

# Lecture 7: Logistic Regression, Generalized Linear Models

## CS480/680 Intro to Machine Learning

2023-1-31

Pascal Poupart  
David R. Cheriton School of Computer Science





# Exponential Family

- More generally, when  $\Pr(\mathbf{x}|c_k)$  are members of the exponential family (e.g., Gaussian, exponential, Bernoulli, categorical, Poisson, Beta, Dirichlet, Gamma)

$$\Pr(\mathbf{x}|\boldsymbol{\theta}_k) = \exp\left(\boldsymbol{\theta}_k^T T(\mathbf{x}) - A(\boldsymbol{\theta}_k) + B(\mathbf{x})\right)$$

where  $\boldsymbol{\theta}_k$ : parameters of class  $k$

$T(\mathbf{x}), A(\boldsymbol{\theta}_k), B(\mathbf{x})$ : arbitrary functions of the inputs and parameters

- the posterior is a sigmoid logistic linear function in  $\mathbf{x}$

$$\Pr(c_k|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

# Probabilistic Discriminative Models

- Instead of learning  $\Pr(c_k)$  and  $\Pr(\mathbf{x}|c_k)$  by maximum likelihood and finding  $\Pr(c_k|\mathbf{x})$  by Bayesian inference, why not learn  $\Pr(c_k|\mathbf{x})$  directly by maximum likelihood?
- We know the general form of  $\Pr(c_k|\mathbf{x})$ :
  - **Logistic sigmoid** (binary classification)
  - **Softmax** (general classification)

# Logistic Regression

- Consider a single data point  $(\mathbf{x}, y)$ :  $\mathbf{w}^* = \operatorname{argmax}_{\mathbf{w}} \sigma(\mathbf{w}^T \bar{\mathbf{x}})^y (1 - \sigma(\mathbf{w}^T \bar{\mathbf{x}}))^{1-y}$
- Similarly, for an entire dataset  $(\mathbf{X}, \mathbf{y})$ :

$$\mathbf{w}^* = \operatorname{argmax}_{\mathbf{w}} \prod_n \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n)^{y_n} (1 - \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n))^{1-y_n}$$

Objective: negative log likelihood (minimization)

$$L(\mathbf{w}) = - \sum_n y_n \ln \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n) + (1 - y_n) \ln(1 - \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n))$$

Tip:  $\frac{\partial \sigma(a)}{\partial a} = \sigma(a)(1 - \sigma(a))$

# Logistic Regression

- NB: Despite the name, logistic regression is a form of classification.
- However, it can be viewed as regression where the goal is to estimate the posterior  $\Pr(c_k | \mathbf{x})$ , which is a continuous function

# Maximum likelihood

- Convex loss: set derivative to 0

$$0 = \frac{\partial L}{\partial \mathbf{w}} = - \sum_n y_n \frac{\cancel{\sigma(\mathbf{w}^T \bar{\mathbf{x}}_n)} (1 - \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n)) \bar{\mathbf{x}}_n}{\cancel{\sigma(\mathbf{w}^T \bar{\mathbf{x}}_n)}} - \sum_n (1 - y_n) \frac{\cancel{(1 - \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n))} \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n) (-\bar{\mathbf{x}}_n)}{\cancel{1 - \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n)}}$$

$$\Rightarrow 0 = - \sum_n y_n \bar{\mathbf{x}}_n - \sum_n y_n \cancel{\sigma(\mathbf{w}^T \bar{\mathbf{x}}_n)} \bar{\mathbf{x}}_n + \sum_n \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n) \bar{\mathbf{x}}_n + \sum_n y_n \cancel{\sigma(\mathbf{w}^T \bar{\mathbf{x}}_n)} \bar{\mathbf{x}}_n$$

$$\Rightarrow 0 = \sum_n [\sigma(\mathbf{w}^T \bar{\mathbf{x}}_n) - y_n] \bar{\mathbf{x}}_n$$

- Sigmoid prevents us from isolating  $\mathbf{w}$ , so we use an iterative method instead

# Gradient descent

- Iterative reweighted least square:

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla L(\mathbf{w})$$

where  $\nabla L$  is the gradient:  $\nabla L(\mathbf{w}) = \begin{bmatrix} \frac{\partial L}{\partial w_0} \\ \vdots \\ \frac{\partial L}{\partial w_m} \end{bmatrix}$

and  $\eta$  is the learning rate  
(scalar that determines the step length)



# Newton's method

- Iterative reweighted least square:

$$\mathbf{w} \leftarrow \mathbf{w} - \mathbf{H}^{-1} \nabla L(\mathbf{w})$$

where  $\nabla L$  is the gradient (column vector)

and  $H$  is the Hessian (matrix)

$$H = \begin{bmatrix} \frac{\partial L}{\partial^2 w_0} & \cdots & \frac{\partial L}{\partial w_0 \partial w_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial L}{\partial w_m \partial w_0} & \cdots & \frac{\partial L}{(\partial w_m)^2} \end{bmatrix}$$

# Hessian

$$\mathbf{H} = \nabla(\nabla L(\mathbf{w})) = \sum_{n=1}^N \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n) (1 - \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n)) \bar{\mathbf{x}}_n \bar{\mathbf{x}}_n^T = \bar{\mathbf{X}} \mathbf{R} \bar{\mathbf{X}}^T$$

$$\text{where } \mathbf{R} = \begin{bmatrix} \sigma_1(1 - \sigma_1) & & \\ & \ddots & \\ & & \sigma_N(1 - \sigma_N) \end{bmatrix}$$

$$\text{and } \sigma_1 = \sigma(\mathbf{w}^T \bar{\mathbf{x}}_1), \quad \sigma_N = \sigma(\mathbf{w}^T \bar{\mathbf{x}}_N)$$

# Case study

- Applications: recommender systems, ad placement
- Used by all major companies
- Advantages: logistic regression is **simple, flexible and efficient**

# App Recommendation

- Flexibility: millions of features (binary & numerical)

- Example:

- Efficiency: classification by dot products

Multiple classes:

$$c^* = \operatorname{argmax}_k \frac{\exp(\mathbf{w}_k^T \bar{\mathbf{x}})}{\sum_{k'} \exp(\mathbf{w}_{k'}^T \bar{\mathbf{x}})}$$
$$= \operatorname{argmax}_k \mathbf{w}_k^T \bar{\mathbf{x}}$$

Two classes:

$$c^* = \begin{cases} 1 & \sigma(\mathbf{w}^T \bar{\mathbf{x}}) \geq 0.5 \\ 0 & \text{otherwise.} \end{cases}$$
$$= \begin{cases} 1 & \mathbf{w}^T \bar{\mathbf{x}} \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

- Sparsity:

- Parallelization:

# Numerical Issues

- Logistic Regression is subject to overfitting
  - Without enough data, logistic regression can classify each data point arbitrarily well (i.e.,  $\Pr(\text{correct class}) \rightarrow 1$ )
- Problems:  $weights \rightarrow \pm\infty$   
Hessian  $\rightarrow$  singular
- Picture

# Regularization

- Solution: penalize large weights
- Objective:

$$\begin{aligned} \min_{\mathbf{w}} L(\mathbf{w}) + \frac{1}{2} \lambda \|\mathbf{w}\|_2^2 \\ = \min_{\mathbf{w}} - \sum_n y_n \ln \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n) + (1 - y_n) \ln(1 - \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n)) + \frac{1}{2} \lambda \mathbf{w}^T \mathbf{w} \end{aligned}$$

- Gradient:  $\nabla L(\mathbf{w}) = \bar{\mathbf{X}}(\sigma(\bar{\mathbf{X}}^T \mathbf{w}) - \mathbf{y}) + \lambda \mathbf{w}$
- Hessian:  $\mathbf{H} = \bar{\mathbf{X}} \mathbf{R} \bar{\mathbf{X}}^T + \lambda \mathbf{I}$

where  $R_{nn} = \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n)(1 - \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n))$

the term  $\lambda \mathbf{I}$  ensures that  $\mathbf{H}$  is not singular (eigenvalues  $\geq \lambda$ )

# Generalized Linear Models

- How can we do non-linear regression and classification while using the same machinery?
- Idea: map inputs to a different space and do linear regression/classification in that space

# Example

- Suppose the underlying function is quadratic



# Basis functions

- Use non-linear basis functions:

- Let  $\phi_i$  denote a basis function:  $\phi_0(x) = 1$

$$\phi_1(x) = x$$

$$\phi_2(x) = x^2$$

- Let the hypothesis space  $H$  be

$$H = \{x \rightarrow w_0\phi_0(x) + w_1\phi_1(x) + w_2\phi_2(x) \mid w_i \in \mathbb{R}\}$$

- If the basis functions are non-linear in  $x$ , then a non-linear hypothesis can still be found by linear regression

# Common basis functions

- Polynomial:  $\phi_j(x) = x^j$
- Gaussian:  $\phi_j(x) = e^{-\frac{(x-\mu_j)^2}{2s^2}}$
- Sigmoid:  $\phi_j(x) = \sigma\left(\frac{x-\mu_j}{s}\right)$  where  $\sigma(a) = \frac{1}{1+e^{-a}}$
- Also Fourier basis functions, wavelets, etc.

# Generalized Linear Models

- Linear regression:

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \sum_{n=1}^N (y_n - \mathbf{w}^T \bar{\mathbf{x}}_n)^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

- Generalized linear regression:

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \sum_{n=1}^N (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

- Linear separator (classification):

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} - \sum_n y_n \ln \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n) + (1 - y_n) \ln(1 - \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n)) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

- Generalized linear separator (classification):

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} - \sum_n y_n \ln \sigma(\mathbf{w}^T \phi(\mathbf{x}_n)) + (1 - y_n) \ln(1 - \sigma(\mathbf{w}^T \phi(\mathbf{x}_n))) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$