

1/11/2024 Linear Regression

Predicting a real number

linear function of input

Features (d total), each house is feature vector $x \in \mathbb{R}^d$

target (y)

Sale price	Area	#bedrooms	has garage?	...	Constant feature
\$500k = $y^{(1)}$	1200	2	0		1
\$700k	2000	3	1		1
	\vdots				

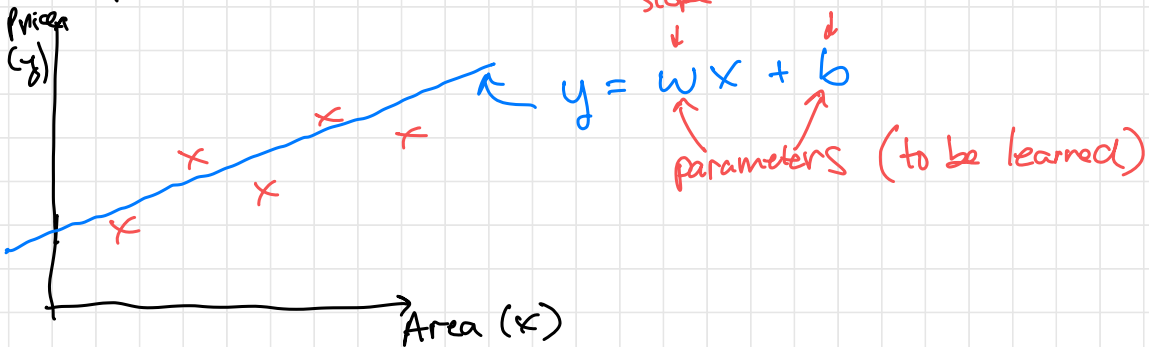
$x^{(1)}$

\downarrow makes param optional

Dataset $D = \{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$

of size n

Simplest case: $d=1$



When $d > 1$:

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

$$= \begin{bmatrix} 1 & x_1 & x_2 & \dots & x_d \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} + b$$

$1 \times d$ $d \times 1$ $= \mathbb{R}$

Parameters:
 $w \in \mathbb{R}^d$
 $b \in \mathbb{R}$

Q: How to choose good w & b ?

A: Define a loss function

$$L(w, b) = \text{[How bad do } w \& b \text{ fit our observed data?]}$$

$$= \frac{1}{n} \sum_{i=1}^n \left(\underbrace{w^T x^{(i)} + b}_{\text{model prediction}} - \underbrace{y^{(i)}}_{\text{true output}} \right)^2$$

Goal: Find w & b that minimize $L(w, b)$
optimization problem?

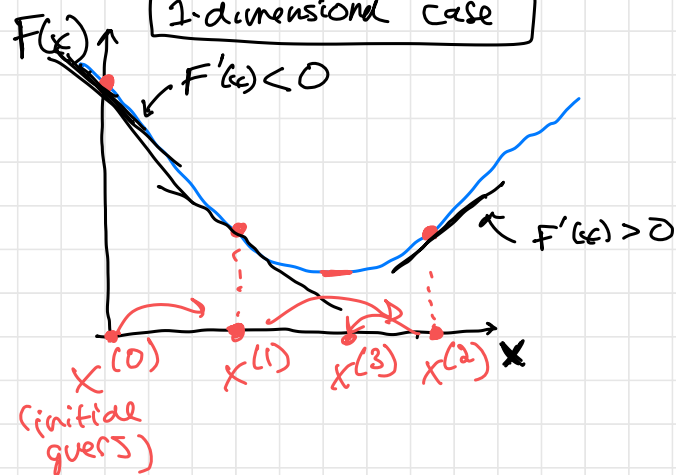
Gradient Descent

General method for minimizing function

Given: Function F from $\mathbb{R}^d \rightarrow \mathbb{R}$, differentiable

Gradient will try to find x that minimizes $F(x)$

1-dimensional case



Have current guess $x^{(t)}$
If $F'(x^{(t)}) < 0$,
increase $x^{(t)}$
to yield $x^{(t+1)}$

If $F'(x^{(t)}) > 0$
decrease $x^{(t)}$
to yield $x^{(t+1)}$

d-dimensional case optimizing w.r.t. $x \in \mathbb{R}^d$

Partial Derivative: $\frac{\partial F}{\partial x_i}$ \leftarrow Take derivative wrt x_i holding all other x_j 's constant

New G.D. Rule:

For each $i = 1, \dots, d$:

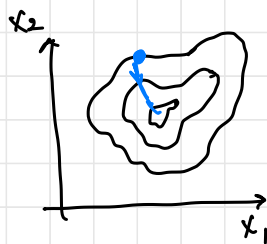
IF $\frac{\partial F}{\partial x_i} \Big|_{x=x^{(t)}} < 0$, increase $x_i^{(t)}$

IF $\frac{\partial F}{\partial x_i} \Big|_{x=x^{(t)}} > 0$, decrease $x_i^{(t)}$

Gradient $\nabla_x F(x) = \left[\frac{\partial F}{\partial x_1}, \frac{\partial F}{\partial x_2}, \dots, \frac{\partial F}{\partial x_d} \right]$

Starting at $x^{(t)}$, best direction to go (to minimize F) is direction of negative gradient

Fact: Negative gradient is direction of steepest descent



Gradient Descent Algorithm:

$x^{(0)} \leftarrow (0, 0, \dots, 0) \in \mathbb{R}^d$

for t in $1, \dots, T$ \leftarrow total # steps

$x^{(t)} \leftarrow x^{(t-1)} - \eta \nabla_x F(x^{(t-1)})$

return $x^{(T)}$

\leftarrow learning rate (e.g. 0.01)

Gradient Descent for Linear Regression

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

$$\nabla_w L(w) = \frac{1}{n} \sum_{i=1}^n 2 \cdot \underbrace{(w^T x^{(i)} - y^{(i)})}_{\text{Scalar}} \cdot \underbrace{x^{(i)}}_{\text{Vector}}$$

$$\frac{d}{dw} 8w = 8$$

GD for Linear Regression:

$$w^{(0)} \leftarrow [0, \dots, 0] \in \mathbb{R}^d$$

for $t = 1, \dots, T$:

$$w^{(t)} \leftarrow w^{(t-1)} - \eta \cdot \frac{1}{n} \sum_{i=1}^n 2 \cdot \underbrace{(w^{(t-1)T} x^{(i)} - y^{(i)})}_{\text{Determining if we add or subtract multiple of } x^{(i)}} \cdot x^{(i)}$$

return $w^{(T)}$

If $w^T x^{(i)} - y^{(i)} > 0$: prediction too large
subtract multiple of $x^{(i)}$ from w
 $\Rightarrow w^T x^{(i)}$ smaller

If $w^T x^{(i)} - y^{(i)} < 0$: prediction too small,
add multiple of $x^{(i)}$
 $\Rightarrow w^T x^{(i)}$ bigger