

1. Logistic Regression: Training Stability.

(a) A converged quickly while B can not converge.

(b) Dataset B is linearly separable while A is not.

Check the code in "col-grad"
we can write out the loss function

$$\nabla_{\theta} J(\theta) = -\frac{1}{m} \sum_{i=1}^m \frac{y^{(i)} x^{(i)}}{1 + e^{-y^{(i)} x^{(i)} \theta}} \quad \theta := \theta - \text{lr} \cdot \nabla_{\theta} J(\theta)$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-(y^{(i)} x^{(i)} \theta)})$$

The dataset B is linearly separable, so we can update θ that satisfy $y^{(i)} x^{(i)} \theta > 0$ for all cases.

In order to minimize $J(\theta)$, we need parameter θ close to positive infinite.
 $|\nabla_{\theta} J(\theta)|$ increases during this process, so $|\theta_{\text{old}} - \theta_{\text{new}}|$ will get larger

thus it can't converge.

However for unseparable cases, as θ increase, $J(\theta)$ for examples in the other side of classification boundary increases as well. $|\nabla_{\theta} J(\theta)|$ decreases which leads $|\theta_{\text{old}} - \theta_{\text{new}}|$ converges within a small value.

v) ii) No. learning rate is a fixed constant and can not affect $\nabla_{\theta} J(\theta)$.

iii) Yes. There will be some iteration that $\text{lr} \cdot \nabla_{\theta} J(\theta) < 10^{-5}$

iv) No. linear scale will not change the separability of data.

v) Yes. This constrains θ can not be updated to infinite.

vi) Yes. This allows data to be unseparable linearly.

vi) hinge loss: $J(y) = \max(0, 1 - y \hat{y})$, $\hat{y} = w^T x + b$

when dataset is linearly separable, $y \hat{y} > 0$ for all cases.

we can increase parameter w to make $y \hat{y} > 1$.

thus $J(y)$ can be minimize to 0.

2. Model Calibration.

(a) the log likelihood for logistic regression is

$$l(\theta) = \sum_{i=1}^m p(y_i | x_i; \theta) = \sum_{i=1}^m y_i \log p(y_i=1 | x_i; \theta) + (1-y_i) \log p(y_i=0 | x_i; \theta)$$

$$= \sum_{i=1}^m y_i \log h(x^{(i)}) + (1-y_i) \log (1-h(x^{(i)}))$$

$$\frac{\partial l(\theta)}{\partial \theta_j} = \sum_{i=1}^m y_i \frac{1}{h(x^{(i)}) \cdot h'(x^{(i)})} + \frac{(1-y_i)}{1-h(x^{(i)})} (-h'(x^{(i)}))$$

$$= \sum_{i=1}^m y_i \frac{1}{h(x^{(i)})} h(x^{(i)}) (1-h(x^{(i)})) = \sum_{i=1}^m \frac{(1-y_i)}{1-h(x^{(i)})} h(x^{(i)}) (1-h(x^{(i)})) x_j$$

$$= \sum_{i=1}^m y_i (1-h(x^{(i)})) - (1-y_i) h(x^{(i)}) = \sum_{i=1}^m (y_i - h(x^{(i)})) x_j$$

set $f = 0$, because of intercept term $x_0 = 1$.

$$\sum_{i=1}^m y_i = \sum_{i=1}^m h_\theta(x^{(i)})$$

$$y = \mathbb{I}\{\sum y^{(i)} = 1\}. \quad h_\theta(x^{(i)}) = P(y^{(i)} = 1 | x^{(i)}; \theta), |\{i \in I_{0,1} | y^{(i)} = 1\}| = |\{i \in I_{0,1} | y^{(i)} = 1\}| = m.$$

$$\frac{\sum_{i \in I_{0,1}} P(y^{(i)} = 1 | x^{(i)}; \theta)}{|\{i \in I_{0,1} | y^{(i)} = 1\}|} = \frac{\sum_{i \in I_{0,1}} \mathbb{I}\{\sum y^{(i)} = 1\}}{|\{i \in I_{0,1} | y^{(i)} = 1\}|} \text{ holds true.}$$

(b) As a binary classification problem,

$$0.5 < P(y^{(i)} = 1 | x^{(i)}; \theta) < 1.$$

so, for $I \in I_{0.5, 1}$ when model is perfectly calibrated.

$$\frac{\sum_{i \in I_{0.5, 1}} P(y^{(i)} = 1 | x^{(i)}; \theta)}{|\{i \in I_{0.5, 1} | y^{(i)} = 1\}|} = \frac{\sum_{i \in I_{0.5, 1}} \mathbb{I}\{\sum y^{(i)} = 1\}}{|\{i \in I_{0.5, 1} | y^{(i)} = 1\}|} < \frac{\sum_{i \in I_{0,1}} \mathbb{I}\{\sum y^{(i)} = 1\}}{|\{i \in I_{0,1} | y^{(i)} = 1\}|} = \text{precision.}$$

from above we can see the model doesn't achieve its perfect accuracy.
The converse is not necessarily true as well.

(c) Applying L_2 regularization to the loss function

$$J(\theta) = -\frac{1}{m} \sum_i^m y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log (1-h_\theta(x^{(i)})) + \frac{1}{2} \lambda \|\theta\|_2^2$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = -\frac{1}{m} \sum_i^m (y^{(i)} - h_\theta(x^{(i)})) x_j + \lambda \theta_j = 0.$$

set $f = 0$, $x_f = 1$.

$$\sum_i^m y^{(i)} = \sum_i^m h_\theta(x^{(i)}) - \lambda \theta_0. \text{ so the model will not be well-calibrated, it can achieve}$$

the best precision.

3. Bayesian Interpretation of Regularization

MAP: maximum a posteriori estimation

$$\theta_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} p(\theta | x, y).$$

MLE: maximum likelihood estimation

$$\theta_{\text{MLE}} = \underset{\theta}{\operatorname{argmax}} p(y | x, \theta)$$

$$(a) p(\theta | x, y) = \frac{p(x, y | \theta)}{p(x, y)} = \frac{p(y | x, \theta) p(x, \theta)}{p(x, y)} = \frac{p(y | x, \theta) p(\theta | x) p(x)}{p(x, y)}$$

$$= \frac{p(y | x, \theta) p(\theta | x) p(x)}{p(x, y)}$$

Assume that $p(\theta | x) = p(\theta)$

$$\therefore \underset{\theta}{\operatorname{argmax}} p(\theta | x, y) = \underset{\theta}{\operatorname{argmax}} \frac{p(y | x, \theta) p(\theta) p(x)}{p(x, y)} = \underset{\theta}{\operatorname{argmax}} p(y | x, \theta) p(\theta).$$

(b) negative log likelihood loss function is that:

$$J(\theta) = -\log p(y | x, \theta) + \lambda \|\theta\|_2^2$$

MLE estimation is

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

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

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

Remark:

Different regularization stands for different assumptions for θ .

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

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

$$\theta \sim \mathcal{N}(0, \sigma^2 I)$$

$$p(\theta) = \frac{1}{(2\pi)^{\frac{m}{2}} \sigma^m} \exp\left\{-\frac{1}{2\sigma^2} \theta^\top \theta\right\} = \frac{1}{(2\pi)^{\frac{m}{2}} \sigma^m} \exp\left\{-\frac{1}{2\sigma^2} \|\theta\|^2\right\}$$

$$\log p(\theta) = -\frac{1}{2\sigma^2} \|\theta\|^2 - \frac{m}{2} \log 2\pi - m \log \sigma.$$

$$\begin{aligned} \theta_{MAP} &= \underset{\theta}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|\theta\|^2 + \frac{m}{2} \log \sigma + m \log \sigma - \log p(y|x, \theta) \\ &= \underset{\theta}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|\theta\|^2 - \log p(y|x, \theta) \end{aligned}$$

$$\lambda = \frac{1}{2\sigma^2}.$$

(c) For a specific instance, $\theta \sim \mathcal{N}(0, \sigma^2 I)$

$$\begin{aligned} y^{(i)} &= \theta^\top x^{(i)} + \epsilon \\ y^{(i)} &\sim \mathcal{N}(\theta^\top x^{(i)}, \sigma^2) \end{aligned}$$

$$\begin{aligned} \theta_{MAP} &= \underset{\theta}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|\theta\|^2 - \log p(y|x, \theta) \\ &= \underset{\theta}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|\theta\|^2 - \sum_{i=1}^m \log p(y^{(i)}|x^{(i)}, \theta) \\ &= \underset{\theta}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|\theta\|^2 - \sum_{i=1}^m \log \frac{1}{(2\pi)^{\frac{1}{2}} \sigma} \exp\left\{-\frac{1}{2\sigma^2} (y^{(i)} - \theta^\top x^{(i)})^2\right\} \\ &= \underset{\theta}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|\theta\|^2 + \sum_{i=1}^m \frac{1}{2\sigma^2} (y^{(i)} - \theta^\top x^{(i)})^2 \\ &= \underset{\theta}{\operatorname{argmin}} \frac{1}{2\sigma^2} (\vec{y} - X\theta)^\top (\vec{y} - X\theta) + \frac{1}{2\sigma^2} \theta^\top \theta \end{aligned}$$

$$J(\theta) = \frac{1}{2\sigma^2} (\vec{y} - X\theta)^\top (\vec{y} - X\theta) + \frac{1}{2\sigma^2} \theta^\top \theta.$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{\sigma^2} \sum_{i=1}^m (y^{(i)} - \theta^\top x^{(i)}) x_j^{(i)} + \frac{1}{\sigma^2} \theta_j.$$

$$\nabla_{\theta} J(\theta) = \frac{1}{\sigma^2} X^\top (X\theta - \vec{y}) + \frac{1}{\sigma^2} \theta.$$

Set $\nabla_{\theta} J(\theta) = 0$, we have $\theta = (X^\top X + \frac{\sigma^2}{\sigma^2} I)^{-1} X^\top \vec{y}$

(d) $\theta \sim \mathcal{L}(0, b^2)$

$$p(\theta) = \frac{1}{2b} \exp(-\frac{|\theta|}{b}) \quad y|x, \theta \sim \mathcal{N}(\theta^\top x, b^2)$$

$$\begin{aligned} \theta_{MAP} &= \underset{\theta}{\operatorname{argmax}} p(\theta|x, y) = \underset{\theta}{\operatorname{argmax}} p(y|x, \theta) p(\theta) \\ &= \underset{\theta}{\operatorname{argmax}} \log p(y|x, \theta) + \log p(\theta) \\ &= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log \frac{1}{(2\pi)^{\frac{1}{2}} b} \exp\left\{-\frac{1}{2b^2} (y^{(i)} - \theta^\top x^{(i)})^2\right\} + \log \frac{1}{2b} \exp\left\{-\frac{|\theta|}{b}\right\} \\ &= \underset{\theta}{\operatorname{argmax}} -\sum_{i=1}^m \frac{1}{2b^2} (y^{(i)} - \theta^\top x^{(i)})^2 - \frac{|\theta|}{b} \\ &= \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^m \frac{1}{2b^2} (y^{(i)} - \theta^\top x^{(i)})^2 + \frac{|\theta|}{b} \\ &= \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^m (y^{(i)} - \theta^\top x^{(i)})^2 + \frac{2b^2}{b} |\theta| \\ &= \underset{\theta}{\operatorname{argmin}} J(\theta). \quad \gamma = \frac{2b^2}{b} \end{aligned}$$

4. Constructing Kernels.

(a) Yes.

K_1, K_2 are symmetric. So $K_1 + K_2$ is symmetric as well.

K_1, K_2 are PSD.

$$z^T K z = z^T (K_1 + K_2) z = z^T K_1 z + z^T K_2 z \geq 0. \text{ So } K \text{ is PSD as well.}$$

(b) No.

counterexample, $K_2 = 2K_1$

$$z^T K z = z^T (K_1 - 2K_1) z = -z^T K_1 z, \text{ NSD.}$$

(c) Yes. α is a positive real number.

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

(d) No.

$$z^T K z = -\alpha z^T K_1 z \leq 0, \text{ NSD.}$$

(e) Yes.

$$\begin{aligned} z^T K z &= \sum_i \sum_j z_i K_{ij} z_j \\ &= \sum_i \sum_j z_i K_1(x^{(i)}, x^{(j)}) K_2(x^{(i)}, x^{(j)}) z_j \\ &= \sum_i \sum_j z_i (\sum_a \phi_{1a}(x^{(i)})^\top \phi_{1a}(x^{(j)}) \phi_{2a}(x^{(i)})^\top \phi_{2a}(x^{(j)})) z_j \\ &= \sum_a \sum_b z_i (\sum_i \phi_{1a}(x^{(i)}) \phi_{1a}(x^{(i)})) (\sum_j \phi_{2b}(x^{(j)}) \phi_{2b}(x^{(j)})) z_j \\ &= \sum_a \sum_b (\phi_{1a}(x^{(i)}) \phi_{2b}(x^{(i)})) z_i)^2 \geq 0 \end{aligned}$$

(f) Yes

$$\begin{aligned} z^T K z &= \sum_i \sum_j z_i K_{ij} z_j \\ &= \sum_i \sum_j z_i f(x^{(i)}) f(x^{(j)}) z_j \\ &= (\sum_i z_i f(x^{(i)}))^2 \geq 0. \end{aligned}$$

(g) Yes

$$\begin{aligned} z^T K z &= \sum_i \sum_j z_i K_{ij} z_j \\ &= \sum_i \sum_j z_i K_3(\phi(x^{(i)}), \phi(x^{(j)})) z_j \\ &= \sum_i \sum_j z_i \phi_3(\phi(x^{(i)}))^\top \phi_3(\phi(x^{(j)})) z_j \\ &= \sum_i \sum_j z_i \sum_a \phi_{3a}(\phi(x^{(i)})) \phi_{3a}(\phi(x^{(j)})) z_j \\ &= \sum_a \sum_b (\phi_{3a}(\phi(x^{(i)})) z_i)^2 \geq 0. \end{aligned}$$

(h) Yes.

$$p(x) = a_0 + a_1 x + \dots + a_n x^n, \quad a_0, \dots, a_n \in \mathbb{R}^+$$

$$k(x, z) = p(K(x, z)) = a_0 + a_1 K_1(x, z) + a_2 K_2(x, z)^2 + \dots + a_n K_n(x, z)^n.$$

As I proved before, $\circ k(x, z) = K_1(x, z) K_2(x, z)$ is a valid kernel.

generalize this conclusion to multinomial case:

$$k(x, z) = (K_1(x, z))^n \text{ is a valid kernel.}$$

② $k(x_1z) = \alpha k_1(x_1z)$, $\alpha \in \mathbb{R}^+$ is a valid kernel.

③ $k(x_1z) = k_1(x_1z) + k_2(x_1z)$ is a valid kernel

So $k(x_1z) = p(k_1(x_1z))$ is a valid kernel.

5. Kernelizing the Perceptron.

(a) Apply kernel trick to perceptron to make it work in high dimensional space (eg. ϕ)

(i) Recall the update rule:

$$\theta^{(i+1)} := \theta^{(i)} + \alpha (y^{(i+1)} - h_{\theta^{(i)}}(x^{(i+1)})) x^{(i+1)}$$
$$\theta^{(i)} = \sum_{j=0}^{\infty} \beta_j \phi(x^{(j)}) , \quad \theta^0 = \vec{0}$$

$$(ii) h_{\theta^{(i)}}(x^{(i+1)}) = g(\theta^{(i)\top} \phi(x^{(i+1)}))$$
$$= g\left(\sum_{j=0}^{\infty} \beta_j \phi^T(x^{(j)}) \phi(x^{(i+1)})\right)$$
$$= \text{sign}\left(\sum_{j=0}^{\infty} \beta_j \langle \phi(x^{(j)}), \phi(x^{(i+1)}) \rangle\right)$$
$$= \text{sign}\left(\sum_{j=0}^{\infty} \beta_j K(x^{(j)}, \phi(x^{(i+1)}))\right)$$

$$(iii) \theta^{(i+1)} := \theta^{(i)} + \alpha (y^{(i+1)} - h_{\theta^{(i)}}(\phi(x^{(i+1)}))) \phi(x^{(i+1)})$$
$$= \sum_{j=1}^{\infty} \beta_j \phi(x^{(j)}) + \underbrace{\alpha (y^{(i+1)} - \text{sign}\left(\sum_{j=0}^{\infty} \beta_j K(x^{(j)}, x^{(i+1)}))\right) \phi(x^{(i+1)})}_{\beta^{(i+1)}}$$
$$= \sum_{j=1}^{\infty} \beta_j \phi(x^{(j)})$$
$$\therefore \beta^{(i+1)} := \alpha (y^{(i+1)} - \text{sign}\left(\sum_{j=0}^{\infty} \beta_j K(x^{(j)}, x^{(i+1)}))\right))$$

(c) Dot kernel performs badly. For the dot kernel, $\phi(x)=x$. So it doesn't map the data to higher dimension, the data is still not linearly separable.

6. Spam Classification