

CSCI 416/516. Linear Regression & optimization

dimensionality

$$y = f(x) = \sum_j w_j x_j + b, \text{ whereas } x = (x_1, x_2, \dots, x_j, \dots, x_D) \in \mathbb{R}^D$$

$$y \text{ is linear in } x: y = wx + b; y = w^T x + b.$$

the decision boundary is linear; could be a line (when only 1 feature)
or a hyperplane (when there are D features)

Loss Function: $L(y, t) = \frac{1}{2} (y - t)^2$
Coefficient
residual

Cost Function: $J(w, b) = \frac{1}{N} \sum_{i=1}^N L(y_i, t_i) = \frac{1}{2N} \sum_{i=1}^N (y_i - t_i)^2$
 $= \frac{1}{2N} \sum_{i=1}^N (w^T x_i + b - t_i)^2 = \frac{1}{2N} \sum_{i=1}^N \left(\sum_{j=1}^D (w_j x_j + b) - t_i \right)^2$

$$w = (w_1, w_2, \dots, w_j, \dots, w_D); x = (x_1, \dots, x_j, \dots, x_D); y = w^T x + b$$

import numpy as np

$$\Rightarrow y = \text{np.dot}(w, x) + b$$

one feature across all training sample.

$$X = \begin{bmatrix} x^{(1)T} \\ x^{(2)T} \\ x^{(3)T} \end{bmatrix} = \begin{bmatrix} 8 & 0 & 3 & 0 \\ 6 & 1 & 5 & 3 \\ 2 & 5 & -2 & 8 \end{bmatrix} \rightarrow x^{(1)T}, \text{ one training example}$$

entire set of samples

Vectorization: able to calculate the predictions for every single sample in "one go" (matrix manipulation)

$$w^T X + b1 = \begin{bmatrix} w^T x^{(1)} + b \\ \vdots \\ w^T x^{(N)} + b \end{bmatrix} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix} = y$$

prediction(s) of the whole set of samples, represented by $X = (x^{(1)T} \dots x^{(N)T})$

$$\sum_{i=1}^N (y^{(i)} - t^{(i)})^2 = \|y - t\|^2$$

$$\mathcal{J} = \frac{1}{2N} \sum_{i=1}^N (y^{(i)} - t^{(i)})^2 = \frac{1}{2N} \|y - t\|^2$$

originally, $y = w^T X + b$; but we can integrate the bias b into the weight ^{matrix}.

$$X = \begin{bmatrix} 1 & x^{(1)T} \\ 1 & x^{(2)T} \\ \vdots & \vdots \\ 1 & x^{(N)T} \end{bmatrix} \in \mathbb{R}^{N \times (D+1)} \quad \text{together, } X \text{ has the shape of } \mathbb{R}^{N \times (D+1)}$$

with the dimensionality of 1.
 with the dimensionality of D

$$w = \begin{bmatrix} b \\ w_1 \\ w_2 \\ w_3 \end{bmatrix} \in \mathbb{R}^{D+1}$$

After the cost function is defined, there are 2 ways to minimize it.

- algebraic: z^* minimizes $f(z) \Rightarrow \forall z \neq z^*, f(z) > f(z^*)$
- Calculus: find the global minimum for $f(z^*)$ given z^*

$$\frac{\partial y}{\partial w_j} = \frac{\partial}{\partial w_j} (w^T x + b) = \frac{\partial}{\partial w_j} \left(\sum_{j=0}^D w_j x_j + b \right)$$

$$= \frac{\partial}{\partial w_j} (w_0 x_0 + \dots + w_j x_j + \dots + w_D x_D + b) = x_j$$

$$\frac{\partial y}{\partial b} = \frac{\partial}{\partial b} (w^T x + b) = \frac{\partial}{\partial b} (w_0 x_0 + \dots + w_j x_j + \dots + w_D x_D + b) = 1$$

$$\begin{aligned} \text{Chain Rule: } \frac{\partial \mathcal{L}}{\partial w_j} &= \frac{\partial \mathcal{L}}{\partial y} \cdot \frac{\partial y}{\partial w_j} = \frac{\partial}{\partial y} \left(\frac{1}{2} (y - t)^2 \right) \cdot \frac{\partial y}{\partial w_j} \\ &= (y - t) \cdot x_j \end{aligned}$$

$$\frac{\partial f}{\partial b} = \frac{\partial f}{\partial y} \cdot \frac{\partial y}{\partial b} = \frac{\partial}{\partial y} \left(\frac{1}{2} (y - t)^2 \right) \cdot \frac{\partial y}{\partial b} = (y - t) \cdot 1 = y - t$$

to find the critical pts with respect to w_j and b our parameters

$$\frac{\partial \mathcal{J}}{\partial w_j} = \frac{1}{N} \sum_i^N \frac{\partial \mathcal{L}_i}{\partial w_j} = \frac{1}{N} \sum_i^N (y_i - t_i) x_j^i = 0 \quad \text{critical pt, where the}$$

$$\frac{\partial \mathcal{J}}{\partial b} = \frac{1}{N} \sum_i^N \frac{\partial \mathcal{L}_i}{\partial b} = \frac{1}{N} \sum_i^N (y_i - t_i) \cdot 1 = 0 \quad \text{derivative is 0.}$$

Definition for gradient $\nabla = \nabla f(w) = \left(\frac{\partial f(w)}{\partial w_1}, \dots, \frac{\partial f(w)}{\partial w_D} \right)^T$

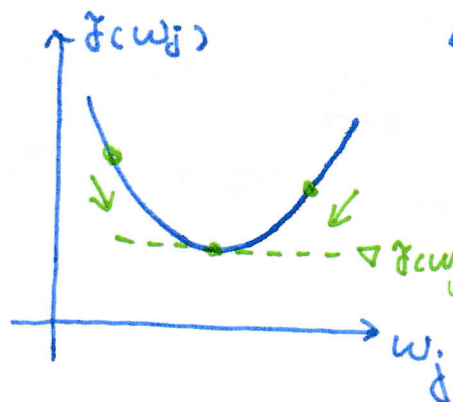
the gradient of the first gradient: $\nabla^2 f(w) = \frac{\partial^2 f(w)}{\partial w_{i,j}^2}$

L^2 or ℓ^2 regularisation: $R(w) = \frac{1}{2} \|w\|_2^2 = \frac{1}{2} \sum_j w_j^2$

$$\begin{aligned} \mathcal{J}_{\text{reg}} &= \mathcal{J}(w) + \lambda R(w) = \mathcal{J}(w) + \frac{\lambda}{2} \sum_j w_j^2 = \frac{1}{N} \sum_i^N \mathcal{L}_i(w) + \frac{\lambda}{2} \sum_j w_j^2 \\ &= \frac{1}{2N} \sum_i^N (y_i - t_i)^2 + \frac{\lambda}{2} \sum_j w_j^2 = \frac{1}{2N} \sum_i^N (w^T x_i - t_i)^2 + \frac{\lambda}{2} \sum_j w_j^2 \end{aligned}$$

λ is the hyperparameter here, which is different from the parameters w_j and b .

Gradient Descent.



we have to initialize the parameters, in this figure, w_j , somewhere on the axis. After the

initialization, we want to guide w_j 's change

so that it (gradually) goes to the "global

minimum, where $\nabla \mathcal{J}(w_j) = 0$.

$\frac{\partial \tilde{J}}{\partial w_j} > 0 \Rightarrow \text{change of } \tilde{J} \text{ corresponds to: increasing } w_j \Rightarrow \text{increasing } \tilde{J}$

$\frac{\partial \tilde{J}}{\partial w_j} < 0 \Rightarrow \text{increasing } w_j \Rightarrow \text{decreasing } \tilde{J}$

gradient descent: always moving against the direction of the gradient

$w_j \leftarrow w_j - \alpha \cdot \frac{\partial \tilde{J}}{\partial w_j}$, whereas α (alpha) is the learning rate.

which is also a hyperparameter.

$$\nabla_w \tilde{J} = \frac{\partial \tilde{J}}{\partial w} = \left(\frac{\partial \tilde{J}}{\partial w_1}, \dots, \frac{\partial \tilde{J}}{\partial w_D} \right)$$

$$\frac{\partial \tilde{J}}{\partial w_j} = \frac{1}{N} \sum_i^N \frac{\partial L_i}{\partial w_j} = \frac{1}{N} \sum_i^N (y^i - t^i) x^i$$

$$w_j \leftarrow w_j - \alpha \cdot \frac{1}{N} \sum_i^N (y^i - t^i) x_j^i; \quad w \leftarrow w - \frac{\alpha}{N} \sum_i^N (y^i - t^i) x^i$$

Gradient Descent under L_2 norm,

$$w \leftarrow w - \alpha \frac{\partial}{\partial w} (\tilde{J} + \lambda R(w))$$

$$w \leftarrow w - \alpha \left(\frac{\partial \tilde{J}}{\partial w} + \lambda \frac{\partial R(w)}{\partial w} \right)$$

$$w \leftarrow w - \alpha \left(\frac{\partial \tilde{J}}{\partial w} + \frac{\partial \lambda}{\partial w} \frac{1}{2} \|w\|_2^2 \right)$$

$$w \leftarrow w - \alpha \left(\frac{\partial \tilde{J}}{\partial w} + \lambda w \right)$$

$$w \leftarrow w - \alpha \cdot \frac{\partial \tilde{J}}{\partial w} - \alpha \cdot \lambda \cdot w$$

$$w \leftarrow w (1 - \alpha \cdot \lambda) - \frac{\partial \tilde{J}}{\partial w} \cdot \alpha$$

$$\frac{\partial \lambda}{\partial w} \frac{1}{2} \|w\|_2^2$$

$$= \frac{\partial \lambda}{\partial w} \frac{1}{2} \sum w_j^2$$

$$= \left(\frac{\partial \lambda}{\partial w_1} \cdot \frac{1}{2} w_1^2, \frac{\partial \lambda}{\partial w_2} \cdot \frac{1}{2} w_2^2 \right.$$

$$\left. , \dots, \frac{\partial \lambda}{\partial w_D} \cdot \frac{1}{2} w_D^2 \right)$$

$$= (\lambda w_1, \lambda w_2, \dots, \lambda w_D)$$

$$= \lambda w$$

Stochastic gradient descent

$J(\theta) = \frac{1}{N} \sum_{i=1}^N L_i = \frac{1}{N} \sum_{i=1}^N L(y(x_i, \theta), t_i)$ whereas L takes t_i and y which takes x_i and θ for computation

instead of $\theta \leftarrow \theta - 2 \cdot \frac{\partial J}{\partial \theta}$, we update the parameter vector $\theta = (w, b)$ every time when we train on each sample in the training dataset.