

LOGISTIC REGRESSION

- CLASSIFICATION
- REGRESSION

$$\min_w \frac{1}{n} \sum_{i=1}^n (w^T x^i - y^i)^2$$

$$(*) \quad \mathbb{E}[Y|X] = w^T X$$

EVEN IF

(*) DOES NOT
HOLD

CLASSIFICATION

WHAT IS LOSS FUNCTION FOR CLASSIFICATION?

OUR GUESS FOR LABEL $y^i \in \{0, 1\}$ IS
GOING TO BE $\text{SIGN}(w^T x^i)$ FOR SOME VECTOR w .

• PENALIZED IF $\text{SIGN}(w^T x^i) \neq y^i$
IF $\text{SIGN}(w^T x^i) = y^i$ (NO PENALTY)

EXAMINE QUANTITY $y^i \cdot (w^T x^i)$ IF QUANT IS POS NO PENALTY
" IS NEG PENALTY

0-1 LOSS
REG

$$\ell(z) = \begin{cases} 1 & \text{IF } z \leq 0 \\ 0 & \text{IF } z > 0 \end{cases}$$

SQ LOSS

$$(w^T x^i - y^i)^2$$

CLASSIFICATION

$$\ell_{0-1}(y^i \cdot w^T x^i)$$

OPTIMIZATION PROBLEM ASSOCIATED WITH CLASSIFICATION

$$\min_w \frac{1}{m} \sum_{i=1}^m \ell_{0-1}(y^i \cdot w^T x^i)$$

↑
WANT TO SOLVE.

$$y^i \in \{-1, 1\}$$

QUESTION: WHAT IF THERE
IS NO MARGIN?

THERE MIGHT NOT EXIST A w
s.t. $\text{SIGN}(w^T x^i) = y^i \forall i$

WHEN DOES
PERCEPTRON FIND
A w WITH SMALL
LOSS?

RECALL THAT PERCEPTRON
REQUIRED $\exists w$ s.t. $\forall x$

$$y \cdot \underline{w^T x} > \rho \Rightarrow$$

$$\text{CONVERGENCE} \\ \# \text{ MISTAKES} < \frac{1}{\rho^2}$$

OPTIMIZATION PROBLEM ASSOCIATED WITH CLASSIFICATION

$$\min_w \frac{1}{m} \sum_{i=1}^m \ell_{0-1}(y^i \cdot w^T x^i)$$

WHEN DOES
PERCEPTRONS FIND
LOSS?

WHAT IS THE COMPUTATIONAL COMPLEXITY OF THIS OPTIMIZATION PROBLEM? A w WITH SMALL LOSS?

BAD NEWS: THIS PROBLEM IS NP-HARD

$$y^i \in \{-1, +1\}$$

RECALL THAT PERCEPTRON
REQUIRED $\exists w$ s.t. $\forall x$

$$y \cdot \underline{w^T x} > \rho \Rightarrow$$

$$\text{CONVERGENCE} \\ \# \text{ MISTAKES} < \frac{1}{\rho^2}$$

QUESTION: WHAT IF THERE
IS NO MARGIN?

THERE MIGHT NOT EXIST A w
s.t. $\text{SIGN}(w^T x^i) = y^i \forall i$

OPTIMIZATION PROBLEM ASSOCIATED WITH CLASSIFICATION

$$\min_w \frac{1}{m} \sum_{i=1}^m \ell_{0-1}(y^i \cdot w^T x^i)$$

WHEN DOES
PERCEPTRONS FIND
LOSS?

WHAT IS THE COMPUTATIONAL COMPLEXITY OF THIS OPTIMIZATION PROBLEM?

BAD NEWS: THIS PROBLEM IS NP-HARD
"AGNOSTICALLY LEARNING A HALFSPACE"
AGNOSTIC LEARNING

RECALL THAT PERCEPTRON
REQUIRED $\exists w$ s.t. $\forall x$

$$y \cdot \underline{w^T x} > \rho \Rightarrow$$

$$\# \text{ MISTAKES} < \frac{1}{\rho^2}$$

QUESTION: WHAT IF THERE
IS NO MARGIN?

THERE MIGHT NOT EXIST A w
s.t. $\text{SIGN}(w^T x^i) = y^i \forall i$

TO SUMMARIZE

REGRESSION \leadsto CONVEX LOSS FUNCTION ✓

CLASSIFICATION \leadsto NON CONVEX LOSS (0-1) LOSS (BAD NEWS)

- IDEA LET'S RELAX THE 0-1 LOSS TO A DIFFERENT NIKER LOSS
"SURROGATE LOSSES" RELATED TO 0-1
LOSS BUT CONVEX

LET'S INTRODUCE A FEW LOSSES:

- φ LOGISTIC LOSS
- φ HINGE
- φ EXP

$$\ell_{\text{LOGISTIC}}(z) = \log(1 + e^{-z})$$

$$\ell_{\text{LOGISTIC}}(y^i \cdot w^T x^i) = \log(1 + e^{-(y^i \cdot w^T x^i)})$$

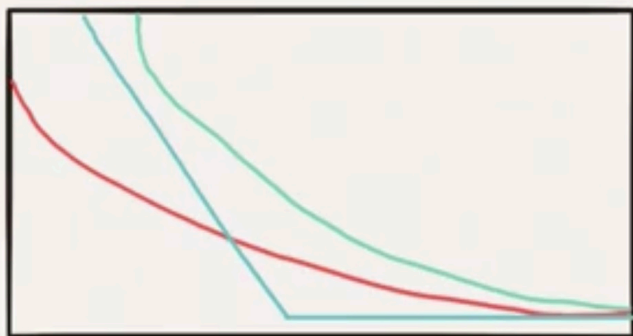
• If $\underbrace{y^i w^T x^i}_{\text{MARGIN}} \ll 0 \Rightarrow \ell_{\text{LOGISTIC}}(y^i \cdot w^T x^i)$ is LARGE
 $\gg 0 \Rightarrow \ell_{\text{LOGISTIC}}(y^i \cdot w^T x^i)$ is SMALL (MOVING TO $\rightarrow 0$)

$$\ell_{\text{Hinge}}(z) = \max\{1 - z, 0\}$$

$$\ell_{\text{Hinge}}(y^i \cdot w^T x^i)$$

LARGE WHEN $y^i \cdot w^T x^i$ IS NEG
 SMALL WHEN $\underline{y^i \cdot w^T x^i}$ IS POS

$$\ell_{\text{EXP}}(z) = e^{-z}$$



$$(z = y^T w^T x) \nearrow$$

LOGISTIC LOSS

HINGE LOSS

EXP LOSS

OPTIMIZATION PROBLEM ASSOCIATED WITH LOGISTIC LOSS?

$$L(w) = \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-y^i \cdot w^T x^i))$$

$$\min_w L(w)$$

ENTER THE SIGMOID FUNCTION:

$$g(z) = \frac{1}{1 + e^{-z}}$$

As z gets large $g(z) \rightarrow 1$
As z gets small $g(z) \rightarrow 0$



$$g(z) = \frac{1}{1+e^{-z}}$$

$$Y \in \{0, +1\}$$

$$\underline{E[Y|X]} = g(Y \cdot \underline{w^T X})$$

FACT: $g(z) + g(-z) = 1$

FOR SOME $\underline{w} \quad \exists \underline{w}$

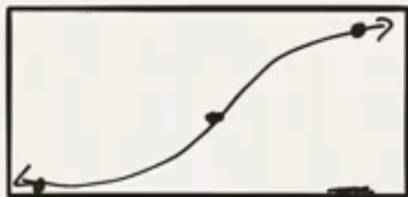
$$\frac{e^z}{e^z + 1} + \frac{1}{1+e^z}$$

$$\Rightarrow \Pr[Y=1 | \underline{X}] = g(\underline{w^T X})$$

$$\frac{e^z}{1+e^z} + \frac{1}{1+e^z} = 1.$$

GIVEN \underline{X} IF $\underline{w^T X}$ IS LARGE

IF $\underline{w^T X}$ IS NEG
SMALL



$$\Pr[Y=y^i | X=x^i; \underline{w}] = g(y^i \cdot \underline{w}^T x^i)$$

"MODEL FOR LOGISTIC REGRESSION"

$$g(z) = \frac{1}{1 + e^{-z}}$$

GIVEN A TRAINING SET $S = \{(x^1, y^1), \dots, (x^m, y^m)\}$

WHAT IS THE MOST LIKELY \underline{w} GIVEN THE TRAINING SET?

$$\text{LIKELIHOOD}(\underline{w}) = \prod_{i=1}^m p(Y=y^i | X=x^i; \underline{w}) = \prod_{i=1}^m g(y^i \cdot \underline{w}^T x^i)$$

MAX \underline{w} \nearrow

$$\text{LOG-LIKELIHOOD}(\underline{w}) = \sum_{i=1}^m \log g(y^i \cdot \underline{w}^T x^i)$$

LOGISTIC
 \nearrow LOSS

$$= - \sum_{i=1}^m \log (1 + \exp(-y^i \cdot \underline{w}^T x^i)) = -m \cdot L(\underline{w})$$

Now our goal will be to minimize logistic loss $L(w)$.

How should we minimize logistic loss $L(w)$?

IDEA: RUN GRADIENT DESCENT ON LOGISTIC LOSS.

↑ THIS IS THE ALGORITHM FOR PERFORMING LOGISTIC REGRESSION!

LET'S SAY WE FIND w^1 ; FOR FUTURE EXAMPLES WE LABEL THEM +1 WITH PROB $g(w^{1T}x)$

LET'S COMPUTE THE GRADIENT OF $L(w)$

LOGISTIC LOSS

$$\varphi_{\text{logistic}}(z) = \log(1 + e^{-z})$$

THIS PRECISELY TELLS US
HOW TO FIND MAX UNBELLANED
 w .

$$1. \quad \varphi'_{\text{logistic}}(z) = \frac{-e^{-z}}{1+e^{-z}} = -\frac{1}{1+e^z} = -g(-z)$$

2. COMPUTE

$$\frac{\partial \varphi_{\text{logistic}}(y \cdot w^T x)}{\partial w_k} = -g(-y \cdot w^T x) \cdot \underbrace{y \cdot x_k}_{\text{CHAIN RULE.}}$$

WITH THIS FORMULA WE CAN DIRECTLY APPLY
GRADIENT DESCENT

WHAT HAPPENS IF WE HAVE MULTIPLE
LABELS FOR y ? $y \in \{0, 1\}$

WHAT IF $y \in \{1, \dots, K\}$

MULTINOMIAL LOGISTIC REGRESSION: w^1, \dots, w^{K-1}

$$P_r[y=1 | X] = \frac{e^{w^1 \cdot x}}{\sum_{i=1}^K e^{w^i \cdot x}}$$

$$P_r[y=i | X] = \frac{e^{w^i \cdot x}}{\sum_{j=1}^K e^{w^j \cdot x}}$$

WHAT HAPPENS IF WE HAVE MULTIPLE
LABELS FOR y ? $y \in \{0, 1\}$

WHAT IF $y \in \{1, \dots, K\}$

MULTINOMIAL LOGISTIC REGRESSION: w^1, \dots, w^{K-1}

$$Pr[y=1|x] \propto e^{w^1 \cdot x}$$

$$Pr[y=k] =$$

$$Pr[y=i|x] \propto e^{w^i \cdot x}$$

$$1 - \sum_{i=1}^{K-1} Pr[y=i]$$

WHAT IS THE ASSOCIATED LOSS?

CROSS-ENTROPY LOSS

GENERALIZATION OF LOGISTIC LOSS

(IMAGINE y IS A VECTOR OF LENGTH K

WITH A 1 IN THE j^{th} POSITION IF

CORRECT LABEL IS j) (ONE-HOT ENCODING OF LABELS).

LET'S SAY OUR GUESS FOR THE PROB y HAS LABEL i

IS $\underline{p_i}$

$$-\sum_{i=1}^K y_i \log(p_i)$$

SUMMARY \rightarrow TAKES REAL-VALUES INTO PROBABILITIES

$W^T x$ VIA SIGMOID $W^T x \rightarrow [\sigma]$
INTO A PROBABILITY

$$\underbrace{(z_1, \dots, z_k)}_{K \text{ COORDINATES}} \rightarrow \underbrace{\left(\frac{e^{z_1}}{z}, \frac{e^{z_2}}{z}, \dots, \frac{e^{z_k}}{z}\right)}$$

$$z = \sum_{i=1}^k e^{z_i}$$