

# Midterm Review Friday

- Come with questions

## Logistic Regression

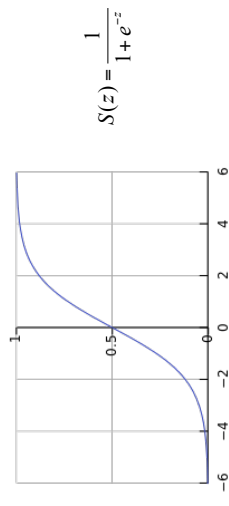
### Logistic Regression

Idea:

- Naïve Bayes allows computing  $P(Y|X)$  by learning  $P(Y)$  and  $P(X|Y)$
- Why not learn  $P(Y|X)$  directly?

### Idea

- Similar to perceptron, but **sigmoid** function applied to linearity.



- **Error function** is based on **Max Likelihood** → as a result we get **genuine probabilities** as output of prediction (not just 1,0 discrete values).

## Models so far, and the new one

$$z = \sum_{i=1}^m w_i x_i$$

Perceptron

$$h(\mathbf{x}, \mathbf{w}) = \text{sign}(z)$$

Linear regression

$$h(\mathbf{x}, \mathbf{w}) = z$$

**Logistic regression**

$$h(\mathbf{x}, \mathbf{w}) = S(z)$$

Output  $h(\mathbf{x}, \mathbf{w})$  will be interpreted as probability.

Why? Because of the range of  $S$ , and particular cost function we'll optimize.

## Probability Interpretation

- Assume tuples  $(\mathbf{x}, y)$ , where  $y$  is binary, are generated from some noisy data-source according to some distribution  $f$ .

$$p(y | \mathbf{x}) = \begin{cases} f(\mathbf{x}) & \text{for } y = 0 \\ 1 - f(\mathbf{x}) & \text{for } y = 1 \end{cases}$$

**For mathematical convenience we will use  $[0, 1]$  as our label set instead of  $[-1, 1]$**

- We will learn  $f(\mathbf{x})$  by approximating it with the sigmoid  $S(z)$ , i.e.  $S(\mathbf{w} \cdot \mathbf{x})$   
 $\mathbf{w} = [w_0, w_1, \dots, w_m]$      $\mathbf{x} = [x_0 = 1, x_1, \dots, x_m]$
- Makes sense as  $S(z)$  is a function from 0 to 1.

## Approximation

- What is the probability of  $y=0$ ?  
(according to our approximation)  $p(y=0|\mathbf{x}) = \frac{1}{1+e^{w\mathbf{x}}}$
- What is the probability of  $y=1$ ?  
(according to our approximation)  $p(y=1|\mathbf{x}) = 1 - \frac{1}{1+e^{w\mathbf{x}}}$   

$$= \frac{1+e^{w\mathbf{x}}-1}{1+e^{w\mathbf{x}}}$$

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

## Training Logistic Regression: MCLE

- Choose parameters  $W = \langle w_0, \dots, w_n \rangle$  to maximize conditional likelihood of training data  
**where**  $p(y=0|\mathbf{x}) = \frac{1}{1+e^{w\mathbf{x}}}$   
 $p(y=1|\mathbf{x}) = \frac{e^{w\mathbf{x}}}{1+e^{w\mathbf{x}}}$
- Training data  $D = \{ \langle X^1, Y^1 \rangle, \dots, \langle X^N, Y^N \rangle \}$
- Data likelihood =  $\prod_{i=1}^N P(\langle X^i, Y^i \rangle > |W)$
- Data conditional likelihood =  $\prod_{i=1}^N P(Y^i | X^i, W)$

$$W_{MCLE} = \underset{w}{\operatorname{argmax}} \prod_{i=1}^N P(Y^i | X^i, W)$$

## Expressing Conditional Log Likelihood

$$\begin{aligned} \mathcal{L}(W) &= \sum_i Y^i \ln P(Y^i = 1 | X^i, W) + \sum_i (1 - Y^i) \ln P(Y^i = 0 | X^i, W) \\ &= \sum_i Y^i \ln \frac{P(Y^i = 1 | X^i, W)}{P(Y^i = 0 | X^i, W)} + \sum_i \ln P(Y^i = 0 | X^i, W) \end{aligned}$$

## Training Logistic Regression: Maximum Conditional Likelihood Estimation (MCLE)

- we have L training examples:  $\{ \langle X^1, Y^1 \rangle, \dots, \langle X^N, Y^N \rangle \}$
- maximum likelihood estimate for parameters  $W$   

$$W_{MLE} = \underset{w}{\operatorname{argmax}} P(\langle X^1, Y^1 \rangle, \dots, \langle X^N, Y^N \rangle > |W)$$

$$= \underset{w}{\operatorname{argmax}} \prod_{i=1}^N P(\langle X^i, Y^i \rangle > |W)$$
- maximum conditional likelihood estimate

## Expressing Conditional Log Likelihood

$$\begin{aligned} \mathcal{L}(W) &\equiv \ln \left( \prod_i P(Y^i | X^i, W) \right) \\ &= \sum_i \ln P(Y^i | X^i, W) \end{aligned}$$

$$\begin{aligned} p(y=0|\mathbf{x}) &= \frac{1}{1+e^{w\mathbf{x}}} \\ p(y=1|\mathbf{x}) &= \frac{e^{w\mathbf{x}}}{1+e^{w\mathbf{x}}} \end{aligned}$$

$$\mathcal{L}(W) = \sum_i Y^i \ln P(Y^i = 1 | X^i, W) + \sum_i (1 - Y^i) \ln P(Y^i = 0 | X^i, W)$$

For the samples with  $y^i=0$

For the samples with  $y^i=1$

## Expressing Conditional Log Likelihood

$$\begin{aligned} \mathcal{L}(W) &= \sum_i Y^i \ln P(Y^i = 1 | X^i, W) + \sum_i (1 - Y^i) \ln P(Y^i = 0 | X^i, W) \\ &= \sum_i Y^i \ln \frac{P(Y^i = 1 | X^i, W)}{P(Y^i = 0 | X^i, W)} + \sum_i \ln P(Y^i = 0 | X^i, W) \\ &= \sum_i Y^i (w^T x^i) - \sum_i \ln(1 + e^{w^T x^i}) \end{aligned}$$

$$\begin{aligned} p(y=0|\mathbf{x}) &= \frac{1}{1+e^{w\mathbf{x}}} \\ p(y=1|\mathbf{x}) &= \frac{e^{w\mathbf{x}}}{1+e^{w\mathbf{x}}} \end{aligned}$$

## Maximizing Conditional Log Likelihood

$$\begin{aligned}\mathcal{L}(W) &= \sum_i Y^i (w'x^i) - \sum_i \ln(1 + e^{w'x^i}) \\ \frac{d}{dw} \mathcal{L}(W) &= \sum_i Y^i x^i - \sum_i \frac{x^i e^{w'x^i}}{1 + e^{w'x^i}} \\ &= \sum_i x^i (Y^i - P(Y = 1|W, x^i))\end{aligned}$$

- **Good news:**  $L(W)$  is concave function of  $W$
- **Bad news:** no closed-form solution to maximize  $L(W)$  [i.e., no canonical equation]

## Gradient Descent Algorithm

Initialize  $w_0 = 0$

For  $t=0, 1, 2, \dots$  do

    Compute the gradient and update the weights

$$\frac{d}{dw} \mathcal{L}(W) = \sum_i x^i (Y^i - P(Y = 1|W, x^i))$$

$$w \leftarrow w + \kappa \left( \sum_i x^i (Y^i - P(Y = 1|W, x^i)) \right)$$

    Iterate with the next step until  $w$  doesn't change too much  
    (or for a fixed number of iterations)

Return final  $w$ .

## Iterative Method

- Start at  $w_0 = 1$ ; take a step along **steepest slope**
- Fixed step size:  
 $w_1 = w_0 + \eta v$
- $v$  is a **vector** in the direction of the **steepest slope**.  
– What's the steepest slope?

## Making predictions

- A new tuple comes:  $(x, ?)$

$$p(y = 0 | x) = \frac{1}{1 + e^{w'x}}$$

$$p(y = 1 | x) = \frac{e^{w'x}}{1 + e^{w'x}}$$

- Fix a threshold in  $[0, 1]$  to make predictions.  
 $p(y = 1 | x) > \text{threshold}$       Predict  $y = 1$   
 $p(y = 1 | x) \leq \text{threshold}$       Predict  $y = 0$

## Example

GPA, GRE, and success:

	Dummy GPA	GRE	y
100, 800, 1	1	1.0	1
90, 800, 1	1	0.9	1
90, 700, 1	1	0.9	0.875
70, 600, 0	1	0.7	0.75
60, 700, 0	1	0.6	0.875
60, 700, 1	1	0.6	0.875
50, 600, 0	1	0.5	0.75
50, 650, 0	1	0.5	0.8125
50, 800, 1	1	0.5	1.0
50, 700, 0	1	0.5	0.875
50, 700, 1	1	0.5	0.875

Scaled so that max GPA/GRE is 1

## Example

Logistic regression:

$$-7.44 + 4.4056 * \text{GPA} + 5.57 * \text{GRE}$$

$$p(y = 0 | gpa, gre) = \frac{1}{1 + e^{(-7.44 + 4.4056 * \text{GPA} + 5.57 * \text{GRE})}}$$

Fix a threshold in  $[0, 1]$  to make predictions.

$p(y = 1 | gpa, gre) > \text{threshold}$       Predict  $y = 1$

$p(y = 1 | gpa, gre) \leq \text{threshold}$       Predict  $y = 0$

## Odds

Definition:  $odds(0 \text{ vs. } 1 \text{ given } \mathbf{x}) = \frac{p(y=0|\mathbf{x})}{p(y=1|\mathbf{x})}$

Formula:

$$odds(0 \text{ vs. } 1 \text{ given } \mathbf{x}) = \frac{1}{1 + e^{\mathbf{w}'\mathbf{x}}}$$

$$= \frac{1}{e^{\mathbf{w}'\mathbf{x}}}$$

$$= e^{-\mathbf{w}'\mathbf{x}}$$

We predict ? if this number is less than ?

## Interpretation

$odds(\text{successful vs. unsuccessful given } gpa \text{ and } gre) = e^{7.44 + 4.405 \cdot GPA + 5.57 \cdot GRE}$

If GPA increases by .1 then the odds of success will increase by 55%

$$e^{0.441} \approx 1.55$$

If GRE increases by .1 then the odds of success will increase by 75%

$$e^{0.557} \approx 1.75$$

## Logistic Regression in Matlab

```
kappa = 0.1;
[n,p] = size(X);
num_its = 500;

L = zeros(1,num_its);
for t=1:num_its,
    L(t) = sum(y.*(X*w) - log(1+exp(X*w)));
    W = W + kappa * ...
        sum( X.*(repmat(y - exp(X*w), [1,p])) );
end;
plot(L)
xlabel('Likelihood')
xlabel('Iteration')
fprintf('After optimizing, accuracy is %.2f\n', 1-sum(abs((exp(-1*X*w)<1)-y))/length(y))

X=[
1 1.0 1.0;
1 0.9 1.0;
1 0.9 0.875;
1 0.7 0.75;
1 0.6 0.875;
1 0.6 0.75;
1 0.5 0.75;
1 0.5 0.8125;
1 0.5 1.0;
1 0.5 0.875;
1 0.5 0.875];

y=[1 1 0 0 1 0 1 0 1];
w=[1; 1; 1];

After optimizing, accuracy is
0.82
```

