

Learning Linear Classifiers

Maximum Likelihood estimation

Quality metric for logistic regression: Maximum likelihood estimation

Learning problem:

Training data: N observations (x_i, y_i)

$$x[1] = \# \text{awesome} \quad x[2] = \# \text{awful} \quad y = \text{sentiment}$$

2
0
3
4
1
2
0
0
2

9
3
4
3
1

+1
-1
-1
+1
+1
-1
-1
+1

Optimize quality metric on training data

$\hat{\omega}$

Finding best coefficients

$$x[1] = \# \text{awesome} \quad x[2] = \# \text{awful} \quad y = \text{sentiment}$$

	$x[1] = \# \text{awesome}$	$x[2] = \# \text{awful}$	$y = \text{sentiment}$
0	2	-1	+1
3	3	-1	+1
2	4	-1	+1
0	3	-1	+1
0	1	-1	+1

$$P(y=+1|x_i, \hat{\omega}) = 0.0$$

$$P(y=+1|x_i, \hat{\omega}) = 1.0$$

Dick $\hat{\omega}$ that makes

① quality metric = L, likelihood function

Negative data points

$$P(y=+1 | x_i, \omega) = 0.0$$

Positive data points

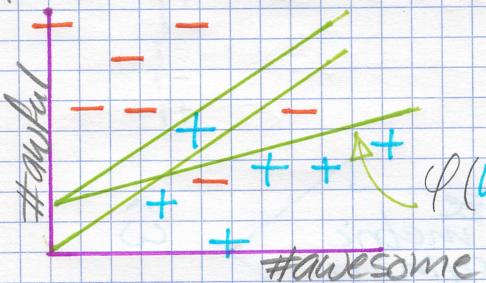
$$P(y=+1 | x_i, \omega) = 1.0$$

No $\hat{\omega}$ achieves perfect predictions (usually)

Likelihood $\varphi(\omega)$: Measures quality of fit for model with coefficients ω

Find "best" classifier

Maximize likelihood over all possible $\omega_0, \omega_1, \omega_2$



Best model:

Highest likelihood $\varphi(\omega)$
 $\omega = (w_0 = 1, w_1 = 0.5, w_2 = -1.5)$

gradient ascent to find $\hat{\omega}$

Data likelihood

Quality metric: probability of data

$$x[1] = \#awesome \quad x[2] = \#awful \quad y = \text{sentiment}$$

If model good, should predict $\hat{y}_i = +1$

Pick ω to maximize:

$$P(y=+1 | x, \omega) = P(y=+1 | x[1]=2, x[2]=1, \omega)$$

$$x[1] = \# \text{awesome}$$

$$x_2 = 0$$

$$x[2] = \# \text{awful}$$

$$y_2 = 2$$

$$\begin{aligned} y &= \\ y_2 &= -1 \end{aligned}$$

if model good, should predict
 $y_2 = -1$

Pick w to maximize
 $P(y = -1 | x_2, w)$

Maximizing (likelihood
 probability of data)

Data
 point

$x[1]$

$x[2]$

y

Choose w to maximize

x_1, y_1	2	1	+1
x_2, y_2	0	2	-1
x_3, y_3	3	3	-1
x_4, y_4	4	1	+1
x_5, y_5	1	1	+1
x_6, y_6	2	4	-1
x_7, y_7	0	3	-1
x_8, y_8	0	1	-1
x_9, y_9	2	1	+1

$$P(y = +1 | x_1, w) = P(y = +1 | x[1] = 2, x[2] = 1, w)$$

$$P(y = -1 | x_2, w)$$

$$P(y = -1 | x_3, w)$$

$$P(y = +1 | x_4, w)$$

Must combine into single measure of quality?

Multiply probability

$$P(y = +1 | x_1, w) P(y = -1 | x_2, w) P(y = -1 | x_3, w) \dots$$

Learn logistic regression model with maximum likelihood estimation (MLE)

Data point $x[1] \ x[2] \ y$ Choose w to maximize

$$x_1, y_1 \quad 2 \quad 1 \quad y_1 = +1 \quad P(y=+1|x[1]=2, x[2]=1, w)$$

$$x_2, y_2 \quad 0 \quad 2 \quad -1 \quad P(y=-1|x[1]=0, x[2]=2, w)$$

$$x_3, y_3 \quad 3 \quad 3 \quad y_3 = -1 \quad P(y=-1|x[1]=3, x[2]=3, w)$$

$$x_4, y_4 \quad 4 \quad 1 \quad +1 \quad P(y=+1|x[1]=4, x[2]=1, w)$$

$$\ell(w) = P(y=+1|x[1]=2, x[2]=1, w) P(y=-1|x[1]=0, x[2]=2, w) \cdot$$

$$\cdot P(y=-1|x[1]=3, x[2]=3, w) P(y=+1|x[1]=4, x[2]=1, w)$$

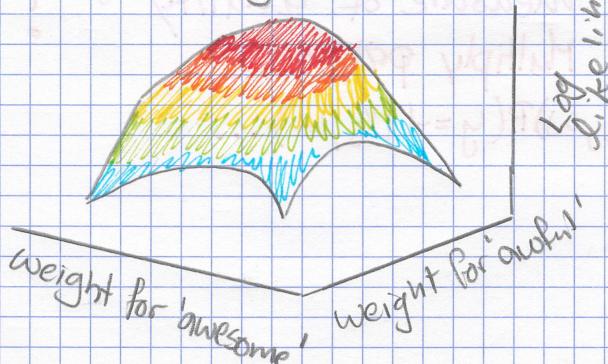
$$P(y_3|x_3, w)$$

$$\ell(w) = \prod_{i=1}^N P(y_i|x_i, w) \leftarrow \text{pick } w \text{ to make this fn. as large as possible}$$

num of data points

Finding best linear classifier with gradient ascent

Maximizing likelihood



No closed-form solution \Rightarrow use gradient ascent

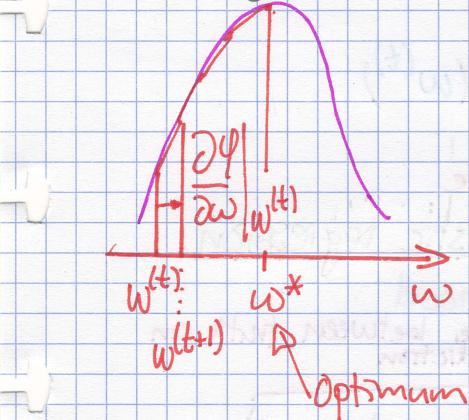
Maximize function over all possible w_0, w_1, w_2

$$\max_{w_0, w_1, w_2} \prod_{i=1}^N P(y_i|x_i, w)$$

$\ell(w_0, w_1, w_2)$ is a function of 3 variables

Review of gradient ascent

Finding the max via hill climbing



Algorithm

While not converged

$$w^{(t+1)} \leftarrow w^{(t)} + \eta \frac{\partial \phi}{\partial w} \Big|_{w^{(t)}}$$

stepsize

Convergence criteria

For convex functions, optimum occurs when

$$\frac{\partial \phi}{\partial w} = 0$$



In practice, stop when

$$\left| \frac{\partial \phi}{\partial w} \right|_{w^{(t)}} < \epsilon$$

↑ tolerance

Moving to multiple dimension:

Gradients

$$\nabla \phi(w) = \begin{bmatrix} \frac{\partial \phi}{\partial w_0} \\ \frac{\partial \phi}{\partial w_1} \\ \vdots \\ \frac{\partial \phi}{\partial w_D} \end{bmatrix}$$

D+1 dim. vector

Gradient ascent

Algorithm $w^{(0)} = 0$, random, or something else

while not converged:

$$w^{(t+1)} \leftarrow w^{(t)} + \eta \nabla \varphi(w^{(t)})$$

↑
stepsize

Learning algorithm for logistic regression

Derivative of (log-)likelihood

$$\frac{\partial \ell(w)}{\partial w_j} = \sum_{i=1}^n h_j(x_i) \left(\mathbb{1}[y_i = +1] - P(y = +1 | x_i, w) \right)$$

↓
Sum over data points
Feature value
Difference between truth and prediction
predict x_i is positive

Indicator function

$$\mathbb{1}[y_i = +1] = \begin{cases} 1 & \text{if } y_i = +1 \\ 0 & \text{if } y_i = -1 \end{cases}$$

Example of computing derivative for logistic regression

Computing derivative

$$\frac{\partial \varphi(w^{(t)})}{\partial w_j} = \sum_{i=1}^n h_j(x_i) \left(\mathbb{1}[y_i = +1] - P(y = +1 | x_i, w^{(t)}) \right)$$

$w_i^{(t)}$	$w_0^{(t)}$	$\frac{\partial \varphi}{\partial w_0}$
	0	
	-1	
	-2	

$h_1(x) = \# \text{ awesome}$

$x[1]$	$x[2]$	y	$P(y=+1 x_i, w)$	Contribution to derivative for w_1
2	1	+1	0.5	$2(1 - 0.5) = 1$
0	2	-1	0.02	$0(0 - 0.02) = 0$
3	3	-1	0.05	$3(0 - 0.05) = -0.15$
4	1	+1	0.88	$4(1 - 0.88) = 0.48$

Total derivative:

$$\frac{\partial \phi(w^{(t)})}{\partial w_1} = 1 + 0 - 0.15 + 0.48 = 1.33$$

$$w_i^{(t+1)} = w_i^{(t)} + \eta \frac{\partial \phi(w^{(t)})}{\partial w_1} \quad | \quad \eta = 0.1$$

$$= 1 + 0.1 * 1.33 = 1.133$$

Interpreting derivative for logistic regression

Derivative of (log-)likelihood interpretation

\downarrow Sum over data points Feature value Difference between truth and prediction

$$\frac{\partial \phi(w)}{\partial w_j} = \sum_{i=1}^N h_j(x_i) \left(\mathbb{1}[y_i=+1] - P(y=+1|x_i, w) \right) \Delta_i$$

If $h_j(x_i) = 1$: $P(y=+1|x_i, w) \approx 1$ $P(y=+1|x_i, w) \approx 0$

$y_i = +1$ $\Delta_i = (1 - 1) \approx 0$
 \hookrightarrow don't change anything

$\Delta_i \approx 1 \Rightarrow$ increase w_j
 \Rightarrow score(x_i) becomes larger
 $\Rightarrow P(y=+1|x_i, w)$ increases

$y_i = -1$ $\Delta_i = -1 \Rightarrow w_j$ to decrease
 \Rightarrow score(x_i) decreases
 $\Rightarrow P(y=+1|x_i, w)$ decrease

$\Delta_i \approx 0$
 \Rightarrow don't change anything

Summary of gradient ascent for Logistic regression

init $w^{(1)} = 0$ (or randomly, or smartly), $t=1$
 while $\|\nabla \varphi(w^{(t)})\| > \epsilon$ tolerance
 for $j=0, \dots, D$: coefficient

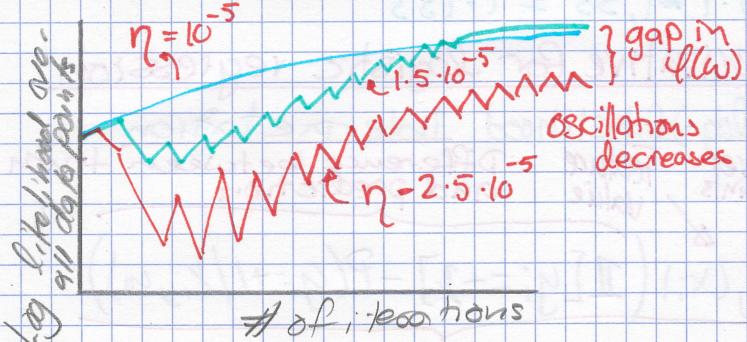
$$\text{partial}[j] = \sum_{i=1}^N h_j(x_i) (I[y_i = +1] - P(y = +1 | x_i, w^{(t)}))$$

$$w_j^{(t+1)} \leftarrow w_j^{(t)} + \eta \text{partial}[j]$$

$$\eta \leftarrow \frac{\partial \varphi(w^{(t)})}{\partial w_j}$$

Choosing step size η

- Careful with step sizes that are too large



- Very large step sizes can even cause divergence or wild oscillations

Rule of thumb for choosing step size

Simple rule of thumb for picking step size η

- Unfortunately, picking step size requires a lot of trial and error

- Try several values, exponentially spaced

- Goal: plot learning curves to find one η that is too small (smooth but moving ^{too slowly})

- find one η that is too large (oscillating or divergence)

- Try values in between to find "best" η
↳ exponentially spaced, pick one that leads best training data likelihood
- Advanced tip: can also try step size that decreases w.r.t iterations, e.g.,

$$\eta_t = \frac{\eta_0}{t}$$