

Lecture 17: Logistic Regression contd.

Instructor: Prof. Ganesh Ramakrishnan

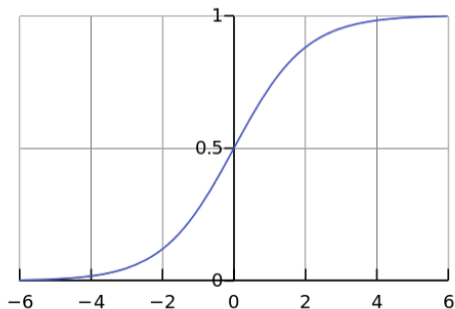
Sigmoidal (perceptron) Classifier

- 1 **(Binary) Logistic Regression**, abbreviated as **LR** is a single node perceptron-like classifier, but with....

▶ $\text{sign}((\mathbf{w}^*)^T \phi(\mathbf{x}))$ replaced by $g((\mathbf{w}^*)^T \phi(\mathbf{x}))$ where $g(s)$ is sigmoid function: $g(s) = \frac{1}{1+e^{-s}}$

- 2 $f_{\mathbf{w}}(\mathbf{x}) = g((\mathbf{w}^*)^T \phi(\mathbf{x})) = \frac{1}{1+e^{-(\mathbf{w}^*)^T \phi(\mathbf{x})}} \in [0, 1]$ can be interpreted as $Pr(y = 1|\mathbf{x})$

▶ Then $Pr(y = 0|\mathbf{x}) = 1 - f_{\mathbf{w}}(\mathbf{x})$



Logistic Regression: The Sigmoidal (perceptron) Classifier

- ① Estimator $\hat{\mathbf{w}}$ is a function of the dataset

$$\mathcal{D} = \left\{ (\phi(\mathbf{x}^{(1)}), y^{(1)}), (\phi(\mathbf{x}^{(2)}), y^{(2)}), \dots, (\phi(\mathbf{x}^{(m)}), y^{(m)}) \right\}$$

- ▶ Estimator $\hat{\mathbf{w}}$ is meant to approximate the parameter \mathbf{w} .
- ② Maximum Likelihood Estimator: Estimator $\hat{\mathbf{w}}$ that maximizes the likelihood $L(\mathcal{D}; \mathbf{w})$ of the data \mathcal{D} .
- ▶ Assumes that all the instances $(\phi(\mathbf{x}^{(1)}), y^{(1)}), (\phi(\mathbf{x}^{(2)}), y^{(2)}), \dots, (\phi(\mathbf{x}^{(m)}), y^{(m)})$ in \mathcal{D} are all independent and identically distributed (iid)
 - ▶ Thus, Likelihood is the probability of \mathcal{D} under iid assumption: $\hat{\mathbf{w}} = \max_{\mathbf{w}} L(\mathcal{D}, \mathbf{w}) =$

Logistic Regression: The Sigmoidal (perceptron) Classifier

- ① Estimator $\hat{\mathbf{w}}$ is a function of the dataset

$$\mathcal{D} = \left\{ (\phi(\mathbf{x}^{(1)}), y^{(1)}), (\phi(\mathbf{x}^{(2)}), y^{(2)}), \dots, (\phi(\mathbf{x}^{(m)}), y^{(m)}) \right\}$$

- ▶ Estimator $\hat{\mathbf{w}}$ is meant to approximate the parameter \mathbf{w} .

- ② Maximum Likelihood Estimator: Estimator $\hat{\mathbf{w}}$ that maximizes the likelihood $L(\mathcal{D}; \mathbf{w})$ of the data \mathcal{D} .

- ▶ Assumes that all the instances $(\phi(\mathbf{x}^{(1)}), y^{(1)}), (\phi(\mathbf{x}^{(2)}), y^{(2)}), \dots, (\phi(\mathbf{x}^{(m)}), y^{(m)})$ in \mathcal{D} are all independent and identically distributed (iid)
- ▶ Thus, Likelihood is the probability of \mathcal{D} under iid assumption: $\hat{\mathbf{w}} = \max_{\mathbf{w}} L(\mathcal{D}, \mathbf{w}) =$

$$\operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^m p(y^{(i)} | \phi(\mathbf{x}^{(i)})) = \operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^m \left(\frac{1}{1 + e^{-(\mathbf{w})^T \phi(\mathbf{x}^{(i)})}} \right)^{y^{(i)}} \left(\frac{e^{-(\mathbf{w})^T \phi(\mathbf{x}^{(i)})}}{1 + e^{-(\mathbf{w})^T \phi(\mathbf{x}^{(i)})}} \right)^{1 - y^{(i)}}$$

if $y^{(i)} = 1$ then $P(y=1|x^{(i)})$

if $y^{(i)} = 0$ then $P(y=0|x^{(i)})$

Training LR

Recall Entropy

$$-\sum p_i \log(p_i)$$

- 1 Thus, Maximum Likelihood Estimator for \mathbf{w} is

$$\hat{\mathbf{w}} = \operatorname{argmax}_{\mathbf{w}} L(\mathcal{D}, \mathbf{w}) = \operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^m p(y^{(i)} | \phi(\mathbf{x}^{(i)}))$$

$$= \operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^m \left(\frac{1}{1 + e^{-\mathbf{w}^T \phi(\mathbf{x}^{(i)})}} \right)^{y^{(i)}} \left(\frac{e^{-\mathbf{w}^T \phi(\mathbf{x}^{(i)})}}{1 + e^{-\mathbf{w}^T \phi(\mathbf{x}^{(i)})}} \right)^{1-y^{(i)}}$$

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

$$\left. \begin{aligned} &-\frac{1}{m} \sum_{i=1}^m y^{(i)} \log f_{\mathbf{w}}(\mathbf{x}^{(i)}) \\ &+ (1-y^{(i)}) \log(1-f_{\mathbf{w}}(\mathbf{x}^{(i)})) \end{aligned} \right\}$$

- 2 Maximizing the likelihood $\Pr(\mathcal{D}; \mathbf{w})$ w.r.t \mathbf{w} , is the same as minimizing the negative log-likelihood $E(\mathbf{w}) = -\frac{1}{m} \log \Pr(\mathcal{D}; \mathbf{w})$ w.r.t \mathbf{w} .

- ▶ Derive the expression for $E(\mathbf{w})$.
- ▶ $E(\mathbf{w})$ is called the cross-entropy loss function

Minimizing negative Log-likelihood for LR

1 The Cross-entropy Loss function:

$$E(w) = - \frac{1}{m} \sum_{i=1}^m P_{Y_0}(Y=1|x^{(i)}) \log(P_{Y_M}(Y=1|x^{(i)}))$$

then this is low & this $-\log(\cdot)$ is high

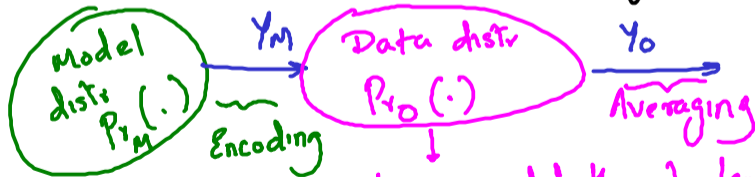
If $P_{Y_0} \approx P_{Y_M}$
then $E(w)$ will
be low

$$+ P_{Y_0}(Y=0|x^{(i)}) \log(P_{Y_M}(Y=0|x^{(i)}))$$

If this prob is also high

If this prob is close to 1

$-\log(\cdot)$ low



Viewing model through lens of data

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

$$-\frac{1}{m} \sum_i y^i \log(f_{\omega}(x^i)) + (1-y^i) \log(1-f_{\omega}(x^i))$$

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

⇓

$$-\frac{1}{m} \sum_i y^i \left(\log(\cancel{1}) - \log(1 + e^{-\omega^T \phi(x^i)}) \right)$$

$$+ (1-y^i) \left(\log(\underbrace{e^{-\omega^T \phi(x^i)}}_{-\omega^T \phi(x^i)}) - \log(1 + e^{-\omega^T \phi(x^i)}) \right)$$

$$= -\frac{1}{m} \sum_i y_i \omega^T \phi(x^i) - \log \left((1 + e^{-\omega^T \phi(x^i)}) e^{\omega^T \phi(x^i)} \right)$$

Minimizing negative Log-likelihood for LR

- ① The Cross-entropy Loss function:

"Observed/Empirical distribution"
 $P_{r_0}(y=1|x^{(i)})$

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

Model distribution
 $P_{r_M}(y=1|x^{(i)})$

or with some simplification,

Like unsigned distance error of perceptron

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

- ② 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 $\Pr(y|\mathbf{w}, \phi(\cdot))$, rather than the 'true' distribution $\Pr(y|\mathcal{D})$.

$$E(\mathbf{w}) = \mathbf{E}_{\Pr(y|\mathcal{D})} \left[-\log \Pr(y|\mathbf{w}, \phi(\cdot)) \right] \quad (3)$$

empirical distr

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

Gradient descent for LR

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

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

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

Handwritten notes:

$$f_{\mathbf{w}}(\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \phi(\mathbf{x})}}$$
$$\log(f_{\mathbf{w}}(\mathbf{x})) = -\log(1 + e^{-\mathbf{w}^T \phi(\mathbf{x})})$$
$$\nabla E(\mathbf{w}^k) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \nabla_{\mathbf{w}^k} (\log f_{\mathbf{w}}(\mathbf{x}^{(i)})) + (1 - y^{(i)}) \nabla_{\mathbf{w}^k} (\log(1 - f_{\mathbf{w}}(\mathbf{x}^{(i)})))$$
$$\nabla_{\mathbf{w}^k} (\log f_{\mathbf{w}}(\mathbf{x}^{(i)})) = \left(\frac{\phi(\mathbf{x}) e^{-\mathbf{w}^T \phi(\mathbf{x})}}{1 + e^{-\mathbf{w}^T \phi(\mathbf{x})}} \right) = (1 - f_{\mathbf{w}^k}(\mathbf{x}^{(i)})) \phi(\mathbf{x}^{(i)})$$

Gradient descent for LR

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

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

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

- 3 The descent update

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

Gradient descent for LR

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

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

- 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

Gradient descent for LR

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

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

- 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

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

Descent update for LR

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

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

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

2 \Rightarrow The final descent update is

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

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

*misclassification cost of
 $x^{(i)}$ higher if $y^{(i)} = 1$
& $f_{\mathbf{w}}(\mathbf{x}^{(i)}) \rightarrow 0$
and vice versa*

*\hookrightarrow like
perceptron
but weighted*

Gradient descent for LR

- 1 The final descent update

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

- 2 The iterative update rule: (Weka package implements gradient & also LBFGS)

$$\mathbf{w}^{(k+1)} = \mathbf{w}^k + \eta \left[\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} - f_{\mathbf{w}^k}(\mathbf{x}^{(i)}) \right) \phi(\mathbf{x}^{(i)}) \right] \quad (9)$$

- 3 Stochastic version of the same: (Mahout implements batch stochastic on each data batch B_j)

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

$$\mathbf{w}^{(k+1)} = \mathbf{w}^k + \eta \frac{1}{|B_j|} \sum_{\mathbf{x}^{(i)} \in B_j} \left(y^{(i)} - f_{\mathbf{w}^k}(\mathbf{x}^{(i)}) \right) \phi(\mathbf{x}^{(i)})$$

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

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 $\text{sign}\left(\left(\mathbf{w}^{(k)}\right)^T \phi\left(\mathbf{x}^{(i)}\right)\right) \neq y^{(i)}$.

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

penalty based on sign mismatch

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

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

$\omega^{(k+1)} = \omega^k + \eta \sum (y^{(i)} - (\phi^T(x^i) \omega^k + b)) \phi(x^i)$

penalty is based on extent of digression

- 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.

$\hookrightarrow g(\omega^T \phi(x) + b)$

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

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?

$$Pr(y=1 | \phi(\mathbf{x})) = \frac{1}{1 + e^{-\mathbf{w}^T \phi(\mathbf{x})}}$$

↓
To help generalize well
to new data points

Need: $Pr(\mathbf{w}) \dots$ s.t. $Pr(\mathbf{w} | \mathcal{D}) \propto \underbrace{Pr(\mathcal{D} | \mathbf{w})}_{\text{has similar form as } Pr(\mathbf{w})} Pr(\mathbf{w}) = \underbrace{L(\mathbf{w} | \mathcal{D})}_{\text{has similar form as } Pr(\mathbf{w})} Pr(\mathbf{w})$

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

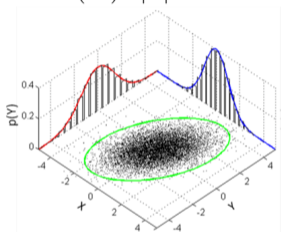
$$P_Y(\mathbf{w}) = \mathcal{N}\left(0, \frac{1}{\lambda} \mathbf{I}\right) \Leftrightarrow \underbrace{w_i \in \left[-\frac{3}{\sqrt{\lambda}}, \frac{3}{\sqrt{\lambda}}\right]}_{\text{roughly}}$$

Bayesian Inference For Logistic Regression

MAP Estimation and regularized LR

- 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 $\Rightarrow \Pr(\mathbf{w}) = \frac{1}{(\frac{2\pi}{\lambda})^{\frac{m}{2}}} e^{-\frac{\lambda}{2} \|\mathbf{w}\|_2^2}$. Write expression for $\Pr(D|\mathbf{w})$ & then derive $\Pr(\mathbf{w}|D) \propto \Pr(D|\mathbf{w}) \Pr(\mathbf{w})$

$$\log(\Pr(\mathbf{w}|D)) = \log(\Pr(D|\mathbf{w})) + \log(\Pr(\mathbf{w})) = (\text{cross entropy}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$$

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: $\tilde{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}} \Pr(\mathbf{w})L(\mathcal{D}; \mathbf{w}) =$

$$\underset{\mathbf{w}}{\operatorname{argmax}} \log \Pr(\mathbf{w}) + \log L(\mathcal{D}, \mathbf{w})$$

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: $\tilde{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}} \Pr(\mathbf{w})L(\mathcal{D}; \mathbf{w}) =$

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

$= \underset{\mathbf{w}}{\operatorname{argmax}} \underbrace{\log\left(e^{-\lambda/2\|\mathbf{w}\|_2^2}\right)}_{\text{can be ignored}} - \underbrace{m \operatorname{Cross Entropy}}_{\text{can be ignored}} = \underset{\mathbf{w}}{\operatorname{argmax}} -m \operatorname{Cross Entropy} - \frac{\lambda}{2}\|\mathbf{w}\|_2^2$

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:

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

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*

Tut 7: Posterior $\Pr(\mathbf{w}|\mathcal{D})$ is Gaussian!

Gradient descent for Regularized LR

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:

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

shrinking

- 3 Stochastic version of the same:

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

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

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

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 $c = 1, \dots, K$.