# Lecture 18: Logistic Regression, Neural Networks

Instructor: Prof. Ganesh Ramakrishnan

## MAP estimation and regularized LR

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

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

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

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

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

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

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

*shrinkage.*

2. The iterative update rule:

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

3. Stochastic version of the same:

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

# One Extension to multi-class logistic

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

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

$$= \frac{e^{-\omega^T \phi(x)}}{1 + e^{-\omega^T \phi(x)}} \quad (c=1)$$

For $k=2$

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

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

$$= \frac{1}{1 + e^{-\omega^T \phi(x)}} \quad (c=2)$$

special case: $k=2$

$\omega_k = 0$

# Alternative (equivalent) extension to multi-class logistic

(More intuitive)

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

For weights $\omega$ in second form, $\exists \, \tilde{\omega}$ in first form with same distr & vice versa.

Tut 1: show equivalence of these two formulations!

Hint: In the second formulation divide each numerator & denominator by $e^{-\omega_k^T \phi(x)}$

Formulation ② to ①

$$\tilde{w}_c = w_c - w_k$$

Formulation ① to ②

$$w_c = \tilde{w}_c \text{ for } c = 1 .. K-1$$

$$w_K = 0$$

} Assume $w_c$ for ② & $\tilde{w}_c$ for ①

**Tut 7:** Are there other ways of generalizing Logistic Regression to multiple classes?

# 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 $\mathbf{w} \in \Re^p$):

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

3. Equivalent dual kernelized objective[1] (minimized wrt $\alpha \in \Re^m$):

$$E(\omega) = \frac{-1}{m} \left[ \sum_{i=1}^{m} y^{(i)} \omega^T \phi(x^{(i)}) - \log\left(1 + e^{\omega^T \phi(x^{(i)})}\right) \right]$$

---

[1] http://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 $\mathbf{w} \in \Re^p$):

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

3. Equivalent dual kernelized objective[1] (minimized wrt $\alpha \in \Re^m$):

$$w^T\phi(x_i) = \sum_j \alpha_j k(x_i, x_j)$$

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

Decision function $f_{\mathbf{w}}(\mathbf{x}) = \dfrac{1}{1 + \sum_{j=1}^{m}\alpha_j K\left(\mathbf{x}, \mathbf{x}^{(j)}\right)} = P(y=1 \mid x)$

---

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

# Some Tutorial 7 and 8 Questions

① $E(\omega) = - \sum_{i=1}^{m} y^{(i)} \omega^T \phi(x^{(i)}) + \log \left( 1 + e^{\omega^T \phi(x^{(i)})} \right)$

② Matrix multiplication
equality

- 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: $c \in [1 \ldots K]$ has weight vector $[w_{c,1} \ldots w_{c,p}]$

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

$$\begin{bmatrix} \omega_1 \\ \omega_2 \\ \omega_c \\ \omega_k \end{bmatrix} \underbrace{\quad}_{k+p} \rightarrow \omega^T = [\underbrace{\omega_{1,1} \cdots \omega_{1,p}}\ \ \omega_{c,1} - \ - \ \omega_{c,p} \ - \cdots \ \omega_{k,1} - \cdot \omega_{k,p}]$$

$$kp \text{ dim vector}$$

$$\tilde{\phi}(x, c) \in [\quad \bigcirc \quad \phi(x) \quad \bigcirc \quad]$$

# Logistic Regression Generalized to CRF

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

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

where $\tilde{w} = [w_{1,1} \ldots w_{1,p}, \ldots w_{c,1} \ldots w_{c,p}, \ldots w_{K,1} \ldots w_{K,p}]$ and
$\tilde{\phi}(\mathbf{x}, y) = [\underbrace{\delta(y, 1)\phi(\mathbf{x})}, \ldots, \underbrace{\delta(y, c)\phi(\mathbf{x})} \ldots \underbrace{\delta(y, K)\phi(\mathbf{x})}]$ and $\delta(a, b) = 1$ if $a = b$ and $= 0$
otherwise
   *(handwritten: If $y \neq 1$, will be 0)*   *(handwritten: $\phi(\alpha)$)*   *(handwritten: 0)*

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

---

[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 $[w_{c,1} \dots w_{c,p}]$

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

*(handwritten:)* ① structure over y can be used directly ② shared params

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 Fields[2] with $\mathbf{x} = [\mathbf{x}^{(1)} \dots \mathbf{x}^{(n)}]$ and $\mathbf{y} = [y^{(1)} \dots y^{(n)}]$:

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

*(handwritten annotations:)* [AD, VB, UN]  [Today, is, Thursday]  $w_i \rightarrow -\infty$ [discouraging]  $\phi_i(-, [-AD, JJ]..) \leq 1$  [JJ, JJ, NN, VB, VB]  [Thursday, afternoon, lecture, is, boring]

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

# Conditional Random Fields (Linear)[3]

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



$$\phi_{c_1 c_2}(y_1, y_2)$$

$$\phi_{c_1, c_2}(y_2, y_3)$$
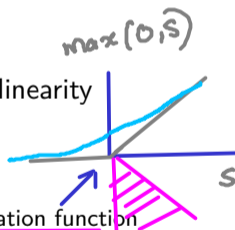
Eg: Hierarchical classification

---

# Non-linear perceptron?

- Kernelized perceptron: $f(x) = sign \left( \sum_i \alpha_i y_i K(x, x_i) + b \right)$

  - INITIALIZE: $\alpha$=zeroes()
  - REPEAT: for $< x_i, y_i >$
    - ★ If $sign \left( \sum_j \alpha_j y_j K(x_i, x_j) + b \right) \neq y_i$
    - ★ then, $\alpha_i = \alpha_i + 1$
    - ★ endif

- Neural Networks: Cascade of layers of perceptrons giving you non-linearity
  - $sign \left( (w^*)^T \phi(x) \right)$ 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: $g(s) = \ln \left( 1 + e^s \right)$

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

OR using (step) perceptron $(x \lor y)$
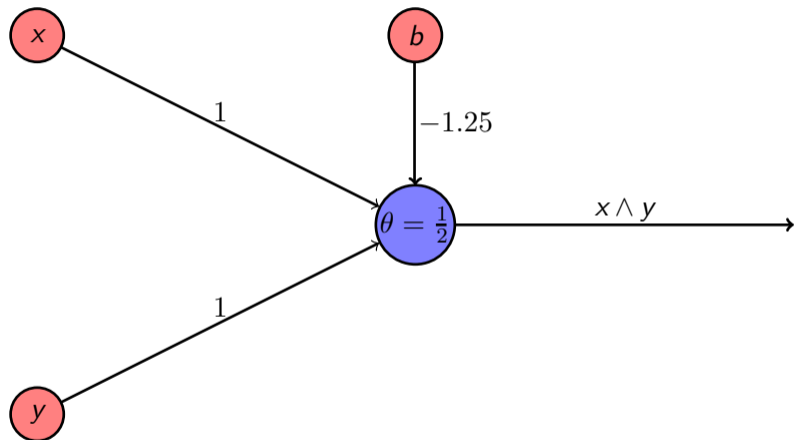


$\theta \in (-0.25, 0.75)$

$x w_1 + y w_2 + b$

$x=1, y=1 \Rightarrow 1.75 > \theta \Rightarrow \text{o/p} = 1$

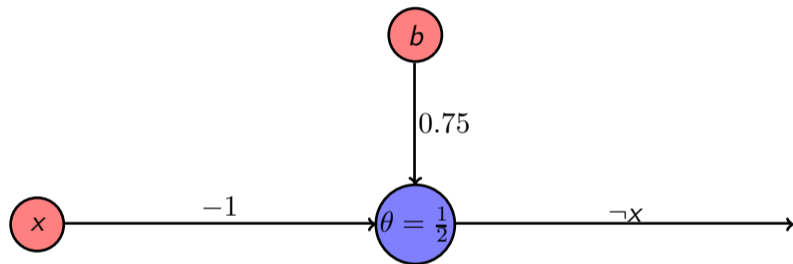$x=0, y=1 \underline{or} x=1, y=0 \Rightarrow 0.75 > \theta \Rightarrow \text{o/p} = 1$

$x=0, y=0 \Rightarrow -0.25 < \theta \Rightarrow \text{o/p} = 0$

# AND using (step) perceptron

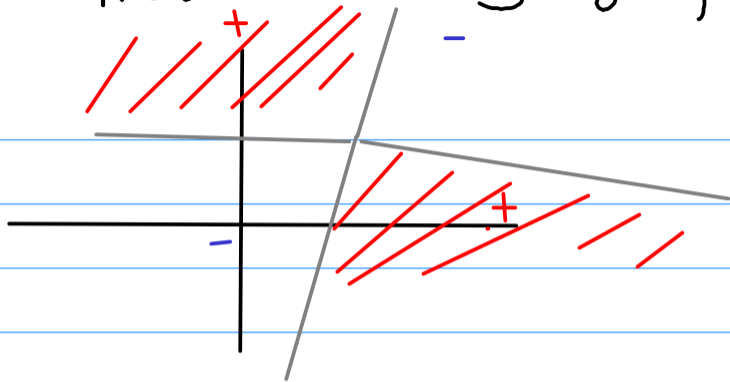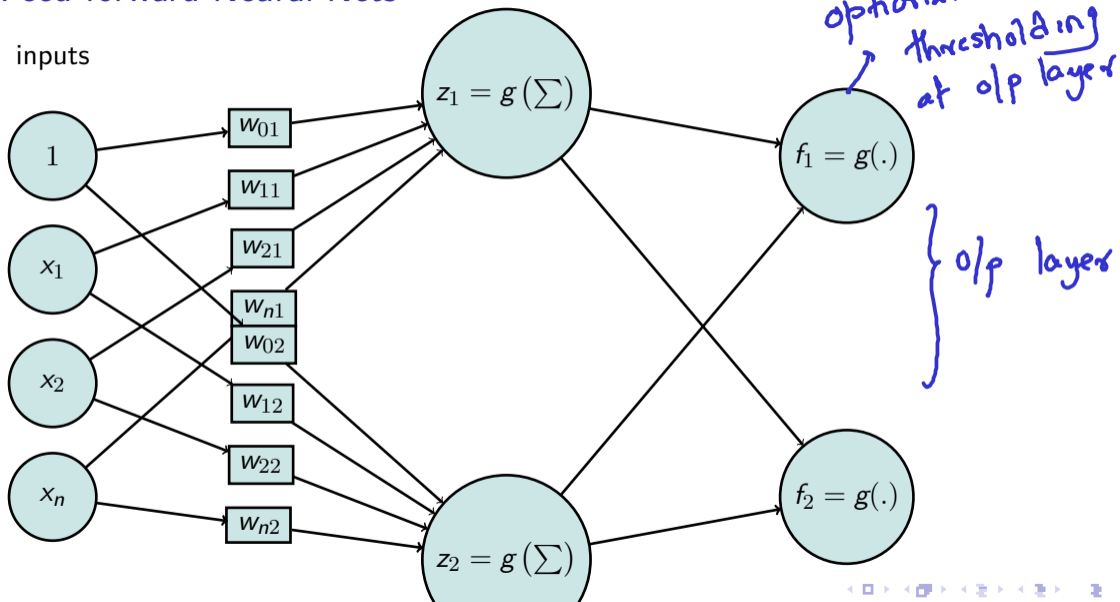$(x \wedge y)$

# NOT using perceptron

How about XOR using single perceptron



No linear seperator possible!
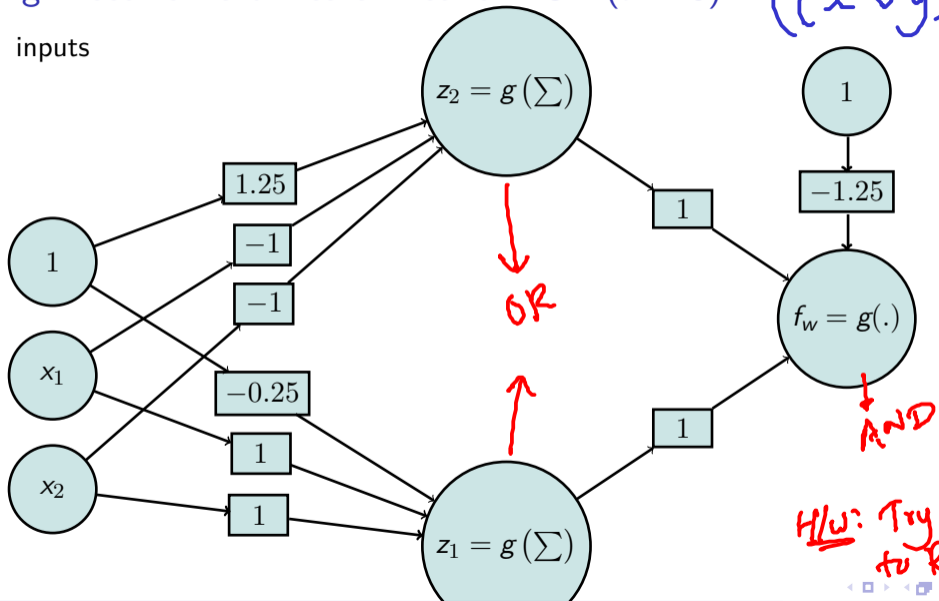
# Feed-forward Neural Nets

inputs



optional thresholding at o/p layer

o/p layer

$z_1 = g\left(\sum\right)$

$z_2 = g\left(\sum\right)$

$f_1 = g(.)$

$f_2 = g(.)$

$1$

$x_1$

$x_2$

$x_n$

$w_{01}$
$w_{11}$
$w_{21}$
$w_{n1}$
$w_{02}$
$w_{12}$
$w_{22}$
$w_{n2}$

Eg: Feed-forward Neural Net for XOR ($\theta = 0$)   $((x \vee y) \wedge (\neg x \vee \neg y))$

inputs

safe1        education-1
cs224minoriit

user=passwd=anon to login to SAFE

cs725:419 is the quizid