

Contents

lec 1 : Linear Regression: I	1
1.1 notation and expression	1
1.2 loss function: sum of squared error	2
1.3 gradient of sum of squared error	3
1.4 batch v.s. stochastic GD	4
lec 2 : Linear Regression: II	5
2.1 vectorization	5
2.2 closed-form solution	5

Lec 1 Linear Regression: I

Def 1.1 (supervised learning)

Supervise Learning: Given data X in feature space and labels Y , learn to predict Y from given X .



Label: 可以是 discrete 的 or continuous 的.

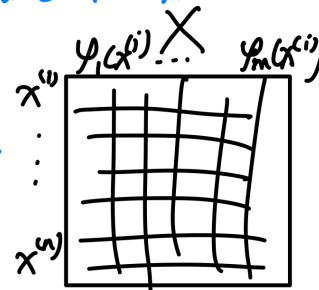
对于 **discrete** 的 label, 这类问题称为 **classification**.

对于 **continuous** 的 label, 这类问题称为 **regression**.

1.1 notation and expression

我们使用以下 notation:

Notation $x \in \mathbb{R}^d$: data
 $\varphi(x) \in \mathbb{R}^m$: features for x
 $\varphi_j(x) \in \mathbb{R}$: the j th feature for x
 $y \in \mathbb{R}$: ctn label
 $x^{(n)}$: the n th training example
 $y^{(n)}$: the n th training label



Def 1.2 ((generalized) linear regression)

给定 N 个 data points $\{(x^{(n)}, y^{(n)})\}_{n=1, \dots, N}$ where each $x^{(n)} \in \mathbb{R}^d, y^{(n)} \in \mathbb{R}$,
以及预先设定好的 M 个 basis functions $\{\phi_i(x)\}_{i=1, \dots, M}$ 用以表示 M 个 features,

我们通过建立一个 $h(x, w) : \mathbb{R}^d \times \mathbb{R}^M \rightarrow \mathbb{R} = \sum_{i=1}^{M-1} w_i \phi_i(x)$, 使其关于 w 线性, 以找到一组参数 $w \in \mathbb{R}^M$, 使得 $h(x^{(n)}, w)$ 能够近似 $y^{(n)}$ for each n , with respect to the loss function we define to measure the distance between two vectors.



1-d case $D=1$

$$\{x^{(1)}, \dots, x^{(N)}\}, \{y^{(1)}, \dots, y^{(N)}\}$$

$$\text{We want: } h(x; w) \approx y$$

\uparrow \uparrow
 input parameter

general case

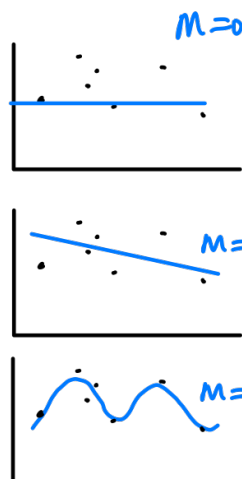
bias term

$$h(x; w) = w_0 + \sum_{j=1}^{M-1} w_j \varphi_j(x)$$

$$= w^T \varphi(x) \rightarrow \varphi_0(x) = 1 \quad (= w \cdot \varphi(x))$$

$$w = (w_0, \dots, w_{M-1})^T, \quad \varphi(x) = (1, \varphi_1(x), \dots, \varphi_{M-1}(x))^T$$

$$w = \begin{bmatrix} w_0 \\ \vdots \\ w_{M-1} \end{bmatrix}, \quad \varphi(x) = \begin{bmatrix} 1 \\ \varphi_1(x) \\ \vdots \\ \varphi_{M-1}(x) \end{bmatrix}$$



Remark 注意: linear regression 指的是 $y \in \mathbb{R}$ 和参数 $w \in \mathbb{R}^M$ 之间是 linear 的, 而不是说 y 和 input x 之间是 linear 的. 我们可以选择 nonlinear 的 basis functions 来 encode x 来表示 features 的特性, 比如我们可以选择:

or $\varphi_j(x) = x^j$ (polynomial)

$\varphi_j(x) = \exp\left(-\frac{(x-\mu_j)^2}{2s^2}\right)$ (Gaussian)

$\varphi_j(x) = \frac{1}{1 + \exp\left(-\frac{x-\mu_j}{s}\right)}$ (Sigmoid)

Hyperparameter s

Diagrams illustrating the Gaussian and Sigmoid basis functions. The Gaussian plot shows three bell curves centered at μ_1, μ_2, μ_3 with width parameter s . The Sigmoid plot shows three S-shaped curves passing through different points, with width parameter s .

1.2 loss function: sum of squared error

这个 loss function 衡量两个 vectors 之间的距离, 目的是衡量 $y \in \mathbb{R}^N$ 和 $h(x, w) \in \mathbb{R}^N$ 这两个 vectors 的差距. 实际上就是它们 difference 的 L_2 -norm 的平方.

We use: sum of squares error

$$E(w) = \frac{1}{2} \sum_{n=1}^N \|h(x^{(n)}, w) - y^{(n)}\|_{L_2}^2$$

$$= \frac{1}{2} \sum_{n=1}^N \left(\sum_{i=0}^{m-1} w_i \varphi_i(x^{(n)}) - y^{(n)} \right)^2$$



$$\nabla E(w) = \begin{bmatrix} \frac{\partial E}{\partial w_0}(w) \\ \vdots \\ \frac{\partial E}{\partial w_m}(w) \end{bmatrix}$$

$$\text{where } \frac{\partial E}{\partial w_k}(w) = \frac{\partial}{\partial w_k} \left[\frac{1}{2} \sum_{n=1}^N \left(\sum_{i=0}^{m-1} w_i \varphi_i(x^{(n)}) - y^{(n)} \right)^2 \right]$$

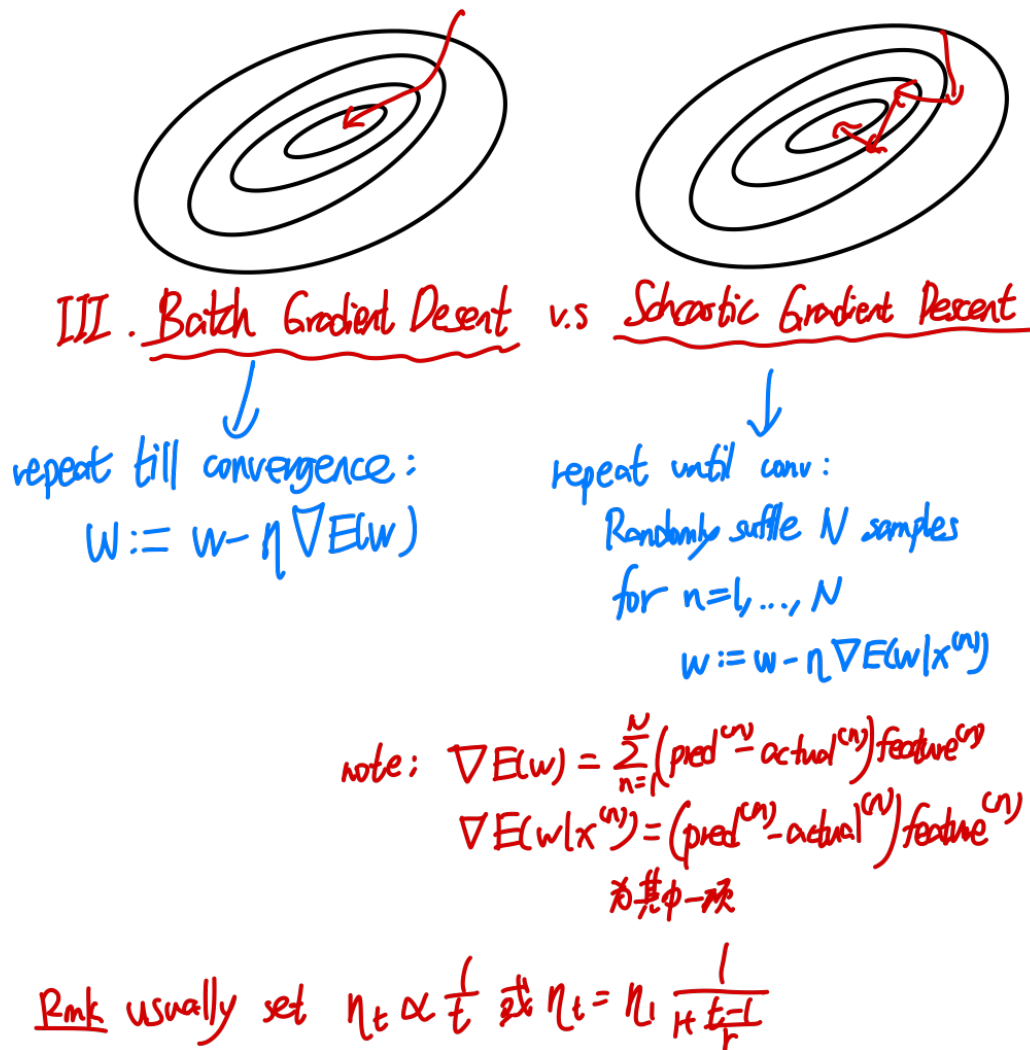
1.3 gradient of sum of squared error

我们下面首先通过求 $\nabla E(w)$ 的每个 entry $\frac{\partial E}{\partial w_k}(w)$ 来写出这个 gradient.

$$\begin{aligned} \frac{\partial E}{\partial w_k}(w) &= \frac{\partial}{\partial w_k} \left[\frac{1}{2} \sum_{n=1}^N \left(\sum_{i=0}^{m-1} w_i \varphi_i(x^{(n)}) - y^{(n)} \right)^2 \right] \\ &= \frac{1}{2} \sum_{n=1}^N \frac{\partial}{\partial w_k} \left(\underbrace{w_0 \varphi_0(x^{(n)})}_{\text{or } \frac{d}{dw_k}(aw_k+b)^2 = 2a(aw_k+b)} + \dots + \underbrace{w_k \varphi_k(x^{(n)})}_{\text{where } a = \varphi_k(x^{(n)})} + \dots + \underbrace{w_{m-1} \varphi_{m-1}(x^{(n)})}_{\text{const}} - y^{(n)} \right)^2 \\ &= \sum_{n=1}^N \left[\left(\sum_{i=0}^{m-1} w_i \varphi_i(x^{(n)}) - y^{(n)} \right) \varphi_k(x^{(n)}) \right] \\ &= \sum_{n=1}^N \left[\left(\sum_{i=0}^{m-1} w_i \varphi_i(x^{(n)}) - y^{(n)} \right) \varphi_k(x^{(n)}) \right] \\ \text{So } \nabla E(w) &= \sum_{n=1}^N \left\{ \left(\sum_{i=0}^{m-1} w_i \varphi_i(x^{(n)}) - y^{(n)} \right) \begin{bmatrix} \varphi_0(x^{(n)}) \\ \vdots \\ \varphi_{m-1}(x^{(n)}) \end{bmatrix} \right\} \\ &= \sum_{n=1}^N \left(w^T \varphi(x^{(n)}) - y^{(n)} \right) \varphi(x^{(n)}) \\ &= \sum_{n=1}^N \left(\text{pred}^{(n)} - \text{actual}^{(n)} \right) \text{feature vector}^{(n)} \end{aligned}$$

1.4 batch v.s. stochastic GD

我们通过迭代降低 gradient 来降低 loss function 的值, 从而优化 weight vector.



More practically, 我们可以采用 minibatch SGD: 即在 batch GD 和 SGD 之间, 每次选择一小部分 samples, 称为一个 **minibatch**, 在这个 minibatch 上进行 GD.

Lec 2 Linear Regression: II

2.1 vectorization

我们可以把每个 $x^{(n)}$ 的 features 写成一个 row vector, 并 stack up N 个 row vectors, 成为一个 $N \times M$ 的 matrix Φ . 从而:

$$h(x, w) = \Phi w$$

vectorization 的好处是: 1. 便于手算; 2. computer 可以进行并行计算.

let $\Phi = \begin{pmatrix} \phi_0(x^{(1)}) & \dots & \phi_{M-1}(x^{(1)}) \\ \vdots & & \vdots \\ \phi_0(x^{(N)}) & \dots & \phi_{M-1}(x^{(N)}) \end{pmatrix} \Rightarrow h(x, w) = \underbrace{\Phi}_{\in \mathbb{R}^{N \times M}} \underbrace{w}_{\in \mathbb{R}^M}$

$E(w) = \frac{1}{2} \sum_{n=1}^N (w^T \phi(x^{(n)}) - y^{(n)})^2$

\downarrow
 $E: \mathbb{R}^M \rightarrow \mathbb{R}$

$= \frac{1}{2} \|\Phi w - y\|^2$

$= \frac{1}{2} (\Phi w - y)^T (\Phi w - y)$

$= \frac{1}{2} (w^T \Phi^T - y^T) (\Phi w - y)$

$= \frac{1}{2} (w^T \Phi^T \Phi w - \underbrace{w^T \Phi^T y}_{\text{the same, scalar}} - y^T \Phi w + y^T y)$

$= \frac{1}{2} w^T \Phi^T \Phi w - w^T \Phi^T y + \frac{1}{2} y^T y$

(then recall chain rule to differentiate $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$)

计算得 linear regression 的 loss function 为:

$$E(w) = \frac{1}{2} w^T \Phi^T \Phi w - w^T \Phi^T y + \frac{1}{2} y^T y$$

2.2 closed-form solution

如果

$$\nabla E(w) = y$$

有一个 closed form solution, 那么这个 solution 一定是一个 local min/max, 从而 possibly 成为一个 global min.

为了计算 closed form solution, 我们首先要给出 $\nabla E(w)$ 的 matrix form 表达式.

这里首先引入 linear form 和 quadratic form ($\mathbb{R}^m \rightarrow \mathbb{R}$) 的 gradient:

linear form:
 $f(x) = b^T x \Rightarrow \nabla_x f(x) = \underline{b}$ since $\frac{\partial f(x)}{\partial x_k} = b_k$

quadratic form: $(= \sum_{i,j=1}^n x_i A_{ij} x_j)$
 $f(x) = x^T A x \Rightarrow \nabla_x f(x) = \underline{2Ax}$ since $\frac{\partial f}{\partial x_k}(x) = 2 \sum_{j=1}^n A_{kj} x_j = 2(Ax)_k$

我们发现: $E(w)$ 就是一个 w 的 quadratic form, 一个 w 的 linear form 和一个 const 的组合. 从而可以求出:

$$\begin{aligned} \nabla E(w) &= \frac{1}{2} \nabla (\underbrace{w^T \phi^T \phi w}_{\substack{\text{quadratic form} \\ w^T (\phi^T \phi) w \\ \xrightarrow{\text{diff}} 2\phi^T \phi}}) - \nabla (\underbrace{w^T \phi^T y}_{\substack{\text{linear form} \\ w^T (\phi^T y) \\ = w \cdot (\phi^T y) \\ = (\phi^T y)^T w \\ \xrightarrow{\text{diff}} \phi^T y}}) + \frac{1}{2} \nabla (\underbrace{y^T y}_{\substack{\text{const} \\ \text{diff: 0}}}) \\ &= \phi^T \phi w - \phi^T y \\ &= \phi^T (\underbrace{\phi w - y}_{\hookrightarrow h(x, w)}) = \phi^T (\text{pred}_y - \text{actual}_y) \end{aligned}$$

从而我们得到 closed form solution (if exists):

$$\begin{aligned} \nabla E(w) &= \phi^T (\phi w - y) := 0 \\ \Rightarrow \phi^T \phi w &= \phi^T y \\ \Rightarrow w_{ML} &= (\phi^T \phi)^{-1} \phi^T y \\ &\quad (\text{if exists}) \end{aligned}$$

因而 closed form exists iff $\Phi^T \Phi$ 可逆, iff Φ 可逆.

并且 recall in linear algebra: $\text{rank}(\Phi^T \Phi) = \text{rank}(\Phi)$. 因而, **closed form exists iff** $M \geq N$ 且 $\text{rank}(\Phi) = N$.