

1. (10 points) **Noisy linear regression**

Brayr Tilban Elberier

A real estate company have assigned us the task of building a model to predict the house prices in Westwood. For this task, the company has provided us with a dataset \mathcal{D} :

$$\mathcal{D} = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})\}$$

where $x^{(i)} \in \mathbb{R}^d$ is a feature vector of the i^{th} house and $y^{(i)} \in \mathbb{R}$ is the price of the i^{th} house. Since we just learned about linear regression, so we have decided to use a linear regression model for this task. Additionally, the IT manager of the real estate company has requested us to design a model with small weights. In order to accommodate his request, we will design a linear regression model with parameter regularization. In this problem, we will navigate through the process of achieving regularization by adding noise to the feature vectors. Recall, that we define the cost function in a linear regression problem as:

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - (x^{(i)})^T \theta)^2$$

where $\theta \in \mathbb{R}^d$ is the parameter vector. As mentioned earlier, we will be adding noise to the feature vectors in the dataset. Specifically, we will be adding zero-mean gaussian noise of known variance σ^2 from the distribution

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

where $\mathbf{I} \in \mathbb{R}^{d \times d}$ and $\sigma \in \mathbb{R}$. With the addition of gaussian noise the modified cost function is given by,

$$\tilde{\mathcal{L}}(\theta) = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - (x^{(i)} + \delta^{(i)})^T \theta)^2$$

where $\delta^{(i)}$ are i.i.d noise vectors with $\delta^{(i)} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$.

- (a) (6 points) Express the expectation of the modified loss over the gaussian noise, $\mathbb{E}_{\delta \sim \mathcal{N}}[\tilde{\mathcal{L}}(\theta)]$, in terms of the original loss plus a term independent of the dataset \mathcal{D} . To be precise, your answer should be of the form:

$$\mathbb{E}_{\delta \sim \mathcal{N}}[\tilde{\mathcal{L}}(\theta)] = \mathcal{L}(\theta) + R$$

where R is not a function of \mathcal{D} . For answering this part, you might find the following result useful:

$$\mathbb{E}_{\delta \sim \mathcal{N}}[\delta \delta^T] = \sigma^2 \mathbf{I}$$

- (b) (2 points) Based on your answer to (a), under expectation what regularization effect would the addition of the noise have on the model?
- (c) (1 point) Suppose $\sigma \rightarrow 0$, what effect would this have on the model?
- (d) (1 point) Suppose $\sigma \rightarrow \infty$, what effect would this have on the model?

$$\textcircled{1} a. E[\mathcal{L}(\theta)]$$

$$= E\left[\frac{1}{N} \sum_{i=1}^N (y^{(i)} - (x^{(i)} + f^{(i)})^T \theta)^2\right]$$

$$= \frac{1}{N} \sum_{i=1}^N E[(y^{(i)} - (x^{(i)} + f^{(i)})^T \theta)^2]$$

$$= \underbrace{(y^{(i)} - x^{(i)T} \theta)^2}_{\text{blue circle}} - \underbrace{2(y^{(i)} - x^{(i)T} \theta)(f^{(i)T} \theta)}_{\text{red circle}} + \underbrace{(f^{(i)T} \theta)^2}_{\text{green circle}}$$

$$E[(y^{(i)} - x^{(i)T} \theta)^2] = (y^{(i)} - x^{(i)T} \theta)^2$$

$$E[-2(y^{(i)} - x^{(i)T} \theta)(f^{(i)T} \theta)]$$

$$= -2(y^{(i)} - x^{(i)T} \theta) E[f^{(i)T} \theta]$$

$$= 0$$

$$E[(f^{(i)T} \theta)^2] = E[(f^{(i)T} \theta)^T (f^{(i)T} \theta)]$$

$$= E[\theta^T f^{(i)} f^{(i)T} \theta]$$

$$= \theta^T E[f^{(i)} f^{(i)T}] \theta = \theta^T \sigma^2 I \theta$$

$$= \sigma^2 \theta^T \theta = \sigma^2 \|\theta\|_2^2$$

$$E_{\theta \sim N}[\mathcal{L}(\theta)] = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - x^{(i)T} \theta)^2 + \sigma^2 \|\theta\|_2^2$$

↳ mean square error ↳ L2 reg.

$$E_{f \sim N[\tilde{L}(\theta)]} = \mathcal{L}(\theta) + \sigma^2 \|\theta\|_2^2$$

- b. It would have L2 regularization effect
- c. Close to no effect at all, overfit?
- d. Model could be heavily regularized and potentially underfit.

3. (30 points) **Softmax classifier gradient.** For softmax classifier, derive the gradient of the log likelihood.

Concretely, assume a classification problem with c classes

- Samples are $(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})$, where $\mathbf{x}^{(j)} \in \mathbb{R}^n$, $y^{(j)} \in \{1, \dots, c\}$, $j = 1, \dots, m$
- Parameters are $\theta = \{\mathbf{w}_i, b_i\}_{i=1, \dots, c}$
- Probabilistic model is

$$\Pr(y^{(j)} = i \mid \mathbf{x}^{(j)}, \theta) = \text{softmax}_i(\mathbf{x}^{(j)})$$

where

$$\text{softmax}_i(\mathbf{x}) = \frac{e^{\mathbf{w}_i^T \mathbf{x} + b_i}}{\sum_{k=1}^c e^{\mathbf{w}_k^T \mathbf{x} + b_k}}$$

Derive the log-likelihood \mathcal{L} , and its gradient w.r.t. the parameters, $\nabla_{\mathbf{w}_i} \mathcal{L}$ and $\nabla_{b_i} \mathcal{L}$, for $i = 1, \dots, c$.

Note: We can group \mathbf{w}_i and b_i into a single vector by augmenting the data vectors with an additional dimension of constant 1. Let $\tilde{\mathbf{x}} = \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$, $\tilde{\mathbf{w}}_i = \begin{bmatrix} \mathbf{w}_i \\ b_i \end{bmatrix}$, then $a_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + b_i = \tilde{\mathbf{w}}_i^T \tilde{\mathbf{x}}$.

This unifies $\nabla_{\mathbf{w}_i} \mathcal{L}$ and $\nabla_{b_i} \mathcal{L}$ into $\nabla_{\tilde{\mathbf{w}}_i} \mathcal{L}$.

⑤. Log likelihood ...

from lecture ... $\text{softmax}_i(\mathbf{x}) = \frac{e^{a_i(\mathbf{x})}}{\sum_{j=1}^c e^{a_j(\mathbf{x})}}$

for $a_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + b_i$ *

$c = \#$ of classes

if we let $\Theta = \{\mathbf{w}_i, b_i\}_{i=1, \dots, c}$, then $\text{softmax}_i(\mathbf{x})$ can be interpreted as the probability that belongs to class i

$\Pr(y^{(j)} = i \mid \mathbf{x}^{(j)}, \theta) = \text{softmax}_i(\mathbf{x}^{(j)})$

$\rightarrow \mathcal{L} = \log \prod_{i=1}^m \Pr(y^{(i)} \mid \mathbf{x}^{(i)}, \theta)$

$$= \log \prod_{i=1}^m \text{softmax}_{y^{(i)}}(x^{(i)})$$

$$= \log \prod_{i=1}^m \left[\frac{e^{w_{y^{(i)}}^T x^{(i)} + b_{y^{(i)}}}}{\sum_{j=1}^C e^{w_j^T x^{(i)} + b_j}} \right]$$

$$= \log \prod_{i=1}^m \left[\frac{e^{(a_{y^{(i)}}(x^{(i)}))}}{\sum_{j=1}^C e^{(a_j(x^{(i)}))}} \right]$$

$$= \sum_{i=1}^m \left[a_{y^{(i)}}(x^{(i)}) - \log \sum_{j=1}^C e^{(a_j(x^{(i)}))} \right]$$

Gradient ↴

$$\nabla_{w_i} L = \frac{dL}{dw_i} = \sum_{j=1}^m x^{(j)} (1 \{y^{(j)} = i\} - \text{softmax}_i(x^{(j)}))$$

$$\nabla_{b_i} L = \frac{dL}{db_i} = \sum_{j=1}^m (1 \{y^{(j)} = i\} - \text{softmax}_i(x^{(j)}))$$

$$\text{softmax}_i(x^{(i)}) = \frac{e^{(a_i(x^{(i)}))}}{\sum_{k=1}^C e^{(a_k(x^{(i)}))}}$$

4. (10 points) **Hinge loss gradient.**

Owing to the drastic changes in climate throughout the world, a weather forecasting organization wants our help to build a model that can classify the observed weather patterns as severe or not severe. They have accumulated data on various attributes of the weather pattern such as temperature, precipitation, humidity, wind speed, air pressure, and geographical location along with severity of weather. However, the contribution of the attributes to the classification of weather as severe or not is unknown.

We choose to use a binary support vector machine (SVM) classification model. The SVM model parameters are learned by optimizing a hinge loss. The company has provided us with a data-set

$$\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(K)}, y^{(K)})\}$$

where $\mathbf{x}^{(i)} \in \mathbb{R}^d$ is a feature vector of the i^{th} data sample and $y^{(i)} \in \{-1, 1\}$. We define the hinge loss per training sample as

$$\text{hinge}_{y^{(i)}}(\mathbf{x}^{(i)}) = \max(0, 1 - y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b)) \quad (1)$$

, where $\mathbf{w} \in \mathbb{R}^d$ and bias $b \in \mathbb{R}$ are the model parameters. With the hinge loss per sample defined, we can then formulate the average loss for our model as:

$$\mathcal{L}(\mathbf{w}, b) = \frac{1}{K} \sum_{i=1}^K \text{hinge}_{y^{(i)}}(\mathbf{x}^{(i)}) \quad (2)$$

Find the gradient of the loss function $\mathcal{L}(\mathbf{w}, b)$ with respect to the parameters i.e $\nabla_{\mathbf{w}} \mathcal{L}$ and $\nabla_b \mathcal{L}$.

Hint: An Indicator function, also known as a characteristic function, takes on the value of 1 at certain designated points and 0 at all other points. Mathematically, we can represent it as follows:

$$\mathbb{1}_{\{p < 1\}} = \begin{cases} 1, & \text{if } p < 1 \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

$$\begin{aligned} \text{hinge}_{y^{(i)}}(\mathbf{x}^{(i)}) &= \max(0, 1 - y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b)) \\ \mathcal{L}(\mathbf{w}, b) &= \frac{1}{K} \sum_{i=1}^K \text{hinge}_{y^{(i)}}(\mathbf{x}^{(i)}) \end{aligned}$$

$$\left[\begin{aligned} \nabla_{\mathbf{w}} \mathcal{L} &= \frac{1}{K} \sum_{i=1}^K \nabla_{\mathbf{w}} \max(0, 1 - y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b)) \\ \text{when } 1 - y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) &> 0 \\ \text{when } 1 - y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) &\leq 0, \mathcal{L} = 0 \end{aligned} \right]$$

$\nabla_{\omega} L$ follow same process as above \uparrow

$$\nabla_{\omega} L = \frac{1}{K} \sum_{i=1}^K -y^{(i)} x^{(i)} \mid \sum 1 - y^{(i)} (\omega^T x^{(i)} + b) > 0 \mid$$

$$\nabla_b L = \frac{1}{K} \sum_{i=1}^K -y^{(i)} \mid \sum 1 - y^{(i)} (\omega^T x^{(i)} + b) > 0 \mid$$