

CS229 Problem Set #2

Supervised learning II

1) Logistic Regression: Training Stability

a) The model didn't converge on the dataset B.

b) The model provided tries to find a linear decision boundary with mapping the features to a higher dimension.

It does so by calculating the functional margins of the data $\hat{y}^{(i)} = y^{(i)}(\theta^T x^{(i)}) + b$.

For each correct classification $\hat{y}^{(i)} > 0$

Also, in the implementation provided the margin is inversely proportional to the gradient.

$$\begin{cases} \text{prob} = Y * X.\text{dot}(\theta) \\ \text{grad} = -(1/m) * (X.^T.\text{dot}(\text{prob} * Y)) \end{cases}$$

Since the dataset B is not linearly separable θ keeps changing considerably (the convergence condition over θ is not attained) after each epoch.

c) i) The convergence condition in the provided model is that $\text{prev_theta} - \theta < 1e^{-15}$

Setting a learning rate $\alpha \ll 1e^{-15}$ could cause the model to converge even after the first update. But the θ is not optimal and this is not the goal of the model.

ii) Decreasing the learning rate by $1/t^2$

seems to be a reasonable solution, which makes the algorithm to converge.

But it is not guaranteed that we have an optimal estimation θ .

iii) We can notice that the dataset is already linearly scaled.

iv) Adding L_2 regularization would

would reduce $\|\theta\|^2$ and thus ~~reduce~~ increase the margins in misclassification

2. Model Calibration

A binary classification model is said to be well calibrated given a range (a, b) if

$$* 0 \leq a < b \leq 1$$

$$* \sum_{i \in I_{a,b}} P(y^{(i)} = 1 | x^{(i)}, \theta) = \frac{\sum_{i \in I_{a,b}} \mathbb{I}\{y^{(i)} = 1\}}{|I_{a,b}|}$$

where $I_{a,b} = \{i \mid i \in \{1, \dots, m\}, P(y^{(i)} = 1 | x^{(i)}, \theta) \in (a, b)\}$,

and $|S|$ denotes the size of the set S .

a) * $\theta \in \mathbb{R}^{n+1}$ be the maximum likelihood parameters learned after training a logistic regression

$$\hat{x}_0^{(i)} = 1$$

Problem: Prove that the above property holds true for the described logistic regression model over the range $(a, b) = (0, 1)$

θ are maximum likelihood parameters

$$\Rightarrow \frac{\partial \ell_\theta}{\partial \theta_j} = \sum_{i=1}^m \left(y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \right)$$

$$\Rightarrow \frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)} = 0$$

We know that $x_0^{(i)} = 1$

$$\Rightarrow \text{When } j=0 : \frac{1}{m} \sum_{i=1}^m y^{(i)} = \frac{1}{m} \sum_{i=1}^m h_\theta(x^{(i)})$$

b) Problem

- * Does well-calibrated model have perfect accuracy?
- * Does perfect accuracy make a model well-calibrated?

Let's prove them wrong with a counterexample

- Assumptions: * $a = 0.6$; $b = 0.7$

* all examples are positive

$$\Rightarrow \sum_{i \in I_{a,b}} \mathbb{I}_{\{y^{(i)} = 1\}} = 1$$

$$\frac{|\{i \in I_{a,b} | y^{(i)} = 1\}|}{|\{i \in I_{a,b}\}|}$$

$$\Rightarrow \sum_{i \in I_{a,b}} \frac{P(y^{(i)} = 1 | x^{(i)}, \theta)}{|\{i \in I_{a,b}\}|} \leq 0.7$$

Even though the model is perfectly accurate, it is not well calibrated.

- Assumptions:
 - * $a = 0, b = 0.7$

* Well calibrated model

$$\Rightarrow P(y^{(i)}=1 | x^{(i)}; \theta) \geq 0.5 \quad \forall i \in I_{a,b}$$

\Rightarrow All predictions are positive ①

But: Since the model is well calibrated: $\sum_{i \in I_{a,b}} y^{(i)} = 1 \} \leq 0.7$

\Rightarrow Not all true labels are positive ②

① + ② \Rightarrow Well calibrated models don't necessarily have perfect accuracy.

c) What effect does L2 regularization in logistic regression have on model calibration

$$l(\theta) = \sum_{i=0}^m \left(y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \right) - \frac{\lambda}{2} \|\theta\|^2$$

With MLE parameters:

$$\Rightarrow \frac{\partial l(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x_j - \lambda \theta_j = 0$$

When $j = 0$

$$\Rightarrow \frac{\partial l(\theta)}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) - \lambda \theta_0 = 0$$

$$\Rightarrow \frac{1}{m} \sum_{i=1}^m h_\theta(x^{(i)}) + \lambda \theta_0 = \frac{1}{m} \sum_{i=1}^m y^{(i)}$$

\Rightarrow The model is not well-calibrated.

3 Bayesian Interpretation of Regularization

a) * $\Theta_{\text{map}} = \arg \max_{\Theta} p(\Theta | x, y)$

* $\Theta_{\text{MLE}} = \arg \max_{\Theta} p(y | x, \Theta)$

* We assume $p(\Theta) = p(\Theta | x)$

Problem: Show that $\Theta_{\text{map}} = \arg \max_{\Theta} p(y | x, \Theta) p(\Theta)$

$$\Rightarrow \Theta_{\text{map}} = \arg \max_{\Theta} \frac{p(y | x, \Theta)}{p(x, y)}$$

$$= \arg \max_{\Theta} \frac{p(y | x, \Theta) p(\Theta | x)}{p(x, y)}$$

$$= \arg \max_{\Theta} \frac{p(y | x, \Theta) \cdot p(\Theta | x) p(x)}{p(x) \cdot p(y | x)}$$

$$= \arg \max_{\Theta} \frac{p(y | x, \Theta) p(\Theta)}{p(y | x)}$$

$$= \arg \max_{\Theta} p(y | x, \Theta) p(\Theta)$$

$$b) * \theta \sim \mathcal{N}(\theta, \eta^2 I)$$

$$* \hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmax}} p(y|x, \theta) p(\theta)$$

Problem: Show that

$$\hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmin}} -\log p(y|x, \theta) + \lambda \|\theta\|_2^2.$$

and find the value of λ ?

$$\hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmax}} p(y|x, \theta) p(\theta)$$

$$= \underset{\theta}{\operatorname{argmin}} -\log(p(y|x, \theta) p(\theta))$$

↑
log is a strictly increasing function

$$= \underset{\theta}{\operatorname{argmin}} -\log(p(y|x, \theta)) - \log(p(\theta))$$

$$* \underset{\theta}{\operatorname{argmin}} -\log(p(\theta)) = \underset{\theta}{\operatorname{argmin}} -\log(C) - \log(\exp(-\frac{1}{2} \theta^T \Sigma \theta))$$

$$= \underset{\theta}{\operatorname{argmin}} \frac{1}{2} \frac{\theta^T \Sigma \theta}{\eta^2}$$

$$= \underset{\theta}{\operatorname{argmin}} \frac{1}{2} \frac{\|\theta\|_2^2}{\eta^2}$$

$$\hat{\theta}_{MAP} = \arg \min_{\theta} -\log p(y|x, \theta) + \frac{1}{2\eta^2} \|\theta\|^2$$

$$\Rightarrow \lambda = \frac{1}{2\eta^2}$$

c) * Linear regression: $y = \theta^T x + \epsilon \quad (2)$

$$* \epsilon \sim \mathcal{N}(0, \sigma^2) \quad (1)$$

$$* \theta \sim \mathcal{N}(0, \eta^2 I) \quad (3)$$

$$* X = \begin{bmatrix} -x^{(1)} \\ -x^{(2)} \\ \vdots \\ -x^{(m)} \end{bmatrix}; \vec{y} = \begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{pmatrix}$$

Problem: Come up with a closed form expression for $\hat{\theta}_{MAP}$.

$$* (1) + (2) + (3) \Rightarrow p(y^{(i)} | x^{(i)}, \theta) = \mathcal{N}(x^T \theta, \sigma^2)$$

$$* p(y|x, \theta) = \prod_{i=1}^m p(y^{(i)} | x^{(i)}, \theta)$$

$$\Rightarrow \arg \max_{\theta} p(y|x, \theta) = \arg \min_{\theta} -\log \left(\sum_{i=1}^m p(y^{(i)} | x^{(i)}, \theta) \right)$$

$$= \arg \min_{\theta} -\sum_{i=1}^m \log \left(\frac{1}{\sqrt{2\pi}} \right) - \frac{1}{2} \left(\frac{(y^{(i)} - x^{(i)T} \theta)^2}{\sigma^2} \right)$$

$$= \arg \min_{\theta} \frac{1}{2\sigma^2} (y - X\theta)^T (y - X\theta)$$

$$\Rightarrow \hat{\theta}_{MAP} = \arg \min_{\theta} \frac{1}{2\sigma^2} (Y - X\theta)^T (Y - X\theta) + \underbrace{\frac{1}{2\eta^2} \|\theta\|_2^2}_{\text{Question(b)}}$$

Question(b)

$$* \quad \frac{\partial}{\partial \theta} \left(\frac{1}{2\sigma^2} (Y - X\theta)^T (Y - X\theta) + \frac{1}{2\eta^2} \|\theta\|_2^2 \right)$$

$$= \frac{1}{2\sigma^2} (2X^T X \theta - 2X^T Y) + \frac{1}{2\eta^2} \cdot 2\theta$$

$$= (X^T X \theta - X^T Y) \frac{1}{\sigma^2} + \frac{1}{\eta^2} \theta$$

Setting $\frac{\partial}{\partial \theta}$ to 0

$$\Rightarrow \left(\frac{X^T X}{\sigma^2} - \frac{I}{\eta^2} \right) \theta = X^T Y \cdot \frac{1}{\sigma^2}$$

$$\hat{\theta}_{MAP} = \left(\frac{X^T X}{\sigma^2} - \frac{I}{\eta^2} \right)^{-1} X^T Y \cdot \frac{1}{\sigma^2}$$

- d)
- * $y = \mathbf{x}^T \Theta + \epsilon$
 - * $\epsilon \sim N(0, \sigma^2)$
 - * $\Theta \sim L(0, bI)$

$$* f_L(z | \mu, b) = \frac{1}{2b} \exp\left(-\frac{|z - \mu|}{b}\right)$$

$$* \Theta_{MAP} = \underset{\Theta}{\operatorname{argmin}} -\log(p(y | \mathbf{x}, \Theta)) - \log(p(\Theta))$$

Problem: Show that Θ_{MAP} is equivalent to linear regression with L_1 regularization whose loss is specified as

$$J(\Theta) = \|\mathbf{X}\Theta - \vec{y}\|_2^2 + \gamma \|\Theta\|_1$$

$$\log(p(\Theta)) = -\log\left(\frac{1}{2\sigma b}\right) - \frac{\|\Theta\|_1}{b}$$

$$\Rightarrow \Theta_{MAP} = \underset{\Theta}{\operatorname{argmin}} \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\Theta)^T (\mathbf{y} - \mathbf{X}\Theta) - \log\left(\frac{1}{2\sigma b}\right) + \frac{\|\Theta\|_1}{b}$$

$$= \underset{\Theta}{\operatorname{argmin}} \|\mathbf{X}\Theta - \vec{y}\|_2^2 + \frac{\|\Theta\|_1}{b}$$

$$\Rightarrow \gamma = \frac{1}{b}$$

4) Constructing Kernels

* K_1, K_2 are kernels over $\mathbb{R}^n \times \mathbb{R}^n$

* $a \in \mathbb{R}^+$

* $f: \mathbb{R}^n \rightarrow \mathbb{R}$

* $\phi: \mathbb{R}^n \rightarrow \mathbb{R}^d$

* K_3 kernel over $\mathbb{R}^d \times \mathbb{R}^d$

* $P(x)$ a polynomial over x with positive coefficients

Problem: For each K below, state whether it is a kernel or not.

a) $K(x, z) = K_1(x, z) + K_2(x, z)$

* $K_1 + K_2$ is a symmetric matrix

cause K_1 and K_2 are symmetric

* $z^T(K_1 + K_2)z = z^T K_1 z + z^T K_2 z \geq 0$

$\Rightarrow K$ is a Kernel.

$$b) K(x, z) = K_1(x, z) - K_2(x, z)$$

$$\text{If } K_2 = \alpha K_1$$

$$\Rightarrow z^T K z = z^T K_1 z - \alpha z^T K_1 z < 0$$

K is not necessarily a kernel

$$c) K(x, z) = \alpha K_1(x, z)$$

* αK_1 is symmetric

$$* z^T K z = \alpha z^T K_1 z \geq 0$$

$\Rightarrow K$ is a kernel

$$d) K(x, z) = -\alpha K_1(x, z)$$

$$z^T K z = -\alpha z^T K_1 z < 0$$

K is not a kernel

$$e) K(x, z) = K_1(x, z) K_2(x, z)$$

$$\Rightarrow K_{ij} = \phi_1(x^{(i)})^T \phi_1(x^{(j)}) \phi_2(x^{(i)})^T \phi_2(x^{(j)})$$

$$\Rightarrow Z^T K Z = \sum_i \sum_j z_i \phi_1(x^{(i)})^T \phi_1(x^{(j)}) \phi_2(x^{(i)})^T \phi_2(x^{(j)}) z_j$$

$$= \sum_i \sum_j \sum_k \sum_l z_i \phi_{1k}(x^{(i)})^T \phi_{1l}(x^{(j)}) \phi_{2k}(x^{(i)})^T \phi_{2l}(x^{(j)}) z_j$$

$$= \sum_{lk} \left(\sum_i z_i \phi_{1k}(x^{(i)})^T \phi_{2l}(x^{(i)}) \right) \left(\sum_j z_j \phi_{1h}(x^{(j)})^T \phi_{2h}(x^{(j)}) \right)$$

$$= \sum_{lk} \left(\sum_i z_i \phi_{1k}(x^{(i)})^T \phi_{2l}(x^{(i)}) \right)^2 \geq 0$$

$\Rightarrow K$ is a kernel.

$$f) K(x, z) = f(x) f(z) \quad \text{where } f: \mathbb{R}^n \mapsto \mathbb{R}$$

$$\Rightarrow Z^T K Z = \sum_i \sum_j z_i f(x^{(i)})^T f(x^{(j)}) z_j$$

$$= \left(\sum_i z_i f(x^{(i)}) \right) \left(\sum_j z_j f(x^{(j)}) \right) \geq 0$$

$\Rightarrow K$ is a kernel.

g) The inputs x, z and don't have an effect
on the validity of kernel K_3

h) $K(x, z) = p(K_1(x, z))$

where $p(a)$ is a polynomial with
positive coefficients

$$\Rightarrow p(x) = \sum_{i=0}^D \alpha_i x^i$$

$$\Rightarrow p(K_1(x, z)) = \sum_{i=0}^D \alpha_i K_1(x, z)^i$$

a)

c)

e)

$\Rightarrow K$ is a valid kernel

5) Kernelizing the Perceptron

* $y \in \{0, 1\}$

* $h_\theta(x) = g(\theta^T x)$ where $g(z) = 1$ if $z \geq 0$, -1 otherwise

* Update Rule:

$$\theta^{(i+1)} := \theta^{(i)} + \alpha(y^{(i+1)} - h_{\theta^{(i)}}(x^{(i+1)}))x^{(i+1)}$$

$$* \theta^{(0)} = \vec{0}$$

* $K(x, z) = \phi(x)^T \phi(z)$ where ϕ is a high-dimensional feature mapping

1) According to Representer Theorem

$$\theta^{(i)} = \sum_{j=0}^i \alpha_j y_j \phi(x_j^{(i)}) \quad y_j \in \{-1, 1\}$$

$$\theta_0 = \sum_{j=0}^i \alpha_j \phi(x_j^{(i)}) \quad \alpha_j \in \mathbb{R}$$

$$\theta^{(0)} = \vec{0}$$

(ii) How to calculate $h_{\phi^{(i)}}(x^{(i+1)})$ efficiently?

$$h_{\phi^{(i)}}(x^{(i+1)}) = g(\Theta^{(i)} \phi(x^{(i+1)}))$$

$$= \text{sign} \left(\sum_{j=1}^i \alpha_j \phi(x^{(i)})^T \phi(x^{(i+1)}) \right)$$

$$= \text{sign} \left(\sum_{j=1}^i \alpha_j K(x^{(i)}, x^{(i+1)}) \right)$$

(iii) Modify the update rule mentioned above to use the feature mapping

$$\Theta^{(i+1)} := \Theta^{(i)} + \alpha (y^{(i+1)} - h_{\phi^{(i)}}(x^{(i+1)})) x^{(i+1)}$$

$$\Rightarrow \Theta^{(i+1)} := \sum_{j=1}^i \alpha_j \phi(x^{(i)}) + \alpha (y^{(i+1)} - \text{sign} \left(\sum_{j=1}^i \alpha_j K(x^{(i)}, x^{(i+1)}) \right)) \phi(x^{(i+1)})$$

c) The dot kernel performs badly, cause it doesn't use feature mapping and the data is not linearly separable