

3

1' 17 1/2

277595184

1/1000

1/100

1. Regularized polynomial regression

We derived in class the solution for a zero-degree polynomial regression. Consider the problem of regularized polynomial regression.

$$Err(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n (h_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2 + \lambda \|\mathbf{w}\|^2$$

1. Derive the solution for a polynomial of degree 0: $h_{\mathbf{w}}(\mathbf{x}) = w_0$. Analyze the solution in the limit of $\lambda \rightarrow \infty$ and $\lambda \rightarrow 0$.
2. Derive the solution for a polynomial of degree 1: $h_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1 x$, by computing the derivatives w.r.t. w_0 and w_1 and writing a system of two linear equations in w_0 and w_1 . No need to solve the system. Analyze the solution in the limit of $\lambda \rightarrow \infty$ and $\lambda \rightarrow 0$.

$$Err(w) = \frac{1}{n} \sum_{i=1}^n (h_w(x_i) - y_i)^2 + \lambda \|w\|^2$$

$$\Downarrow \quad H_w(x) = w_0$$

$$Err(w) = \frac{1}{n} \sum_{i=1}^n (w_0 - y_i)^2 + \lambda \|w_0\|^2$$

$$w \quad \text{derivative} \quad \Downarrow$$

$$Err(w)' = \frac{1}{n} \cdot 2 \sum_{i=1}^n (w_0 - y_i) + 2\lambda w_0 = 0$$

$$\Downarrow$$

$$\frac{2}{n} \cdot \sum_{i=1}^n (w_0 - y_i) = -2\lambda w_0$$

$$\Downarrow \quad / \frac{n}{2}$$

$$\sum_{i=1}^n (w_0 - y_i) = -\lambda n w_0$$

$$n w_0 - \sum_{i=1}^n y_i = -\lambda n w_0$$

⇓

$$n w_0 \leftarrow \lambda n w_0 = \sum_{i=1}^n y_i$$

⇓

$$w_0 (n \leftarrow \lambda n) = \sum_{i=1}^n y_i$$

⇓

$$w_0 = \frac{\sum_{i=1}^n y_i}{n \leftarrow \lambda n}$$

$\lambda \rightarrow \infty$

$\uparrow \infty$

∞

$0 \leftarrow \infty$

$\lambda \rightarrow 0$

$\infty \leftarrow 0$

$$w_0 = \frac{\sum_{i=1}^n y_i}{n} = \bar{y}$$

2. Derive the solution for a polynomial of degree 1: $h_w(x) = w_0 + w_1 x$, by computing the derivatives w.r.t. w_0 and w_1 and writing a system of two linear equations in w_0 and w_1 . No need to solve the system. Analyze the solution in the limit of $\lambda \rightarrow \infty$ and $\lambda \rightarrow 0$.

2

$$E_{RL}(w) = \frac{1}{n} \sum_{i=1}^n (h_w(x_i) - y_i)^2 + \lambda \|w\|^2$$

$$h_w(x_i) = w_0 + w_1 x_i$$

$$\|w\|^2 = w_0^2 + w_1^2$$

$$E_{RL}(w) = \frac{1}{n} \sum_{i=1}^n (w_0 + w_1 x_i - y_i)^2 + \lambda (w_0^2 + w_1^2)$$

$$w_0 \quad w_1 \quad x_i \quad y_i$$

$$\frac{1}{n} \sum_{i=1}^n 2(w_0 + w_1 x_i - y_i) + 2\lambda w_0 = 0$$

$$\frac{1}{n} \sum_{i=1}^n 2(w_0 + w_1 x_i - y_i) + 2\lambda w_0 = 0$$

$$\sum_{i=1}^n (w_0 + w_1 x_i - y_i) + \lambda n w_0 = 0$$

$$n w_0 + w_1 \sum_{i=1}^n x_i - \sum_{i=1}^n y_i + \lambda n w_0 = 0$$

$$\frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$$

$$w_0(1 + \lambda) = \sum_{i=1}^n y_i - w_1 \sum_{i=1}^n x_i$$

\Downarrow

$$w_0 = \frac{\sum_{i=1}^n y_i - w_1 \sum_{i=1}^n x_i}{1 + \lambda}$$

vers

$\lambda \rightarrow 0$

$$w_0 = \frac{\sum_{i=1}^n y_i - w_1 \sum_{i=1}^n x_i}{1}$$

$\lambda \rightarrow \infty$

vers

$w_0 = 0$

vers

w_1

✓

2/5/21

~80

$$E_{\text{MSE}}(w) = \frac{1}{n} \sum_{i=1}^n (w_0 + w_1 x_i - y_i)^2 + \lambda (w_0^2 + w_1^2)$$

$$E'_{\text{MSE}}(w) = \frac{1}{n} \sum_{i=1}^n 2 \cdot x_i (w_0 + w_1 x_i - y_i) + 2\lambda w_1 = 0$$

$$\Downarrow \cdot \frac{n}{2}$$

$$\sum_{i=1}^n x_i (w_0 + w_1 x_i - y_i) + \lambda n w_1 = 0$$

✓

$$w_0 \sum_{i=1}^n x_i + w_1 \cdot \left(\sum_{i=1}^n x_i^2 \right) - \sum_{i=1}^n x_i y_i + \lambda n w_1$$

✓

$$w_0 \sum_{i=1}^n x_i + w_1 \cdot \left(\sum_{i=1}^n x_i^2 + \lambda n \right) - \sum_{i=1}^n x_i y_i = 0$$

✓

$$W_1 \cdot \left(\sum_{i=1}^n x_i^2 + \lambda n \right) = \sum_{i=1}^n x_i y_i - W_0 \sum_{i=1}^n x_i$$

↓

$$W_1 = \frac{\sum_{i=1}^n x_i y_i - W_0 \sum_{i=1}^n x_i}{\sum_{i=1}^n x_i^2 + \lambda n}$$

$$\sum_{i=1}^n x_i^2 + \lambda n$$

$\lambda \rightarrow \infty$ $\lambda \rightarrow \infty$

$W_1 \rightarrow 0$ $(\lambda \rightarrow \infty)$

$-\lambda \rightarrow 0$ $\lambda \rightarrow 0$ $\lambda \rightarrow 0$

$$W_1 = \frac{\sum_{i=1}^n x_i y_i - W_0 \sum_{i=1}^n x_i}{\sum_{i=1}^n x_i^2}$$

2. Logistic regression

1. Prove that the logistic regression classifier is equivalent to a softmax over a linear multiclass classifier for two classes $y = "a", y = "b"$, when their separating hyperplanes obey $\mathbf{w}_a = -\mathbf{w}_b$.
2. For a vector $\mathbf{z} \in \mathbb{R}^K$, consider the regular softmax function:

$$\text{softmax}_i(\mathbf{z}) = \frac{\exp(z_i)}{\sum_{k=1}^K \exp(z_k)} \quad (1)$$

For any vector $\mathbf{b} = (b, \dots, b) \in \mathbb{R}^K$ for some $b \in \mathbb{R}$, prove that $\text{softmax}_i(\mathbf{z}) = \text{softmax}_i(\mathbf{z} - \mathbf{b})$ for any $1 \leq i \leq K$.

3. For a vector $\mathbf{z} \in \mathbb{R}^K$, consider the softmax function that is scaled by a constant $T \in \mathbb{R}$:

$$f_i(\mathbf{z}) = \frac{\exp(Tz_i)}{\sum_{k=1}^K \exp(Tz_k)} \quad (2)$$

Further, for a vector $\mathbf{z} = (z_1, \dots, z_K) \in \mathbb{R}^K$, instead of considering the $\arg \max(z_1, \dots, z_K)$ function as a function with categorical output $1, \dots, K$ (corresponding to the index of a vector's largest element), consider the $\arg \max$ function with **one-hot** representation of the output (assuming there is a unique maximum element):

$$\arg \max(z_1, \dots, z_K) = (y_1, \dots, y_K) = (0, \dots, 0, 1, 0, \dots, 0) \quad (3)$$

where $y_i = 1$ if and only if $i = \arg \max(z_1, \dots, z_K)$, meaning that z_i is the unique maximum value of $\mathbf{z} = (z_1, \dots, z_K)$.

- (a) For any vector $\mathbf{z} \in \mathbb{R}^K$ whose maximum element is unique, show that the softmax converges to the $\arg \max$ function as $T \rightarrow \infty$, i.e., prove that:

$$\lim_{T \rightarrow \infty} (f_1(\mathbf{z}), \dots, f_K(\mathbf{z})) = \arg \max(z_1, \dots, z_K) \quad (4)$$

when $\arg \max$ is in **one-hot** encoding.

- (b) For any vector $\mathbf{z} \in \mathbb{R}^K$ whose maximum element is **not necessarily** unique, compute $\lim_{T \rightarrow \infty} (f_1(\mathbf{z}), \dots, f_K(\mathbf{z}))$ and provide a literal interpretation for your result.
 - (c) For any vector $\mathbf{z} \in \mathbb{R}^K$, what happens when $T \rightarrow 0$? Namely, compute the limit $\lim_{T \rightarrow 0} (f_1(\mathbf{z}), \dots, f_K(\mathbf{z}))$.
4. Write the gradient update rule for a logistic regression model, when the usual loss of the negative log likelihood is now regularized with the square of the L_2 norm over the weight vector $\frac{1}{2} \|\mathbf{w}\|^2$.

1. Prove that the logistic regression classifier is equivalent to a softmax over a linear multiclass classifier for two classes $y = "a", y = "b"$, when their separating hyperplanes obey $\mathbf{w}_a = -\mathbf{w}_b$.

$$\mathbf{w}_a = -\mathbf{w}_b \quad \text{yes}$$

$$\mathbf{w}' = 2\mathbf{w}_a \quad \text{yes} \quad \text{Logistic regression} \quad \text{yes}$$

$$P(y = "a" | \mathbf{x}) = \sigma(2\mathbf{w}_a^T \mathbf{x})$$

$$P(y = "b" | \mathbf{x}) = 1 - \sigma(2\mathbf{w}_b^T \mathbf{x})$$

$$\cdot \text{Soft max} \quad \text{yes}$$

$$P(y = 'a' / x) = \frac{e^{w_a^T x}}{e^{w_a^T x} + e^{w_b^T x}}$$

$$\Downarrow \quad w_a = -w_b$$

$$\frac{e^{w_a^T x}}{e^{w_a^T x} + e^{-w_a^T x}}$$

$$\Downarrow$$

$$1$$

$$\frac{e^{2w_a^T x}}{e^{2w_a^T x} + 1}$$

$$\approx \sigma(2w_a^T x)$$

$$\therefore \text{e) ps } \sigma(2w_a^T x) = \sigma(2w_a^T x) - 0 \quad \text{UNJP}$$

(2)

2. For a vector $\mathbf{z} \in \mathbb{R}^K$, consider the regular softmax function:

$$\text{softmax}_i(\mathbf{z}) = \frac{\exp(z_i)}{\sum_{k=1}^K \exp(z_k)} \quad (1)$$

For any vector $\mathbf{b} = (b, \dots, b) \in \mathbb{R}^K$ for some $b \in \mathbb{R}$, prove that $\text{softmax}_i(\mathbf{z}) = \text{softmax}_i(\mathbf{z} - \mathbf{b})$ for any $1 \leq i \leq K$.

$$\begin{aligned} \text{softmax}_i(\mathbf{z} - \mathbf{b}) &= \frac{e^{z_i - b}}{\sum_{k=1}^K e^{z_k - b}} = \frac{e^{z_i} \cdot e^{-b}}{\sum_{k=1}^K e^{z_k} \cdot e^{-b}} = \\ &= \frac{e^{z_i} \cdot \cancel{e^{-b}}}{\cancel{e^{-b}} \sum_{k=1}^K e^{z_k}} = \frac{e^{z_i}}{\sum_{k=1}^K e^{z_k}} = \text{softmax}_i(\mathbf{z}) \end{aligned}$$

3. For a vector $\mathbf{z} \in \mathbb{R}^K$, consider the softmax function that is scaled by a constant $T \in \mathbb{R}$:

$$f_i(\mathbf{z}) = \frac{\exp(Tz_i)}{\sum_{k=1}^K \exp(Tz_k)} \quad (2)$$

Further, for a vector $\mathbf{z} = (z_1, \dots, z_K) \in \mathbb{R}^K$, instead of considering the $\arg \max(z_1, \dots, z_K)$ function as a function with categorical output $1, \dots, K$ (corresponding to the index of a vector's largest element), consider the $\arg \max$ function with **one-hot** representation of the output (assuming there is a unique maximum element):

$$\arg \max(z_1, \dots, z_K) = (y_1, \dots, y_K) = (0, \dots, 0, 1, 0, \dots, 0) \quad (3)$$

where $y_i = 1$ if and only if $i = \arg \max(z_1, \dots, z_K)$, meaning that z_i is the unique maximum value of $\mathbf{z} = (z_1, \dots, z_K)$.

- (a) For any vector $\mathbf{z} \in \mathbb{R}^K$ whose maximum element is unique, show that the softmax converges to the $\arg \max$ function as $T \rightarrow \infty$, i.e., prove that:

$$\lim_{T \rightarrow \infty} (f_1(\mathbf{z}), \dots, f_K(\mathbf{z})) = \arg \max(z_1, \dots, z_K) \quad (4)$$

when $\arg \max$ is in **one-hot** encoding.

- (b) For any vector $\mathbf{z} \in \mathbb{R}^K$ whose maximum element is **not necessarily** unique, compute $\lim_{T \rightarrow \infty} (f_1(\mathbf{z}), \dots, f_K(\mathbf{z}))$ and provide a literal interpretation for your result.
- (c) For any vector $\mathbf{z} \in \mathbb{R}^K$, what happens when $T \rightarrow 0$? Namely, compute the limit $\lim_{T \rightarrow 0} (f_1(\mathbf{z}), \dots, f_K(\mathbf{z}))$.

$$f_i(\mathbf{z}) = \frac{\exp(Tz_i)}{\sum_{k=1}^K \exp(Tz_k)}$$

$$\lim_{T \rightarrow \infty} (f_1(\mathbf{z}), \dots, f_K(\mathbf{z})) = \arg \max(z_1, \dots, z_K) =$$

normalized value of the vector \mathbf{z} in the direction of the maximum element

$$f_i(\mathbf{z}) = \frac{\exp(Tz_i)}{\sum_{k=1}^K \exp(Tz_k)} = \frac{\exp(Tz_i)}{\exp(Tz_i) \sum_{k \neq i} \exp(Tz_k - Tz_i)}$$

$$\cancel{\exp(T_i)}$$

$$\cancel{\exp(T_i)} \left(1 + \sum_{\substack{k=1 \\ k \neq i}}^K \frac{\exp(T_{2k})}{\exp(T_i)} \right) =$$

$$\frac{1}{1 + \sum_{\substack{k=1 \\ k \neq i}}^K \frac{\exp(T_{2k})}{\exp(T_i)}} = \frac{1}{1 + \sum_{k=1}^K \exp(T(2k-i))} \stackrel{\text{is } i}{=} \frac{1}{\infty} = 0$$

$$i = \max \{1, N\} \rightarrow \infty$$

$$f_i(z) = \frac{\exp(T_{\max})}{\sum_{k=1}^K \exp(T_{2k})} = \frac{\exp(T_{\max})}{\exp(T_{\max}) + \sum_{\substack{k=1 \\ k \neq \max}}^K \exp(T_{2k})}$$

$$\cancel{\exp(T_{\max})} \left(1 + \sum_{\substack{k=1 \\ k \neq \max}}^K \frac{\exp(T_{2k})}{\exp(T_{\max})} \right) =$$

(b) For any vector $\mathbf{z} \in \mathbb{R}^K$ whose maximum element is **not necessarily** unique, compute $\lim_{T \rightarrow \infty} (f_1(\mathbf{z}), \dots, f_K(\mathbf{z}))$ and provide a literal interpretation for your result.

2

$$f_i(\mathbf{z}) = \frac{\exp(T_{2max})}{\sum_{k=1}^K \exp(T_{2k})}$$

$\sum_{k=1}^K \exp(T_{2k})$
 $\sum_{k=1}^K \exp(T_{2k})$
 $\sum_{k=1}^K \exp(T_{2k})$

$$\frac{\exp(T_{2max})}{n \cdot \exp(T_{2max}) + \sum_{k=1, k \neq 2max}^K \exp(T_{2k})} =$$

$$\frac{\exp(T_{2max})}{n \cdot \exp(T_{2max}) + \sum_{k=1, k \neq 2max}^K \exp(T_{2k})} = \frac{1}{n + \sum_{k=1, k \neq 2max}^K \frac{\exp(T_{2k})}{\exp(T_{2max})}} = \frac{1}{n}$$

780

ms

0

(j)j

$2; \leq 2 \max$

\log

$\frac{1}{c}$

(c) For any vector $\mathbf{z} \in \mathbb{R}^K$, what happens when $T \rightarrow 0$? Namely, compute the limit $\lim_{T \rightarrow 0} (f_1(\mathbf{z}), \dots, f_K(\mathbf{z}))$.

(2)

$$f_i(z) = \frac{\exp(Tz_i)}{\sum_{k=1}^K \exp(Tz_k)}$$

$$\sim \text{as } T \rightarrow 0 \quad e/1) \quad \sqrt{e}/1 \quad (5)$$

$$\frac{\exp(Tz_i)}{\sum_{k=1}^K \exp(Tz_k)} = \frac{\exp(Tz_i)}{\exp(Tz_i) + \sum_{k=1, k \neq i}^K \exp(Tz_k)}$$

$$\frac{\exp(Tz_i)}{\exp(Tz_i) + \sum_{k=1, k \neq i}^K \exp(Tz_k)} = \frac{1}{1 + \sum_{k=1, k \neq i}^K \exp(Tz_k - Tz_i)} \xrightarrow{T \rightarrow 0} \frac{1}{1 + \sum_{k=1, k \neq i}^K 1} = \frac{1}{K}$$

$$\frac{1}{1 + \sum_{k=1, k \neq i}^K 1} = \frac{1}{K}$$

4. Write the gradient update rule for a logistic regression model, when the usual loss of the negative log likelihood is now regularized with the square of the L_2 norm over the weight vector $\frac{1}{2}||\mathbf{w}||^2$.

4

\mathbf{w} - weight vector

$\mathbf{x}^{(i)}$ - feature vector for the sample

$y^{(i)}$ - true label for the i sample

n - number of training samples

λ - regularization parameter

$$\hat{y}^{(i)} = \sigma(\mathbf{w}^T \mathbf{x}^{(i)}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}^{(i)}}} \quad \text{logistic regression}$$

$$J(\mathbf{w}) = L(\mathbf{w}) + \frac{\lambda}{2} ||\mathbf{w}||^2 = -\frac{1}{n} \sum_{i=1}^n y^{(i)} \cdot \log(\hat{y}^{(i)}) +$$

$$(1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

$$\Delta \mathbf{w} J(\mathbf{w}) = -\frac{1}{n} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)}) \cdot \mathbf{x}^{(i)} + \lambda \mathbf{w}$$

Gradient update rule ✓

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} J(\mathbf{w})$$

α = learning rate ✓

$$w \leftarrow w + \delta \frac{1}{n} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)}) \cdot x^{(i)} - \delta \lambda w$$

- (1.1.1)