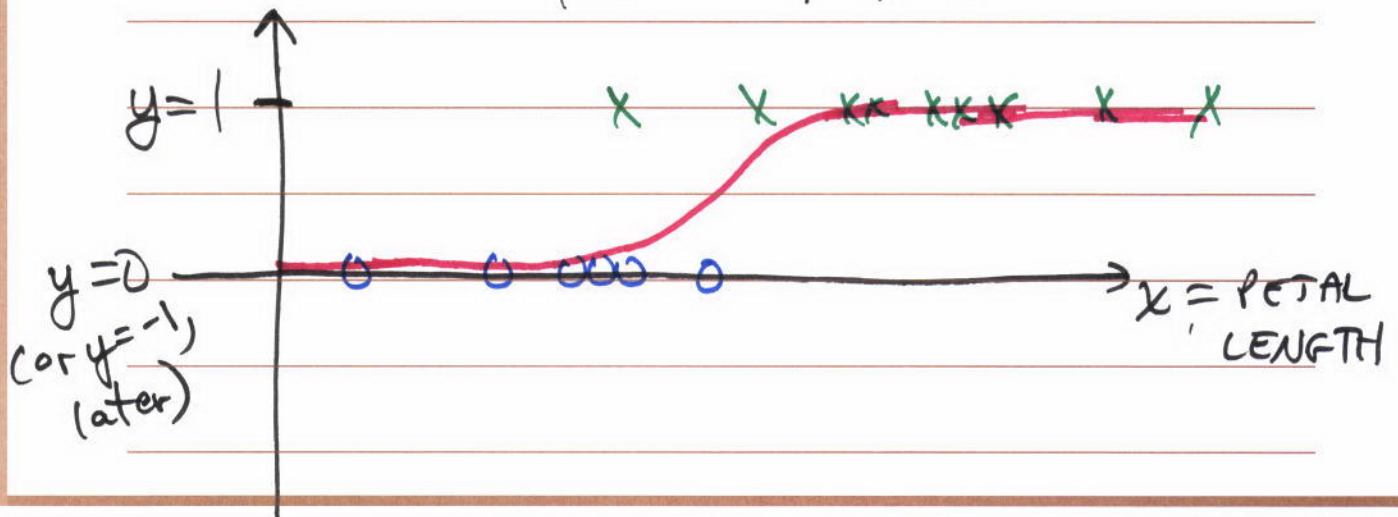
$$g(x) = \underline{w}^T \underline{x}$$

$$\hat{y}(x) = \hat{f}(x) = \text{sgn}\{g(x)\} = \text{sgn}\{\underline{w}^T \underline{x}\}$$

## LOGISTIC REGRESSION

$$\hat{f}(x) = p(y=1 | x, \mathcal{D})$$

$$= p(\text{setosa} | x, \mathcal{D})$$



$$\hat{f}(x) = \text{sigm} \{ \underline{w}^T \underline{x} \} \quad [\text{Murphy notation}]$$

$$= \frac{e^{\underline{w}^T \underline{x}}}{1 + e^{\underline{w}^T \underline{x}}}$$

$$= \theta(\underline{w}^T \underline{x}) \quad [\text{AML notation}]$$

= "logistic" or "sigmoid" function  
(of  $\underline{w}^T \underline{x}$ )

MODEL

$$p(y | \underline{x}, \underline{w}) = \text{Ber}(y | \text{sigm}(\underline{w}^T \underline{x}))$$

$$= \mu^{\mathbb{I}(y=1)} (1-\mu)^{\mathbb{I}(y=0)} \quad [\text{M 8.2}]$$

IN WHICH  $\mu = \text{sigm}(\underline{w}^T \underline{x})$

CHANGE OUTPUT ( $y$ ) REPRESENTATION:

$$\text{LET } \tilde{y} = 2y - 1 \Rightarrow \tilde{y} \in \{-1, +1\}$$

$$p(\tilde{y} | \underline{x}, \underline{w}) = \mu^{\mathbb{I}(\tilde{y}=1)} (1-\mu)^{\mathbb{I}(\tilde{y}=-1)}$$

$$p(\tilde{y} | \underline{x}, \underline{w}) = [\text{sigm}(\tilde{y} \underline{w}^T \underline{x})]^{\mathbb{I}(\tilde{y}=1)} \cdot [\text{sigm}(\tilde{y} \underline{w}^T \underline{x})]^{\mathbb{I}(\tilde{y}=-1)}$$

can show:  $\text{sigm}(-s) = 1 - \text{sigm}(s)$

$$p(\tilde{y} | \underline{x}, \underline{w}) = \text{sigm}(\tilde{y} \underline{w}^T \underline{x})$$



# OBJECTIVE FCN.

## MAX. LIKELIHOOD

$$\text{LIKELIHOOD} = p(\mathcal{D} | \underline{w})$$

$$= \prod_{i=1}^N p(\tilde{y}_i | x_i, \underline{w})$$

$$= \prod_{i=1}^N \left( \frac{e^{\tilde{y}_i \underline{w}^T x}}{1 + e^{\tilde{y}_i \underline{w}^T x}} \right) \cdot \frac{e^{-()}}{e^{-()}}$$

$$= \prod_{i=1}^N \left( \frac{1}{e^{-\tilde{y}_i \underline{w}^T x} + 1} \right)$$

$$-l(\underline{w}) = \text{NLL}(\underline{w}) = \sum_{i=1}^N \ln[1 + e^{-\tilde{y}_i \underline{w}^T x}]$$

$$= J(\underline{w}, \mathcal{D})$$

↑ ~~the~~ OBJ. FCN. FOR MLE OF  $\underline{w}$  IN LOGISTIC REGRESSION.

$J$  IS CONVEX.

$$\text{LET: } -l(\underline{w}) = \sum_{i=1}^N E_i, \quad E_i = \ln[1 + e^{-\tilde{y}_i \underline{w}^T x}]$$

$\tilde{y}_i$	$\underline{w}^T \underline{x}_i$	$E_i = \ln[1 + e^{-\tilde{y}_i \underline{w}^T \underline{x}_i}]$
+1	> 0	$0 < E_i < \ln 2$
+1	>> 0	$E_i \approx 0$
+1	< 0	$E_i > \ln 2$
+1	<< 0	$E_i >> \ln 2$
:		$\hookrightarrow E_i \approx -\underline{w}^T \underline{x}_i$

How to minimize J?  $\rightarrow$  J is CONVEX.

CAN WE  $\nabla_{\underline{w}} J(\underline{w}, \mathcal{D}) = 0$ , SOLVE  
ALGEBRAICALLY?

$\rightarrow$  NOT TRACTIBLE.

How ELSE?

- ITERATIVE, GRADIENT-BASED TECHNIQUES.  
 $\rightarrow$  GRADIENT DESCENT (STEEPEST DESCENT)

$$\underline{w}(k+1) = \underline{w}(k) - \eta(k) \nabla_{\underline{w}} J$$

(BATCH UPDATE)

## > STOCHASTIC GRADIENT DESCENT

(i) INITIALIZING  $\underline{w}(0)$

(ii) ITERATE OVER DATA POINTS UNTIL

CONVERGENCE:

(a) RANDOMLY PICK A DATA POINT  
(WITH REPLACEMENT)

(b) PERFORM A SINGLE-SAMPLE

UPDATE:

$$\underline{w}(i+1) = \underline{w}(i) - \eta(i) \nabla_{\underline{w}} E_i$$

( $i$  = ITERATION NUMBER)

CAN SHOW:

$$\nabla_{\underline{w}} E_i = \frac{\tilde{y}_i \underline{x}_i}{1 + e^{\tilde{y}_i \underline{w}^T(i) \underline{x}_i}}$$

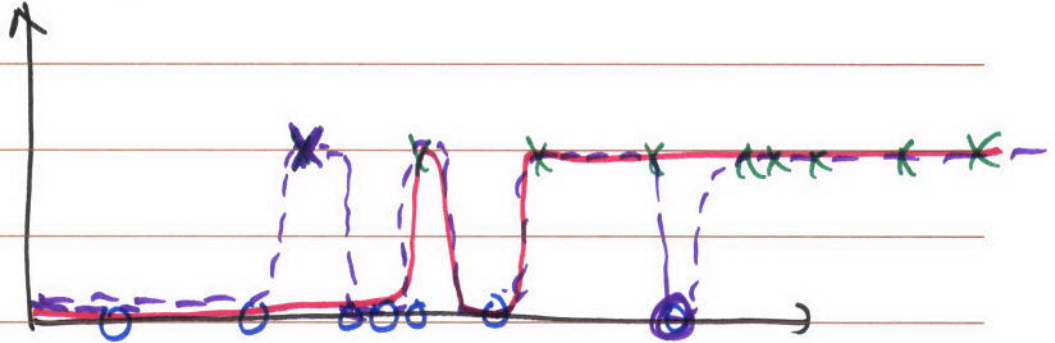
(e.g., see  
EES59 or  
Discussion3).

OTHER TECHNIQUES (e.g., Murphy  
8.3.2-8.3.5).



- IS REGULARIZATION USEFUL (AND POSSIBLE) IN LOGISTIC REGRESSION?

- CAN L.R. OVERFIT?



$$\hat{f}(x) = \text{sign} \{ \underline{w}^T \underline{\phi}(x) \}$$

↑ POLYNOMIAL IN  $x$  TERMS

→ YES, WOULD BE USEFUL.

→ YES, WE CAN ADD A PRIOR TERM  $p(\underline{w})$ ,

⇒ MAP EST. [↔ RIDGE REGRESSION]

OR BAYESIAN INFERENCE.

→ LOGISTIC REGRESSION:

$$\rightarrow J(\underline{w}, \mathcal{D}) = \underbrace{NLL(\underline{w})}_{\text{from logistic regr., MLE}} + \lambda \|\underline{w}\|_2^2$$

from  
logistic  
regr., MLE.

(IF  $p(\underline{w}) = \text{GAUSSIAN}$   
 $= N(\underline{w} | 0, \tau^2)$ )