

Lecture 19: Logistic Regression and Regularization, Kernelized Logistic Regression, Neural Networks

Instructor: Prof. Ganesh Ramakrishnan

Recap: Minimizing negative Log-likelihood for LR

- ① Cross-entropy¹ is the average number of bits needed to identify an event (example \mathbf{x}) drawn from the (data) set \mathcal{D} , if a coding scheme is used that is optimized for a modeled probability distribution $\Pr(y|\mathbf{w}, \phi(.))$, rather than the 'true' distribution $\Pr(y|\mathcal{D})$.

$$H(P(\mathcal{D})) \leq E(\mathbf{w}) = \mathbf{E}_{\Pr(y|\mathcal{D})} [-\log \Pr(y|\mathbf{w}, \phi(.))]$$

$$f_{\mathbf{w}}(\mathbf{x}^{(i)}) = P(y=1|\mathbf{x}^{(i)}) \\ = \frac{1}{1 + \exp(-\mathbf{w}^T \phi(\mathbf{x}^{(i)}))}$$

- ② The Cross-entropy Loss function:

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

with some simplification,

(slides have gradient descent in this form)

Hint: Useful for Tut 7

$$E(\mathbf{w}) = - \left[\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \mathbf{w}^T \phi(\mathbf{x}^{(i)}) - \log (1 + \exp(\mathbf{w}^T \mathbf{x}^{(i)})) \right) \right] \quad (3)$$

(used for grad descent in last class)

¹https://en.wikipedia.org/wiki/Cross_entropy

Recap: Gradient descent for LR

- 1 No closed form solution to the cross-entropy loss

$$\hat{\mathbf{w}}^{MLE} = \arg \min_{\mathbf{w}} - \left[\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \log f_{\mathbf{w}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - f_{\mathbf{w}}(\mathbf{x}^{(i)})) \right) \right] \quad (4)$$

- 2 Apply gradient descent with $\mathbf{w}^{(k+1)} = \mathbf{w}^k - \eta \nabla E(\mathbf{w}^k)$
- 3 The descent update

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

- 4 $\nabla f_{\mathbf{w}}(\mathbf{x}^{(i)}) = \phi(\mathbf{x}^{(i)}) \left(\frac{e^{-(\mathbf{w})^T \phi(\mathbf{x}^{(i)})}}{1 + e^{-(\mathbf{w})^T \phi(\mathbf{x}^{(i)})}} \right)$
 \Rightarrow

- 5 $\nabla \log f_{\mathbf{w}}(\mathbf{x}^{(i)}) = \phi(\mathbf{x}^{(i)}) e^{-(\mathbf{w})^T \phi(\mathbf{x}^{(i)})} \left(\frac{1}{1 + e^{-(\mathbf{w})^T \phi(\mathbf{x}^{(i)})}} \right)^2$ and
 $\nabla \log (1 - f_{\mathbf{w}}(\mathbf{x}^{(i)})) = -\phi(\mathbf{x}^{(i)}) \left(\frac{1}{1 + e^{-(\mathbf{w})^T \phi(\mathbf{x}^{(i)})}} \right)^2$

Recap: Descent update for LR

$$-\eta \nabla E(\mathbf{w}) = -\eta \left[\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \nabla \log f_{\mathbf{w}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \nabla \log (1 - f_{\mathbf{w}}(\mathbf{x}^{(i)})) \right) \right] \quad (6)$$

① $\nabla \log f_{\mathbf{w}}(\mathbf{x}^{(i)}) = \phi(\mathbf{x}^{(i)}) e^{-(\mathbf{w})^T \phi(\mathbf{x}^{(i)})} \left(\frac{1}{1 + e^{-(\mathbf{w})^T \phi(\mathbf{x}^{(i)})}} \right)^2$ and

$$\nabla \log (1 - f_{\mathbf{w}}(\mathbf{x}^{(i)})) = -\phi(\mathbf{x}^{(i)}) \left(\frac{1}{1 + e^{-(\mathbf{w})^T \phi(\mathbf{x}^{(i)})}} \right)^2$$

② \Rightarrow The final descent update is

$$-\eta \nabla E(\mathbf{w}) = \eta \left[\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} - f_{\mathbf{w}}(\mathbf{x}^{(i)}) \right) \phi(\mathbf{x}^{(i)}) \right] \quad (7)$$

Recap: Gradient descent for LR

- 1 The final descent update

LBFAS etc approximate $(\nabla^2 E)^{-1}$

$$-\eta \nabla E(\mathbf{w}) = \eta \left[\frac{1}{m} \sum_{i=1}^m (y^{(i)} - f_{\mathbf{w}}(\mathbf{x}^{(i)})) \phi(\mathbf{x}^{(i)}) \right] \quad (8)$$

Perceptron $\sum_{i \in M} y^{(i)} \phi(\mathbf{x}^{(i)})$ (Weka has both implementations)

- 2 The iterative update rule:

Newton update (faster to converge)

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

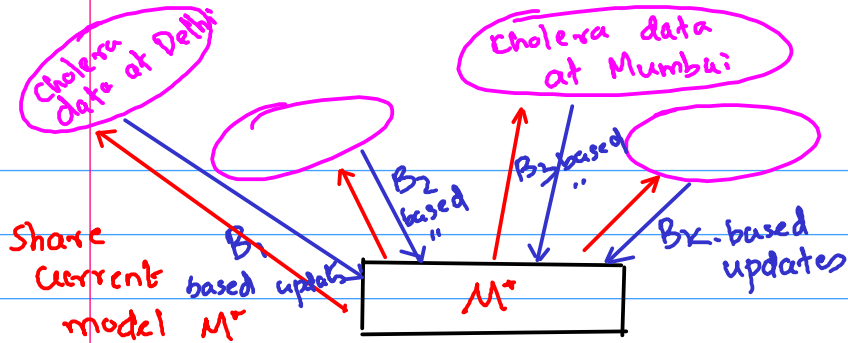
Accounting for extent of misclassification! Deterministically over all pts

- 3 Stochastic version of the same:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \eta (y^{(i)} - f_{\mathbf{w}^k}(\mathbf{x}^{(i)})) \phi(\mathbf{x}^{(i)}) \quad (10)$$

Batch stochastic: $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \eta \frac{1}{|B|} \sum_B \dots$

- 4 How would you contrast the updates with sigmoid (LR) against those with the step function (perceptron)?



Mahout : Distributed machine learning platform. (performs batch stochastic for LR)

Sigmoid (LR) vs. step function (perceptron)

- ① Stochastic update for step fn (perceptron) with $y^{(i)} \in \{-1, 1\}$: Pick any example $(\mathbf{x}^{(i)}, y^{(i)})$, for which $\text{sign} \left((\mathbf{w}^{(k)})^T \phi(\mathbf{x}^{(i)}) \right) \neq y^{(i)}$.

$$\mathbf{w}^{(k+1)} = \mathbf{w}^k + \eta y^{(i)} \phi(\mathbf{x}^{(i)}) \quad (11)$$

- ② Stochastic update for sigmoid fn (LR) with $y^{(i)} \in \{0, 1\}$: Pick any example $(\mathbf{x}^{(i)}, y^{(i)})$, for which $|f_{\mathbf{w}^k}(\mathbf{x}^{(i)}) - y^{(i)}| > 0.5$.

$$\mathbf{w}^{(k+1)} = \mathbf{w}^k + \eta \left(y^{(i)} - f_{\mathbf{w}^k}(\mathbf{x}^{(i)}) \right) \phi(\mathbf{x}^{(i)}) \quad (12)$$

- ③ Recall: (12) is also the stochastic update for linear regression! (12) is a characteristic update for **generalized linear models**² of which perceptron, linear regression and logistic are special cases.

²https://en.wikipedia.org/wiki/Generalized_linear_model

Direct contrast with perceptron

soft update

Prob of error > 0.5 minimizes risk (Bayes risk)

$g(\mathbf{w}^T \phi(\mathbf{x}) + b)$ for nonparam models g is a fn of x

Regularized LR and its Probabilistic Interpretation

- 1 The Regularized (Logistic) Cross-Entropy Loss function:

$$E(\mathbf{w}) = - \left[\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \log f_{\mathbf{w}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - f_{\mathbf{w}}(\mathbf{x}^{(i)})) \right) \right] + \frac{\lambda}{2m} \|\mathbf{w}\|_2^2 \quad (13)$$

- 2 Motivations: Avoiding overfitting by discouraging large values of w_j for every j .
- 3 Probabilistic Explanation?

ie Bayesian interpretation?

$$P(w_i) = \mathcal{N}(0, \frac{1}{\lambda})$$

Regularized LR and its Probabilistic Interpretation

- 1 The Regularized (Logistic) Cross-Entropy Loss function:

$$E(\mathbf{w}) = - \left[\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \log f_{\mathbf{w}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - f_{\mathbf{w}}(\mathbf{x}^{(i)})) \right) \right] + \frac{\lambda}{2m} \|\mathbf{w}\|_2^2 \quad (13)$$

- 2 Motivations: Avoiding overfitting by discouraging large values of w_j for every j .
- 3 Probabilistic Explanation? A Bayesian Posterior probabilistic explanation to regularized LR (next)
- 4 We will reinvoke Bayesian (Parameter) Estimation

Bayesian Inference For Logistic Regression

MAP Estimation and regularized LR

$$-\log p(\mathbf{w}) = \underbrace{+\log \left(\frac{2\pi}{\lambda} \right)^{m/2}}_{\text{can be ignored (independent of } \mathbf{w})} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

- ① Recall the multivariate Gaussian (Normal) Distribution:

$$\mathcal{N}(\mathbf{w}; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{m}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(\mathbf{w}-\mu)^T \Sigma^{-1}(\mathbf{w}-\mu)} \text{ when } \Sigma \in \mathbb{R}^{m \times m} \text{ is positive-definite and } \mu \in \mathbb{R}^m$$

- ② Suppose we want each $|w_i|$ to be bounded roughly by $\pm \frac{3}{\lambda}$

- ③ Then by the 3- σ rule we let $\mathbf{w} \sim \mathcal{N}(\mathbf{w}; 0, \frac{1}{\lambda} I)$ where I is an $m \times m$ identity matrix

④ $\Rightarrow \Pr(\mathbf{w}) = \frac{1}{\left(\frac{2\pi}{\lambda}\right)^{\frac{m}{2}}} e^{-\frac{\lambda}{2} \|\mathbf{w}\|_2^2}$

Now derive MAP estimate for \mathbf{w}

$$\begin{aligned} \mathbf{w}_{\text{MAP}} &= \underset{\mathbf{w}}{\operatorname{argmax}} \Pr(\mathbf{w}|\mathcal{D}) = \underset{\mathbf{w}}{\operatorname{argmax}} \frac{\Pr(\mathbf{w}) \Pr(\mathcal{D}|\mathbf{w})}{\Pr(\mathcal{D})} = \underset{\mathbf{w}}{\operatorname{argmax}} \Pr(\mathbf{w}) L(\mathbf{w}; \mathcal{D}) \\ &= \underset{\mathbf{w}}{\operatorname{argmax}} \log \Pr(\mathbf{w}) + \log(L(\mathbf{w}; \mathcal{D})) = \underset{\mathbf{w}}{\operatorname{argmin}} \left(-\log \Pr(\mathbf{w}) - \log(L(\mathbf{w}; \mathcal{D})) \right) \end{aligned}$$

MAP estimation and regularized LR

① $\Pr(\mathbf{w}) = \frac{1}{\left(\frac{2\pi}{\lambda}\right)^{\frac{m}{2}}} e^{-\frac{\lambda}{2} \|\mathbf{w}\|_2^2}$

② Recall the MLE for LR: $\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}} L(\mathcal{D}; \mathbf{w})$

$$= \underset{\mathbf{w}}{\operatorname{argmax}} \prod_{i=1}^m \left(f_{\mathbf{w}}(\mathbf{x}^{(i)}) \right)^{y^{(i)}} \left(1 - f_{\mathbf{w}}(\mathbf{x}^{(i)}) \right)^{1-y^{(i)}}$$

③ Now the MAP for LR: $\tilde{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}} \Pr(\mathbf{w}) L(\mathcal{D}; \mathbf{w}) =$

$$\underset{\mathbf{w}}{\operatorname{argmax}} \frac{1}{m} \left(\log \Pr(\mathbf{w}) + \log L(\mathcal{D}, \mathbf{w}) \right)$$

MAP estimation and regularized LR

① $\Pr(\mathbf{w}) = \frac{1}{\left(\frac{2\pi}{\lambda}\right)^{\frac{m}{2}}} e^{-\frac{\lambda}{2}\|\mathbf{w}\|_2^2}$

② Recall the MLE for LR: $\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}} L(\mathcal{D}; \mathbf{w})$

$$= \underset{\mathbf{w}}{\operatorname{argmax}} \prod_{i=1}^m \left(f_{\mathbf{w}}(\mathbf{x}^{(i)}) \right)^{y^{(i)}} \left(1 - f_{\mathbf{w}}(\mathbf{x}^{(i)}) \right)^{1-y^{(i)}}$$

③ Now the MAP for LR: $\tilde{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}} \Pr(\mathbf{w}) L(\mathcal{D}; \mathbf{w}) =$

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

MAP estimation and regularized LR

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

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

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

② TO Min of the Regularized cross entropy

$$= \underset{\mathbf{w}}{\operatorname{argmin}} -\frac{1}{m} \sum y^{(i)} \log f_{\mathbf{w}}(\mathbf{x}^{(i)}) + (1-y^{(i)}) \log(1-f_{\mathbf{w}}(\mathbf{x}^{(i)})) + \frac{\lambda}{2m} \|\mathbf{w}\|_2^2$$

MAP estimation and regularized LR

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

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

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

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

$$\tilde{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} - \left[\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \log f_{\mathbf{w}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - f_{\mathbf{w}}(\mathbf{x}^{(i)})) \right) \right] + \underbrace{\frac{\lambda}{2} \|\mathbf{w}\|_2^2}_{\text{red wavy underline}} \quad (14)$$

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 update:

Gradient descent for Regularized LR

$$w^{(k+1)} = w^{(k)} + \frac{\eta}{n} \sum_i (y^{(i)} - f_w(x^{(i)})) \phi(x^{(i)}) - \frac{\lambda \eta}{n} w^{(k)}$$

Gradient descent for Regularized LR

- ① The final descent update

$$-\eta \nabla E(\mathbf{w}) = \eta \left[\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} - f_{\mathbf{w}}(\mathbf{x}^{(i)}) \right) \phi(\mathbf{x}^{(i)}) - \lambda \mathbf{w} \right] \quad (15)$$

- ② The iterative update rule:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^k + \eta \left[\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} - f_{\mathbf{w}^k}(\mathbf{x}^{(i)}) \right) \phi(\mathbf{x}^{(i)}) - \lambda \mathbf{w}^k \right] \quad (16)$$

- ③ Stochastic version of the same:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^k + \eta \left(y^{(i)} - f_{\mathbf{w}^k}(\mathbf{x}^{(i)}) \right) \phi(\mathbf{x}^{(i)}) - \eta \lambda \mathbf{w}^k \quad (17)$$

Inspires other methods of update that work well in practice

Extension to multi-class logistic

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

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

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

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

Assuming:

Classification based on
 $\arg\max_k P_Y(Y=k | \phi(\mathbf{x}))$

- ① Why not $\mathbf{w}_1, \dots, \mathbf{w}_K$ (K vectors?)
② What is geometric interpretation?

Reserved for K^{th} class: case we used $K=1$

Alternative (equivalent) extension to multi-class logistic

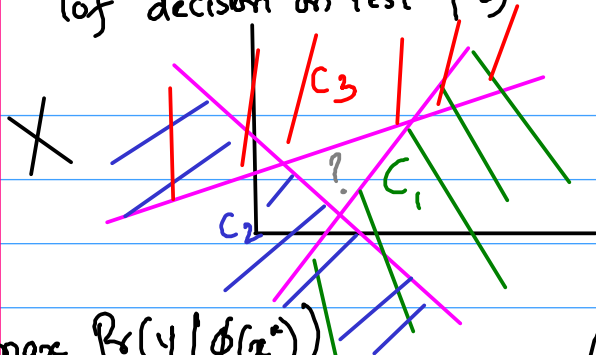
- ① Each class $c = 1, 2, \dots, K$ can have a different weight vector $[\mathbf{w}_{c,1}, \mathbf{w}_{c,2} \dots \mathbf{w}_{c,p}]$ and

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

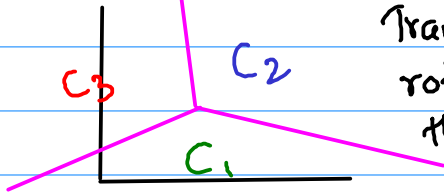
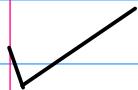
for $c = 1, \dots, K$.

Ans to Q1: The two are equivalent (Tut 7). For grad descent in the second formulation I update K vectors & \therefore more numerical error than updating $K-1$ vectors which are "differences" of \mathbf{w}_c with \mathbf{w}_K (Tut 7)

Geometric interpretation?
(of decision on test pt)



$$y^* = \underset{y}{\operatorname{argmax}} P(y | \phi(x^*))$$



Gradient descent:
Translation (for b) and
rotation of limbs of
the starfish!

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 $\mathbf{w} \in \mathbb{R}^p$):

$$E(\mathbf{w}) = - \left[\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \log f_{\mathbf{w}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - f_{\mathbf{w}}(\mathbf{x}^{(i)})) \right) \right] + \frac{\lambda}{2m} \|\mathbf{w}\|_2^2 \quad (18)$$

- 3 Equivalent dual kernelized objective³ (minimized wrt $\alpha \in \mathbb{R}^m$):

$$\sum_i L(\mathbf{w}, \mathbf{x}^{(i)}, y^{(i)})$$

$$\|\mathbf{w}\|_2^2 = \langle \mathbf{w}, \mathbf{w} \rangle$$

$$E(\mathbf{w}) = -\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \log \left(\frac{1}{1 + \exp \left(- \sum_j \alpha_j k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \right)} \right) + (1 - y^{(i)}) \log \left(\dots \right) \right)$$

$$f(\mathbf{x}^i) = \frac{1}{1 + \sum_j \alpha_j k(\mathbf{x}, \mathbf{x}^j)}$$

$$\mathbf{w}^\top \phi(\mathbf{x}) \Rightarrow \sum_j \alpha_j k(\mathbf{x}, \mathbf{x}^{(j)})$$

$$+ \lambda / 2m \|\mathbf{w}\|^2 ?$$

³Representer Theorem and 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 $\mathbf{w} \in \mathbb{R}^p$):

$$E(\mathbf{w}) = - \left[\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \log f_{\mathbf{w}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - f_{\mathbf{w}}(\mathbf{x}^{(i)})) \right) \right] + \frac{\lambda}{2m} \|\mathbf{w}\|_2^2 \quad (18)$$

- 3 Equivalent dual kernelized objective³
(minimized wrt $\alpha \in \mathbb{R}^m$):

$$E_D(\alpha) = \left[\sum_{i=1}^m \left(\sum_{j=1}^m -y^{(i)} K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \alpha_j + \frac{\lambda}{2} \alpha_i K(\mathbf{x}^{(i)}, \mathbf{x}^{(i)}) \alpha_i \right) + \log \left(1 + \sum_{j=1}^m \alpha_j K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \right) \right] \quad (19)$$

How?

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

Handwritten note: $\mathbf{w}^T \phi(\mathbf{x})$

³Representer Theorem and http://perso.telecom-paristech.fr/~clemenco/Projets_ENPC_files/kernel-log-regression-svm-boosting.pdf

Some Tutorial 7 Questions

- Prove that the Kernelized Logistic Regression form is equivalent to the original Logistic Regression minimum regularized cross entropy form: 2 Hints
- Show equivalence of the two formulations of Multiclass Logistic Regression.

Logistic Regression Generalized to CRF

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

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

① Construct \tilde{w}
& $\tilde{\phi}(\cdot)$

① Made \tilde{w} independent of class c
② Expanded ϕ to $\tilde{\phi}(y,x)$ & take class as argument

Useful when y has some structure / periodicity etc!

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 Fields⁴ with $\mathbf{x} = [\mathbf{x}^{(1)} \dots \mathbf{x}^{(n)}]$ and $\mathbf{y} = [y^{(1)} \dots y^{(n)}]$:

⁴ <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 $[w_{c,1} \dots w_{c,p}]$

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

where $\tilde{\mathbf{w}} = [w_{1,1} \dots w_{1,p}, \dots, w_{c,1} \dots w_{c,p}, \dots, w_{K,1} \dots w_{K,p}]$ 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 Fields⁴ with $\mathbf{x} = [\mathbf{x}^{(1)} \dots \mathbf{x}^{(n)}]$ and $\mathbf{y} = [y^{(1)} \dots y^{(n)}]$:

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

⁴ <http://www.tzi.de/~edelkamp/lectures/ml/scripts/loglinearcrfs.pdf>