

CS485/685

Lecture 8: Jan 28, 2016

Classification by Logistic Regression,
Generalized linear models

[RN] Sec 18.6.4, [B] Sec. 4.3, [M] Chapt.
8, [HTF] Sec. 4.4

Beyond Mixtures of Gaussians

- Mixture of Gaussians:
 - Restrictive assumption: each class is Gaussian
 - Picture:
- Can we consider other distributions than Gaussians?

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, etc.)

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

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

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

- the posterior is a sigmoid logistic linear 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) :
$$\mathbf{w}^* = \operatorname{argmax}_{\mathbf{w}} \sigma(\mathbf{w}^T \bar{\mathbf{x}})$$

- 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

Newton's method

- Iterative reweighted least square:

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

where ∇L is the gradient (column vector)

and H is the Hessian (matrix)

$$H = \begin{bmatrix} \frac{\partial L}{\partial^2 w_1} & \cdots & \frac{\partial L}{\partial w_1 \partial w_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial L}{\partial w_m \partial w_1} & \cdots & \frac{\partial L}{(\partial w_m)^2} \end{bmatrix}$$

Hessian

$$\begin{aligned}\mathbf{H} &= \nabla(\nabla L(\mathbf{w})) \\ &= \sum_n \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n) (1 - \sigma(\mathbf{w}^T \bar{\mathbf{x}}_n)) \bar{\mathbf{x}}_n \bar{\mathbf{x}}_n^T \\ &= \bar{\mathbf{X}} \mathbf{R} \bar{\mathbf{X}}^T\end{aligned}$$

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

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

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)
 - Examples:

- Efficiency: classification by dot products

$$\begin{aligned}c^* &= \operatorname{argmax}_k \sigma(\mathbf{w}_k^T \bar{\mathbf{x}}) \\ &= \operatorname{argmax}_k \mathbf{w}_k^T \bar{\mathbf{x}}\end{aligned}$$

- Sparsity:
- Parallelization:

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 ϕ_i denote a basis function

$$\phi_0(x) = 1$$

$$\phi_1(x) = x$$

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

- Let the hypothesis space H be

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

- 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.

Non-linear classification

- More generally, if $\Pr(\mathbf{x}|c_k)$ is not from the exponential family, map \mathbf{x} to a feature space defined by a set of basis functions ϕ_i where $\Pr(\boldsymbol{\phi}(\mathbf{x})|c_k)$ is from the exponential family

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

- In general we apply the sigmoid to a non-linear combination of the inputs