

Lecture 18: Logistic Regression, Neural Networks

Instructor: Prof. Ganesh Ramakrishnan

MAP estimation and regularized LR

① **FROM** MAP for LR: $\tilde{\mathbf{w}} = \operatorname{argmax}_{\mathbf{w}} \Pr(\mathbf{w}) L(\mathcal{D}, \mathbf{w})$

$$= \operatorname{argmax}_{\mathbf{w}} \frac{1}{\left(\frac{2\pi}{\lambda}\right)^{\frac{m}{2}}} e^{-\frac{\lambda}{2} \|\mathbf{w}\|_2^2} \prod_{i=1}^m \left(f_{\mathbf{w}}(\mathbf{x}^{(i)}) \right)^{y^{(i)}} \left(1 - f_{\mathbf{w}}(\mathbf{x}^{(i)}) \right)^{1-y^{(i)}}$$

.....Taking $-\frac{1}{m} \log(\cdot)$ transformation,

② **TO** Min of the Regularized Logistic (Cross-Entropy) Loss function:

$$\tilde{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} - \left[\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \log f_{\mathbf{w}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log \left(1 - f_{\mathbf{w}}(\mathbf{x}^{(i)}) \right) \right) \right] + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \quad (1)$$

where we have ignored $-\frac{1}{m} \log \left(\left(\frac{2\pi}{\lambda} \right)^{\frac{m}{2}} \right)$ since this term is independent of \mathbf{w} .

.....Thus, MAP $\tilde{\mathbf{w}}$ can be found by minimizing the *Regularized Cross Entropy Error*

Gradient descent for Regularized LR

- 1 The final descent update

$$-\eta \nabla E(\mathbf{w}) = \eta \left[\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} - f_{\mathbf{w}}(\mathbf{x}^{(i)}) \right) \phi(\mathbf{x}^{(i)}) - \lambda \mathbf{w} \right] \quad (2)$$

- 2 The iterative update rule:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^k + \eta \left[\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} - f_{\mathbf{w}^k}(\mathbf{x}^{(i)}) \right) \phi(\mathbf{x}^{(i)}) - \lambda \mathbf{w}^k \right] \quad (3)$$

- 3 Stochastic version of the same:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^k + \eta \left(y^{(i)} - f_{\mathbf{w}^k}(\mathbf{x}^{(i)}) \right) \phi(\mathbf{x}^{(i)}) - \eta \lambda \mathbf{w}^k \quad (4)$$

One Extension to multi-class logistic

- ① Each class $c = 1, 2, \dots, K - 1$ can have a different weight vector $[\mathbf{w}_{c,1}, \mathbf{w}_{c,2}, \dots, \mathbf{w}_{c,k}, \dots, \mathbf{w}_{c,p}]$ and

$$p(Y = c | \phi(\mathbf{x})) = \frac{e^{-(\mathbf{w}_c)^T \phi(\mathbf{x})}}{1 + \sum_{k=1}^{K-1} e^{-(\mathbf{w}_k)^T \phi(\mathbf{x})}}$$

for $c = 1, \dots, K - 1$ so that

$$p(Y = K | \phi(\mathbf{x})) = \frac{1}{1 + \sum_{k=1}^{K-1} e^{-(\mathbf{w}_k)^T \phi(\mathbf{x})}}$$

Alternative (equivalent) extension to multi-class logistic

- ① Each class $c = 1, 2, \dots, K$ can have a different weight vector $[\mathbf{w}_{c,1}, \mathbf{w}_{c,2} \dots \mathbf{w}_{c,p}]$ and

$$p(Y = c | \phi(\mathbf{x})) = \frac{e^{-(\mathbf{w}_c)^T \phi(\mathbf{x})}}{\sum_{k=1}^K e^{-(\mathbf{w}_k)^T \phi(\mathbf{x})}}$$

for $c = 1, \dots, K$.

Logistic Regression Kernelized

- 1 We have already seen (a) Cross Entropy loss and (b) Bayesian interpretation for regularization
- 2 The Regularized (Logistic) Cross-Entropy Loss function (minimized wrt $\mathbf{w} \in \mathbb{R}^p$):

$$E(\mathbf{w}) = - \left[\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \log f_{\mathbf{w}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - f_{\mathbf{w}}(\mathbf{x}^{(i)})) \right) \right] + \frac{\lambda}{2m} \|\mathbf{w}\|_2^2 \quad (5)$$

- 3 Equivalent dual kernelized objective¹
(minimized wrt $\alpha \in \mathbb{R}^m$):

¹http://perso.telecom-paristech.fr/~clemenco/Projets_ENPC_files/kernel-log-regression-svm-boosting.pdf

Logistic Regression Kernelized

- 1 We have already seen (a) Cross Entropy loss and (b) Bayesian interpretation for regularization
- 2 The Regularized (Logistic) Cross-Entropy Loss function (minimized wrt $\mathbf{w} \in \mathbb{R}^p$):

$$E(\mathbf{w}) = - \left[\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \log f_{\mathbf{w}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - f_{\mathbf{w}}(\mathbf{x}^{(i)})) \right) \right] + \frac{\lambda}{2m} \|\mathbf{w}\|_2^2 \quad (5)$$

- 3 Equivalent dual kernelized objective¹
(minimized wrt $\alpha \in \mathbb{R}^m$):

$$E_D(\alpha) = \left[\sum_{i=1}^m \left(\sum_{j=1}^m -y^{(i)} K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \alpha_j + \frac{\lambda}{2} \alpha_i K(\mathbf{x}^{(i)}, \mathbf{x}^{(i)}) \alpha_i \right) + \log \left(1 + \sum_{j=1}^m \alpha_j K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \right) \right] \quad (6)$$

$$\text{Decision function } f_{\mathbf{w}}(\mathbf{x}) = \frac{1}{1 + \sum_{j=1}^m \alpha_j K(\mathbf{x}, \mathbf{x}^{(j)})}$$

¹http://perso.telecom-paristech.fr/~clemenco/Projets_ENPC_files/kernel-log-regression-svm-boosting.pdf

Some Tutorial 7 and 8 Questions

- Prove that the Kernelized Logistic Regression form is equivalent to the original Logistic Regression minimum regularized cross entropy form: 2 Hints
- Contrast the Kernelized Logistic Regression with the dual of Support Vector Classifier

Logistic Regression Generalized to CRF

- 1 Multi-class LR: $c \in [1 \dots K]$ has weight vector $[w_{c,1} \dots w_{c,p}]$

$$\Pr(y = c \mid \mathbf{x}) = \frac{e^{-w_c^T \phi(\mathbf{x})}}{\sum_{k=1}^K e^{-w_k^T \phi(\mathbf{x})}} =$$

²<http://www.tzi.de/~edelkamp/lectures/ml/scripts/loglinearcrrfs.pdf>

Logistic Regression Generalized to CRF

- 1 Multi-class LR: $c \in [1 \dots K]$ has weight vector $[w_{c,1} \dots w_{c,p}]$

$$\Pr(y = c \mid \mathbf{x}) = \frac{e^{-w_c^T \phi(\mathbf{x})}}{\sum_{k=1}^K e^{-w_k^T \phi(\mathbf{x})}} = \frac{e^{-\tilde{w}^T \phi(\mathbf{x}, y=c)}}{Z(\mathbf{x}, \tilde{w})}$$

where $\tilde{w} = [w_{1,1} \dots w_{1,p}, \dots, w_{c,1} \dots w_{c,p}, \dots, w_{K,1} \dots w_{K,p}]$ and

$\phi(\mathbf{x}, y) = [\delta(y, 1)\phi(\mathbf{x}), \dots, \delta(y, c)\phi(\mathbf{x}), \dots, \delta(y, K)\phi(\mathbf{x})]$ and $\delta(a, b) = 1$ if $a = b$ and $= 0$ otherwise

- 2 Extended to non-iid inference in Conditional Random Fields² with $\mathbf{x} = [\mathbf{x}^{(1)} \dots \mathbf{x}^{(n)}]$ and $\mathbf{y} = [y^{(1)} \dots y^{(n)}]$:

²<http://www.tzi.de/~edelkamp/lectures/ml/scripts/loglinearcrfs.pdf>

Logistic Regression Generalized to CRF

- 1 Multi-class LR: $c \in [1 \dots K]$ has weight vector $[w_{c,1} \dots w_{c,p}]$

$$\Pr(y = c \mid \mathbf{x}) = \frac{e^{-w_c^T \phi(\mathbf{x})}}{\sum_{k=1}^K e^{-w_k^T \phi(\mathbf{x})}} = \frac{e^{-\tilde{w}^T \phi(\mathbf{x}, y=c)}}{Z(\mathbf{x}, \tilde{w})}$$

where $\tilde{w} = [w_{1,1} \dots w_{1,p}, \dots, w_{c,1} \dots w_{c,p}, \dots, w_{K,1} \dots w_{K,p}]$ and

$\phi(\mathbf{x}, y) = [\delta(y, 1)\phi(\mathbf{x}), \dots, \delta(y, c)\phi(\mathbf{x}), \dots, \delta(y, K)\phi(\mathbf{x})]$ and $\delta(a, b) = 1$ if $a = b$ and $= 0$ otherwise

- 2 Extended to non-iid inference in Conditional Random Fields² with $\mathbf{x} = [\mathbf{x}^{(1)} \dots \mathbf{x}^{(n)}]$ and $\mathbf{y} = [y^{(1)} \dots y^{(n)}]$:

$$\Pr(\mathbf{y} \mid \mathbf{x}) = \frac{e^{-\mathbf{w}^T \phi(\mathbf{x}, \mathbf{y})}}{Z(\mathbf{x}, \mathbf{w})}$$

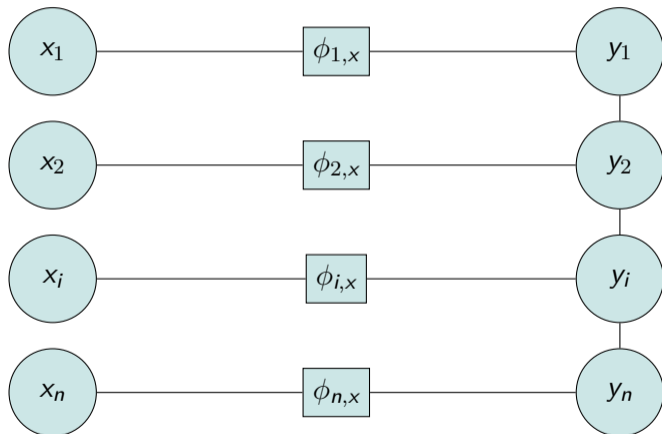
²<http://www.tzi.de/~edelkamp/lectures/ml/scripts/loglinearcrfs.pdf>

Conditional Random Fields (Linear)³

inputs

x -potentials

classes & y -potentials $\phi_{i,y}$



³CRF is an instance of Graphical Models (detailed notes at <https://www.cse.iitb.ac.in/~cs725/notes/classNotes/misc/CaseStudyWithProbabilisticModels.pdf>)

<https://www.cse.iitb.ac.in/~cs725/notes/classNotes/misc/CaseStudyWithProbabilisticModels.pdf>

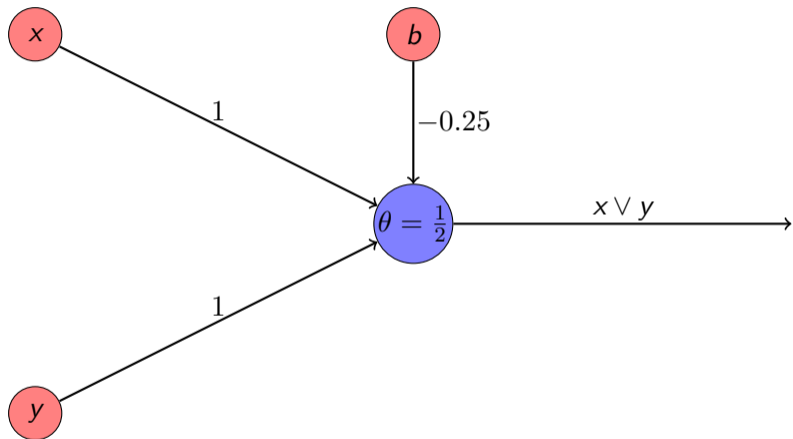


Non-linear perceptron?

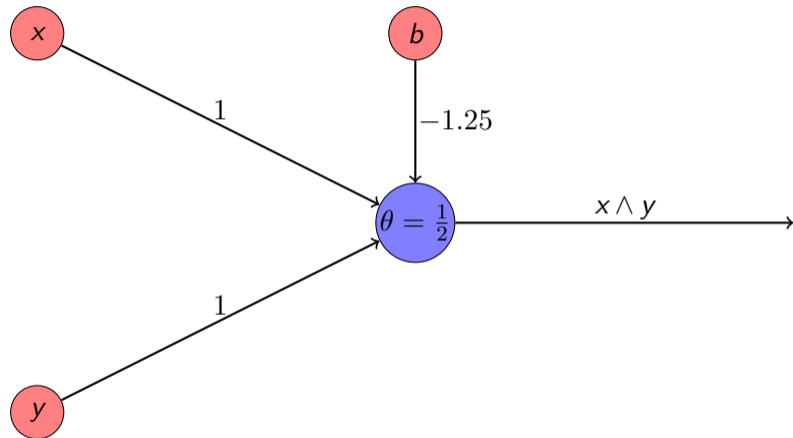
- Kernelized perceptron: $f(x) = \text{sign} \left(\sum_i \alpha_i y_i K(x, x_i) + b \right)$
 - ▶ INITIALIZE: $\alpha = \text{zeroes}()$
 - ▶ REPEAT: for $\langle x_i, y_i \rangle$
 - ★ If $\text{sign} \left(\sum_j \alpha_j y_j K(x_j, x_i) + b \right) \neq y_i$
 - ★ then, $\alpha_i = \alpha_i + 1$
 - ★ endif
- Neural Networks: Cascade of layers of perceptrons giving you non-linearity
 - ▶ $\text{sign} \left((w^*)^T \phi(x) \right)$ replaced by $g \left((w^*)^T \phi(x) \right)$ where $g(s)$ is a
 - 1 step function: $g(s) = 1$ if $s \in [\theta, \infty)$ and $g(s) = 0$ otherwise OR
 - 2 sigmoid function: $g(s) = \frac{1}{1+e^{-s}}$
 - 3 Rectified Linear Unit (ReLU): $g(s) = \max(0, s)$: Most popular activation function
 - 4 Softplus: $g(s) = \ln(1 + e^s)$

Options 2, 3 and 4 have the thresholding step deferred. Threshold changes as bias is changed.

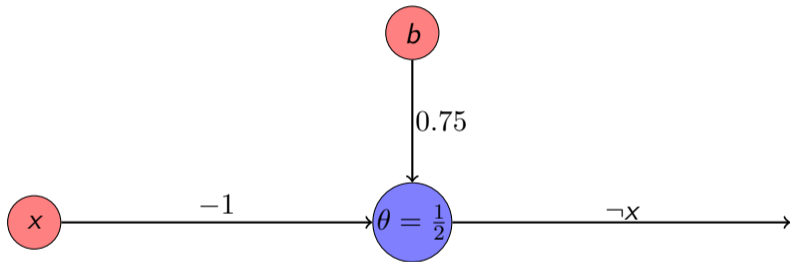
OR using (step) perceptron



AND using (step) perceptron

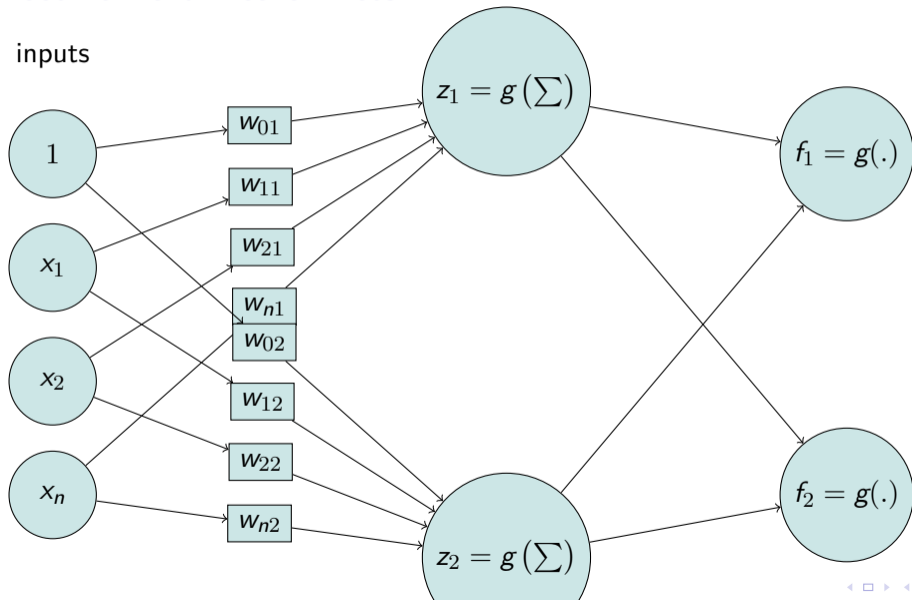


NOT using perceptron



Feed-forward Neural Nets

inputs



Eg: Feed-forward Neural Net for XOR ($\theta = 0$)

inputs

