

# Lecture 11

Thursday, October 5, 2023 10:02 AM

Linear regression

$$\hat{y} = \vec{w}^T \vec{x}$$

$$= w_0 + w_1 x_1 + \dots + w_d x_d$$

hyperplane

$$\vec{1} = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{pmatrix} \quad \vec{x} = \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix}$$

augmented space

$$\mathcal{L} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

SSE  
sum of squared errors

$$P(y=1 | \vec{x})$$

vs

$$P(y=0 | \vec{x}) = 1 - P(y=1 | \vec{x})$$

binary logistic regression

$y$  is categorical

$$y_i \in \{0, 1\}$$

binary class

classification

$$P(y=1 | \vec{x}) \equiv \theta(\vec{w}^T \vec{x})$$

$$\in \mathbb{R}^{d+1}$$

(augmented space)

$\theta \equiv$  Sigmoid

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

$$1 - \theta(z) = \theta(-z)$$

$$\theta(\vec{w}^T \vec{x}) = \frac{e^{\vec{w}^T \vec{x}}}{1 + e^{\vec{w}^T \vec{x}}}$$

$$P(y=0 | \vec{x}) = 1 - \theta(\vec{w}^T \vec{x}) = \theta(-\vec{w}^T \vec{x})$$

$$\theta(z) = \frac{e^z}{1+e^z}, \text{ then } \underbrace{\theta'(z)}_{\text{derivative}} = \frac{\partial \theta(z)}{\partial z} = \theta(z) \cdot (1 - \theta(z)) = \theta(z) \cdot \theta(-z)$$

$$P(y | \vec{x}) = \begin{cases} \theta(z) & \text{if } y=1 \\ 1 - \theta(z) & \text{if } y=0 \end{cases} \quad z = \vec{w}^T \vec{x}$$

$$\mathcal{D} = \left\{ \vec{x}_i, y_i \right\}_{i=1}^n \quad \text{dataset}$$

$$P(y | \vec{x}) = \theta(z)^y (1 - \theta(z))^{1-y} \quad \text{model}$$

Alternative expression

$n$  equations

$$P(y_i | \vec{x}_i) = \theta(\vec{w}^T \vec{x}_i)^{y_i} (1 - \theta(\vec{w}^T \vec{x}_i))^{1-y_i} \quad \forall i=1 \dots n$$

$\vec{w} \in \mathbb{R}^{d+1}$  is the unknown

Q: how to find  $\vec{w}$ ? what's the objective?

Objective: Maximize the likelihood  $\equiv$  Maximum likelihood Estimation (MLE)

$$\underbrace{L(\vec{w})}_{\text{likelihood}} = P(D | \vec{w})$$

$$\max_{\vec{w}} P(D | \vec{w})$$

The  $\vec{w}$  that gives the highest prob to the observed sample is the optimal one!

$$\max_{\vec{w}} L(\vec{w}) = \underbrace{P(D | \vec{w})}_{\text{joint prob over the } n \text{ points}}$$

$$= P(\{ \vec{x}_i, y_i \}_{i=1}^n | \vec{w})$$

Assume all points are independent

$$= \prod_{i=1}^n P(y_i | \vec{x}_i, \vec{w})$$

product of all prob of each point

$$L(\vec{w}) = \prod_{i=1}^n \theta(\vec{w}^T \vec{x}_i)^{y_i} (1 - \theta(\vec{w}^T \vec{x}_i))^{1-y_i}$$

$$\left( \frac{\partial L(\vec{w})}{\partial \vec{w}} \right) \equiv \nabla_{\vec{w}}$$

$$\max_{\vec{w}} \log L(\vec{w}) = \sum_{i=1}^n \log ( \theta(\vec{w}^T \vec{x}_i)^{y_i} (1 - \theta(\vec{w}^T \vec{x}_i))^{1-y_i} )$$

$$\max_{\vec{w}} \log L(\vec{w}) = \sum_{i=1}^n \log \left( \underbrace{\Theta(z_i)^{y_i} \Theta(-z_i)^{1-y_i}} \right); \quad z_i = \vec{w}^T \vec{x}_i$$

$$= \sum_{i=1}^n \log(\Theta(z_i)^{y_i}) + \log(\Theta(-z_i)^{1-y_i})$$

$$\max_{\vec{w}} \log L(\vec{w}) = \sum_{i=1}^n y_i \log \Theta(z_i) + (1-y_i) \log(1-\Theta(z_i))$$

Log-Likelihood (LL)

Convert into minimization objective (ie, turn into a loss/error)

NLL  $\equiv$  binary cross entropy

negative  
log  
likelihood

$$\min_{\vec{w}} \text{BCE}(\vec{w}) \equiv \min_{\vec{w}} \text{NLL}(\vec{w}) \equiv -\log L(\vec{w})$$

$$\max_{\vec{w}} \log L(\vec{w}) = \max_{\vec{w}} \sum_{i=1}^n \underbrace{y_i \log \Theta(z_i) + (1-y_i) \log(1-\Theta(z_i))}$$

$$\nabla_{\vec{w}} = \frac{\partial \log L(\vec{w})}{\partial \vec{w}} =$$

$$y_i \log(\text{sigmoid}(z_i))$$

$$\nabla_{\vec{w}} = \sum_{i=1}^n \left( y_i - \underbrace{\Theta(\vec{w}^T \vec{x}_i)}_n \right) \vec{x}_i \quad \frac{\partial}{\partial \vec{w}} \underbrace{y_i \log(\text{sigmoid}(\vec{w}^T \vec{x}_i))}_{\text{dot}}$$



$$\vec{v}_i = \sum_{i=1}^n (y_i - \theta(\vec{w}^T \vec{x}_i)) \vec{x}_i$$

true  
class

predicted  
prob of  
 $y_i=1$

$$\frac{\partial}{\partial \vec{w}} \sum_{i=1}^n (y_i - \theta(\vec{w}^T \vec{x}_i)) \vec{x}_i$$

$$\frac{1}{\theta(\vec{w}^T \vec{x}_i)}$$

derivative  
of log

$$\theta(\vec{w}^T \vec{x}_i) \cdot (1 - \theta(\vec{w}^T \vec{x}_i)) \vec{x}_i$$

derivative of  
sigmoid

$$\frac{\partial \log L(\vec{w})}{\partial \vec{w}} = 0$$

$$\sum (y_i - \theta(\vec{w}^T \vec{x}_i)) \cdot \vec{x}_i = 0$$

$$\sum \theta(\vec{w}^T \vec{x}_i) \vec{x}_i = \sum y_i \vec{x}_i$$

No closed form solution!

$$\vec{V}_{\vec{w}} = \sum_{i=1}^n \underbrace{(y_i - \theta(\vec{w}^T \vec{x}_i))}_{\text{scalar}} \cdot \underbrace{\vec{x}_i}_{\text{vector}}$$

$$\vec{V}_{\vec{w}} \in \mathbb{R}^{d+1}$$

$$\vec{w}^{(t)} = \vec{w}^{(t-1)} + \eta \cdot \vec{V}_{\vec{w}}$$

step size (eta)

batch gradient  
ascent!

Stochastic Gradient Ascent

based on a single point

Practical

based on a single point

$$\nabla_{\vec{w}}(x_i) = (y_i - \sigma(\vec{w}^T x_i)) \vec{x}_i$$

Practical

pick a mini-batch  
B of points

$$\vec{w}(t) = \vec{w}(t-1) + \eta \cdot \nabla_{\vec{w}}(x_i)$$

$\vec{w}$  : training

testing / inference ?

given an unknown  $\vec{x}_{n+1}$   
 $y_{n+1}$  ?

$$f(\vec{x}_{n+1}) = \begin{cases} \hat{y} = 1 & \text{if } \sigma(\vec{w}^T \vec{x}_{n+1}) \geq 0.5 \\ \hat{y} = 0 & \text{otherwise} \end{cases}$$

Multiclass Logistic Regression

$$\{\vec{x}_i, y_i\}_{i=1}^n$$

$$y_i \in \{c_1, c_2, \dots, c_k\}$$



k - symbolic classes

given  $y_i, \vec{x}_i$

predict k probabilities

$$P(y_i = c_1 | \vec{x}_i) = \pi_1(\vec{x}_i)$$

$$P(y_i = c_2 | \vec{x}_i) = \pi_2(\vec{x}_i)$$

$$P(y_i = c_1 | \vec{x}_i) = \pi_1(\vec{x}_i)$$

$$P(y_i = c_2 | \vec{x}_i) = \pi_2(\vec{x}_i)$$

⋮

$$P(y_i = c_k | \vec{x}_i) = \pi_k(\vec{x}_i)$$

$$\sum_{j=1}^k \pi_j(\vec{x}_i) = 1$$

We need  $k$  different weight vectors

$$\vec{w}_1, \vec{w}_2, \dots, \vec{w}_k \in \mathbb{R}^{d+1}$$

$$\pi_1(\vec{x}_i) = f(\vec{w}_1^T \vec{x}_i)$$

$$\pi_2(\vec{x}_i) = f(\vec{w}_2^T \vec{x}_i)$$

⋮

$$\pi_k(\vec{x}_i) = f(\vec{w}_k^T \vec{x}_i)$$

$f = \text{sigmoid}$

$f$  has to be some function that sums to 1 with all values in  $[0, 1]$

$f \equiv \text{Softmax function}$

$$f(\vec{w}_j^T \vec{x}_i) = \frac{e^{\vec{w}_j^T \vec{x}_i}}{\sum_{a=1}^k e^{\vec{w}_a^T \vec{x}_i}}$$

←  $j$ -th class

← over all classes

$$= \text{softmax}(\vec{w}_j^T \vec{x}_i)$$

$$= \text{softmax}(\vec{w}_j^T \vec{x}_i)$$

$$= \text{softmax}(\vec{w}_j^T \vec{x}_i \mid \vec{w}_1^T \vec{x}_i, \vec{w}_2^T \vec{x}_i, \dots, \vec{w}_k^T \vec{x}_i)$$

Sigmoid

$$\begin{aligned} Q(\vec{w}^T \vec{x}_i) &= \frac{e^{\vec{w}^T \vec{x}_i}}{1 + e^{\vec{w}^T \vec{x}_i}} \\ &= \frac{e^{\vec{w}^T \vec{x}_i}}{e^0 + e^{\vec{w}^T \vec{x}_i}} \end{aligned}$$

Unknown parameters

$\vec{w}_1, \vec{w}_2, \dots, \vec{w}_k$  are all unknown

$$W = \begin{bmatrix} \vec{w}_1 & \vec{w}_2 & \dots & \vec{w}_k \\ | & | & & | \end{bmatrix} \quad (d+1) \times k$$

Weight matrix

$$\text{MLE: } \max_W \log L(W) = \sum_{i=1}^n \log P(y_i \mid W)$$

Under independence assumption

$$P(y_i \mid W) = \begin{cases} \pi_1(\vec{w}_1^T \vec{x}_i) & \text{if } y_i = c_1 \\ \pi_2(\vec{w}_2^T \vec{x}_i) & \text{if } y_i = c_2 \\ \vdots \\ \pi_k(\vec{w}_k^T \vec{x}_i) & \text{if } y_i = c_k \end{cases}$$

$\vec{e}_i$

$k=4$

$$\begin{aligned} \vec{e}_1 &= \begin{bmatrix} 1 & 0 & 0 & 0 \end{bmatrix} \\ \vec{e}_2 &= \begin{bmatrix} 0 & 1 & 0 & 0 \end{bmatrix} \\ \vec{e}_3 &= \begin{bmatrix} 0 & 0 & 1 & 0 \end{bmatrix} \end{aligned}$$

One-hot encoding of each class

One-hot encoding of each class

$\vec{e}_3 = \begin{bmatrix} 0 & 0 & 1 & 0 \end{bmatrix}$   
 $\vec{e}_4 = \begin{bmatrix} 0 & 0 & 0 & 1 \end{bmatrix}$

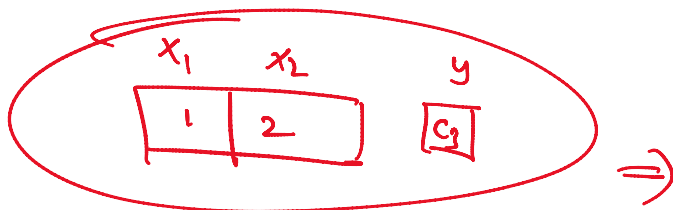
$$(\vec{x}_i, y_i) \rightarrow (\vec{x}_i, \vec{y}_i)$$

$$c_1 \rightarrow e_1$$

$$c_2 \rightarrow e_2$$

$$\vec{y}_i \in \{e_1, e_2, \dots, e_k\}$$

Encode the  $k$  classes into one-hot vectors for the two class.



$$\vec{x}_i = (1, 2) \quad \vec{y}_i = \vec{e}_3 = (0, 0, 1, 0)^T$$

$$\vec{x}_i \in \mathbb{R}^{d+1}$$

after augment

$$\vec{y}_i \in \mathbb{R}^k$$

one-hot

$$P(\vec{y}_i | \vec{x}_i) = \left[ \pi_1(\vec{w}_1^T \vec{x}_i)^{y_{i1}} \cdot \pi_2(\vec{w}_2^T \vec{x}_i)^{y_{i2}} \cdot \dots \cdot \pi_k(\vec{w}_k^T \vec{x}_i)^{y_{ik}} \right]$$

e.g.  $\vec{y}_i = \vec{e}_3 = (0, 0, 1, 0) = (y_{i1}, y_{i2}, y_{i3}, y_{i4})$

$k=4$

$$\log L(w) = \sum_{i=1}^n y_{i1} \log \pi_1(\vec{w}_1^T \vec{x}_i) + y_{i2} \log \pi_2(\vec{w}_2^T \vec{x}_i) + \dots$$

$$\log L(w) = \sum_{i=1}^n y_{i1} \log \pi_1(w_1^T x_i) + y_{i2} \log \pi_2(w_2^T x_i) + \dots + y_{ik} \log \pi_k(w_k^T x_i)$$

Cross-entropy = NLL  
( $\epsilon$ )

Solve for  $W = \begin{bmatrix} \vec{w}_1 & \vec{w}_2 & \dots & \vec{w}_k \end{bmatrix}$        $k$  different weight vectors

$$\frac{\partial \log L(w)}{\partial \vec{w}_1}, \quad \frac{\partial \log L(w)}{\partial \vec{w}_2}, \quad \dots, \quad \frac{\partial \log L(w)}{\partial \vec{w}_k}$$

$\nabla_{w_1} \quad \nabla_{w_2} \quad \nabla_{w_k}$

for the  $a$ -th class

$$\nabla_{w_a} = \sum_{i=1}^n \left( \underbrace{y_{ia}}_{\substack{\uparrow \\ \text{true} \\ \text{prob}}} - \underbrace{\pi_a(\vec{w}_a^T \vec{x}_i)}_{\substack{\text{pred} \\ \text{prob}}} \right) \cdot \vec{x}_i$$

Stochastic version

$$\nabla_{w_a}(\vec{x}_i) = (y_{ia} - \pi_a(w_a^T x_i)) x_i$$

## Gradient Ascent

$w_a^{(0)}$  = random  $(d+1)$  dim vector       $\forall a=1 \dots k$   
Uniform  $(0,1)$  or random normal  $(, \vec{0}, \sigma=0.1)$

Uniform  $(0,1)$  or random normal  $(, 0, \sigma=0.1)$

Epoch {  
for  $i = 1 \dots n$  in random order  $\leftarrow$  stochastic  
for  $j = 1 \dots K \leftarrow$  all classes  
 $\pi_j(\vec{w}_j^T \vec{x}_i) \leftarrow$  softmax values  
 $\nabla_{\vec{w}_j}(\vec{x}_i) = (y_{ij} - \pi_j(\vec{w}_j^T \vec{x}_i)) \cdot \vec{x}_i$   
 $\vec{w}_j^{(t+1)} = \vec{w}_j + \eta \cdot \nabla_{\vec{w}_j}(\vec{x}_i)$

Check for Convergence

$$\delta = \sum_{j=1}^K \left\| \underline{\vec{w}_j^{t+1}} - \underline{\vec{w}_j^t} \right\|$$

Stop if  $\delta < \epsilon$   
 $\nwarrow$  threshold, e.g.  $\epsilon = 0.0001$

Logistic model

$\rightarrow$  Computes log-odd ratios

e.g. binary case

$$\text{log odds ratio} = \log \left( \frac{P(y_i=1 | \vec{x}_i)}{P(y_i=0 | \vec{x}_i)} \right) = \vec{w}^T \vec{x}_i$$

or  
 $(\vec{w} - \vec{0})^T \vec{x}_i$   
 $\nwarrow$   
think of  
this as  
weight vector

Odds ratio

think of  
this as  
weight vector  
for class  
 $y=0$

also true for  $k$ -class situation (softmax)

fix  $k$ -th class

$$\log \text{ odds ratio}_1 \equiv \log \left( \frac{P(y_i = c_1 | x_i)}{P(y_i = c_k | x_i)} \right) = \cancel{\vec{w}_1^T x_i} \quad \text{mistake}$$

Should be  $(\vec{w}_1 - \vec{w}_k)^T x_i$

$$\log_2 = \log \left( \frac{P(y_i = c_2 | x_i)}{P(y_i = c_k | x_i)} \right) = \cancel{\vec{w}_2^T x_i} \quad \text{mistake}$$

should be  $(\vec{w}_2 - \vec{w}_k)^T x_i$

i.e. log odds ratio  
is related to the  
difference of the  
weight vectors