

1-Linear Regression.

$$\hat{y} = w_1 x_1 + w_2 x_2 + w_d x_d + b$$

\vec{w} : weights

\vec{x} : features

\hat{y} : label predicted.

b : bias (offset)

$$\Rightarrow \text{vector form: } \hat{y} = \vec{w}^T \vec{x} + b$$

Now we have n examples.

也就是: $n \times d$ 的维度

$$\begin{matrix} n \\ \left[\begin{array}{c} \vec{x}^{(1)} \\ \vec{x}^{(2)} \\ \vdots \\ \vec{x}^{(n)} \end{array} \right] \\ d \end{matrix}$$

$$\text{So } X_{n \times d} \vec{w}_d = \begin{bmatrix} \vec{x}^{(1)} \cdot \vec{w} \\ \vdots \\ \vec{x}^{(n)} \cdot \vec{w} \end{bmatrix}$$

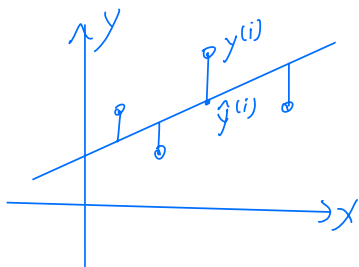
$$\text{即: } \vec{w}^T \vec{x}^{(i)} + b = \begin{bmatrix} w_1 \\ \vdots \\ w_d \\ b \end{bmatrix} \cdot [\begin{matrix} x_1 & x_2 & \dots & x_d & 1 \end{matrix}]$$

于是把这样的每个 $\vec{x}^{(i)} = [x_1 \ x_2 \ \dots \ x_d \ 1]$ 的 vector

$$\text{组成 matrix } X = \begin{bmatrix} \vec{x}^{(1)} \\ \vec{x}^{(2)} \\ \vdots \\ \vec{x}^{(n)} \end{bmatrix}$$

我们可以得到 $L(w, b)$ 的 matrix form:

$$L(\vec{w}) = \frac{1}{2} \|\vec{y} - X\vec{w}\|^2$$



1.1. Matrix form

\Rightarrow 预测结果:

$$\vec{\hat{y}} = X \cdot \vec{w} + b$$

$$\begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_d^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(n)} & x_2^{(n)} & \dots & x_d^{(n)} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} + \begin{bmatrix} b \\ \vdots \\ b \end{bmatrix}$$

我们用一个 loss function 来 measure 这 n 个 examples (data points)

的 performance:

用 squared error 作为 regression loss (第 i 个)

$$l^{(i)}(w, b) = \frac{1}{2} (\hat{y}^{(i)} - y^{(i)})^2$$

1.2 Matrix form of Loss function

\Rightarrow 整体 n 个 examples 的 performance

$$L(w, b) = \frac{1}{n} \sum_{i=1}^n l^{(i)}(w, b) \\ = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (\vec{w}^T \vec{x}^{(i)} + b - y^{(i)})^2$$

注意到这里的 bias b 可以被 absorb 进 weights, 如果给每个 $\vec{x}^{(i)}$ 的最后一个元素为 1 的 dummy feature

我们希望的参数 w 和 b 是:

Denote ideal (w, b) 为 (w^*, b^*)

$$\Rightarrow (\vec{w}^*, b^*) = \underset{w, b}{\operatorname{argmin}} L(\vec{w}, b)$$

如何 minimize the loss 来获得 (w^*, b^*)

1.3 Analytic Solution

即 $\nabla_w L(\vec{w}, b) = \vec{0}$. 虽然没用, 但推导值得一看.

前已陈述: b 可被 absorb 进 \vec{w} , 使

$L(\vec{w}, b)$ 简化为 $L(\vec{w})$ (这里两个 \vec{w} 不一样, 一个多个 b , 实际就是把 b 放进 \vec{w} 里)

所以我们要求的就是 $L(\vec{w}) = \frac{1}{2} \|\vec{y} - X\vec{w}\|^2$

对于 \vec{w} 的 derivative 为 $\vec{0}$ 的时候的 \vec{w} .

\Rightarrow 即 $\nabla_w L(\vec{w}) = \vec{0}$. 时的 \vec{w} 为 \vec{w}^*

$$\Rightarrow \vec{w}^* = (X^T X)^{-1} X^T \vec{y}$$

如何解得这个解：以下为一个简略的推导。

$$\text{Let: } G(\vec{w}) = y - X\vec{w} \quad (\Rightarrow \nabla_{\vec{w}} G(\vec{w}) = -X^T)$$

$$F(\vec{u}) = \frac{1}{2} \|\vec{u}\|^2 \quad (\Rightarrow \nabla_{\vec{u}} F(\vec{u}) = \vec{u})$$

$$\Rightarrow L(\vec{w}) = F(G(\vec{w}))$$

$$\Rightarrow \nabla_{\vec{w}} L(\vec{w}) = \nabla_{\vec{w}} G(\vec{w}) \cdot \nabla_{\vec{u}} F(G(\vec{w}))$$

$$= (-X^T) \times (\vec{y} - X\vec{w})$$

$$= -X^T \vec{y} + X^T X \vec{w}$$

而后 set $\nabla_{\vec{w}} L(\vec{w}) = 0$ 得 $\vec{w}^* = (X^T X)^{-1} X^T y$

而有这样的 analytic solution 的条件是 $(X^T X)$ invertible.
(which is nearly impossible)

* (由于没学过 multi calc, 再详细推一遍)

$$H(\vec{w}) = \vec{w}^T \vec{x} \quad (= \sum_{i=1}^d w_i x_i)$$

$$= w_1 x_1 + w_2 x_2 + \dots + w_d x_d$$

doesn't depend on w_1 !!

$$\Rightarrow \nabla_{\vec{w}} H(\vec{w}) = \begin{bmatrix} \frac{\partial H}{\partial w_1} \\ \frac{\partial H}{\partial w_2} \\ \vdots \\ \frac{\partial H}{\partial w_d} \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} = X^T$$

$$F(\vec{u}) = \frac{1}{2} \|\vec{u}\|^2 = \frac{1}{2} \sum_{i=1}^d u_i^2$$

$$= \frac{1}{2} w_1^2 + \frac{1}{2} w_2^2 + \dots + \frac{1}{2} w_d^2$$

does not depend on w_1

$$\Rightarrow \frac{\partial F}{\partial w_i} = w_i \Rightarrow \nabla F(\vec{w}) = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} = \vec{w}$$

一些其他重要的 gradient:

$$H(\vec{w}) = \frac{1}{2} \vec{w}^T A \vec{w}$$

经过一些 algebra 可得

$$\nabla H(\vec{w}) = \frac{1}{2} (A + A^T) \vec{w}$$

(In particular, when A is symmetric,
 $\nabla H(\vec{w}) = A \vec{w}$.)

现在回归正题:

1.5 Gradient Descent

既然 closed form solution 很少见 (并且如果有就和 Deep Learning 没关系)

那我们仍需要一个能用数值方法来优化 loss function 的方法。

\Rightarrow 从 negative gradient 的方法 update the parameters, 从而 iteratively reduce loss.

For any objective function $J(\vec{w})$,

Gradient descent's form:

$$\vec{w} \leftarrow \vec{w} - \eta \nabla J(\vec{w})$$

Gradient 给出了 J 的 fastest local increase 的方向。
从而快速 minimize J .

这里的 η controls how fast we move
称为 learning rate.

1.6 如何计算 $\nabla_{\vec{w}} L(\vec{w})$

现在我们已经推出

$$\nabla_{\vec{w}} L(\vec{w}) = -X^T \vec{y} + X^T X \vec{w}$$

我们可去掉 matrix 它更易理解和计算。

在此处 (linear regression 下的 least square sum 为 loss func)

$$\nabla_{\vec{w}} L(\vec{w}) = \frac{1}{n} \sum_{i=1}^n \nabla (\vec{w}^T \vec{x}^{(i)} - y^{(i)})^2$$

$$= \frac{1}{n} \sum_{i=1}^n \vec{x}^{(i)} (\vec{w}^T \vec{x}^{(i)} - y^{(i)})$$

$$\text{因而 GD update 为 } \vec{w} \leftarrow \vec{w} - \frac{\eta}{n} \sum_{i=1}^n \vec{x}^{(i)} (\vec{w}^T \vec{x}^{(i)} - y^{(i)})$$

可用matrix推导:

$$X = \begin{bmatrix} -\vec{x}^{(1)}- \\ -\vec{x}^{(2)}- \\ \vdots \\ -\vec{x}^{(n)}- \end{bmatrix} = \begin{bmatrix} | & | & | \\ \vec{x}_1 & \vec{x}_2 & \dots & \vec{x}_d \\ | & | & | \end{bmatrix}$$

(n examples) (d features)

$$X^T = \begin{bmatrix} -\vec{x}_1^T- \\ -\vec{x}_2^T- \\ \vdots \\ -\vec{x}_d^T- \end{bmatrix} = \begin{bmatrix} | & | & | \\ \vec{x}_1^T & \vec{x}_2^T & \dots & \vec{x}_d^T \\ | & | & | \end{bmatrix}$$

这里: $X^T X: d \times d$ $XX^T: n \times n$
 $d \times n \times n \times d$ $n \times d \times d \times n$

$$(X^T X)^T = X^T (X^T)^T = X^T X \text{ (itself)}$$

$$(XX^T)^T = XX^T \text{ (itself)}$$

我们可以发现 XX^T 和 $X^T X$ 都是 symmetric 的.

$$X^T X = \begin{bmatrix} -\vec{x}_1^T- \\ -\vec{x}_2^T- \\ \vdots \\ -\vec{x}_d^T- \end{bmatrix} \begin{bmatrix} | & | & | \\ \vec{x}_1 & \vec{x}_2 & \dots & \vec{x}_d \\ | & | & | \end{bmatrix}$$

(symmetric)

$$= \begin{bmatrix} \vec{x}_1^T \cdot \vec{x}_1 & \vec{x}_1^T \cdot \vec{x}_2 & \dots & \vec{x}_1^T \cdot \vec{x}_d \\ \vec{x}_2^T \cdot \vec{x}_1 & \vec{x}_2^T \cdot \vec{x}_2 & \dots & \vec{x}_2^T \cdot \vec{x}_d \\ \vdots & \vdots & \ddots & \vdots \\ \vec{x}_d^T \cdot \vec{x}_1 & \vec{x}_d^T \cdot \vec{x}_2 & \dots & \vec{x}_d^T \cdot \vec{x}_d \end{bmatrix}$$

$$XX^T = \begin{bmatrix} -\vec{x}^{(1)}- \\ -\vec{x}^{(2)}- \\ \vdots \\ -\vec{x}^{(n)}- \end{bmatrix} \begin{bmatrix} | & | & | \\ \vec{x}_1^T & \vec{x}_2^T & \dots & \vec{x}_d^T \\ | & | & | \end{bmatrix}$$

$$= \begin{bmatrix} \vec{x}^{(1)} \cdot \vec{x}^{(1)T} & \vec{x}^{(1)} \cdot \vec{x}^{(2)T} & \dots & \vec{x}^{(1)} \cdot \vec{x}^{(n)T} \\ \vec{x}^{(2)} \cdot \vec{x}^{(1)T} & \vec{x}^{(2)} \cdot \vec{x}^{(2)T} & \dots & \vec{x}^{(2)} \cdot \vec{x}^{(n)T} \\ \vdots & \vdots & \ddots & \vdots \\ \vec{x}^{(n)} \cdot \vec{x}^{(1)T} & \vec{x}^{(n)} \cdot \vec{x}^{(2)T} & \dots & \vec{x}^{(n)} \cdot \vec{x}^{(n)T} \end{bmatrix}$$

$$\text{因而 } \nabla_w L(\vec{w}) = \boxed{-X^T \vec{y} + X^T X \vec{w}}$$

$$= \begin{bmatrix} \vec{x}_1^T \cdot \vec{x}_1 & \vec{x}_1^T \cdot \vec{x}_2 & \dots & \vec{x}_1^T \cdot \vec{x}_d \\ \vec{x}_2^T \cdot \vec{x}_1 & \vec{x}_2^T \cdot \vec{x}_2 & \dots & \vec{x}_2^T \cdot \vec{x}_d \\ \vdots & \vdots & \ddots & \vdots \\ \vec{x}_d^T \cdot \vec{x}_1 & \vec{x}_d^T \cdot \vec{x}_2 & \dots & \vec{x}_d^T \cdot \vec{x}_d \end{bmatrix} \cdot \vec{w} - \begin{bmatrix} -\vec{x}_1^T- \\ -\vec{x}_2^T- \\ \vdots \\ -\vec{x}_d^T- \end{bmatrix} \cdot \vec{y}$$

$$= \vec{x}^{(i)} (\vec{w}^T \vec{x}^{(i)} - \vec{y}^{(i)})$$

总结:

$$\nabla_w L(w) = \frac{1}{n} \sum_i \vec{x}^{(i)} (\vec{w}^T \vec{x}^{(i)} - \vec{y}^{(i)})$$

$$\text{因而 GD update 为 } \boxed{w \leftarrow w - \frac{\eta}{n} \sum_i \vec{x}^{(i)} (\vec{w}^T \vec{x}^{(i)} - \vec{y}^{(i)})}$$

1.6 Minibatch Stochastic Gradient Descent.

(小批量随机梯度下降)

每次 iteration, 随机选取一个 minibatch B ,
 consisting of a number of training data points.

$$\text{那么 } \boxed{\vec{w} \leftarrow \vec{w} - \frac{\eta}{|B|} \sum_{i \in B} \vec{x}^{(i)} (\vec{w}^T \vec{x}^{(i)} - \vec{y}^{(i)})}$$

此处的 η 和 $|B|$ 就称为 hyperparameters (超参数).

它们在同一次 training 中为 fixed 的,
 但也可由一个 outer loop 来 optimize
 by tracking performance on a validation set.

(下面就不推了)

可得到从上一轮的 \vec{w} (用 $\vec{w}^{(t)}$ 来表示)

更新到这一次的 \vec{w}

Loss function 值的 iteration 为:

$$\tilde{L}(w; w^{(t)})$$

$$= \underline{L(w^{(t)})} + \underline{\nabla L(w^{(t)}) \cdot (w - w^{(t)})} + \frac{\eta}{2} \|w - w^{(t)}\|^2$$

1.7 用 learned model 进行 prediction

我们通过了用 least square sum 的 loss function (6.2)

对 linear regression 的 weights (\vec{w}, b)

用 gradient descent 求出了使 loss function (6.2)

最小的 weight 记作 $\hat{\vec{w}}, \hat{b}$ 。

$\hat{y} = \hat{\vec{w}}^T \vec{x} + \hat{b}$ 就是我们现在的 learned model。

我们可以把由一组 features x_1, x_2, \dots, x_d

组成的一个 testing data point:

$$\vec{x}^{(i)} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \text{ 代入, 来预测这个}$$

data point $\vec{x}^{(i)}$ 的 estimated target $\hat{y}^{(i)}$ 。

这个过程在 deep learning jargon 中叫做 inference。