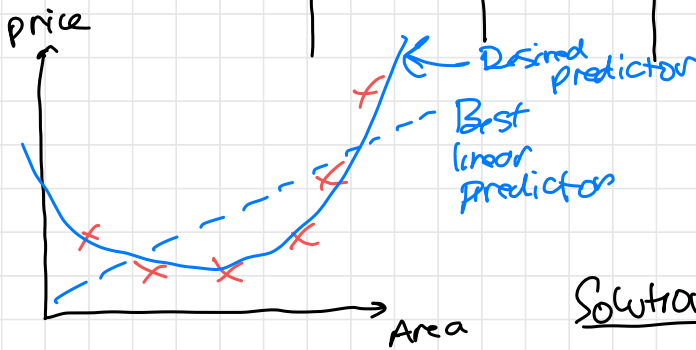


1/16/2024: Linear Regression II

- ① How can we learn more complex functions?
- ② Why does gradient descent work for lin. reg.?
- ③ Why use squared error?

(y)

Sale price	Area - ^{Real number}	#bed ^{Integer}	house type ^{categorical}	Area ² ^{New feature}	Area ³
\$00k	1200	2	Condo	144,000	~

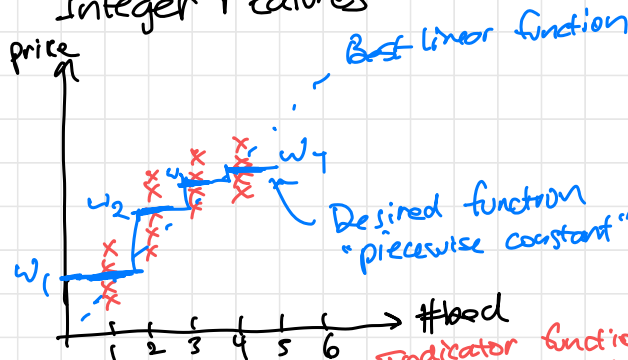


$$\text{prediction} = w_1 \cdot \text{Area} + w_2 \cdot \text{Area}^2 + w_3 \cdot \text{Area}^3$$

Solution: Add more features

Linear regression is linear in the input features
(which we control)

Integer Features



Solution:
Add indicator features!
Some boolean function

y	#bed = 1?	#bed = 2?	#bed = 3?
\$00k	1	0	0
\$00k	0	1	0
i	0	1	0
	0	0	1

Now:

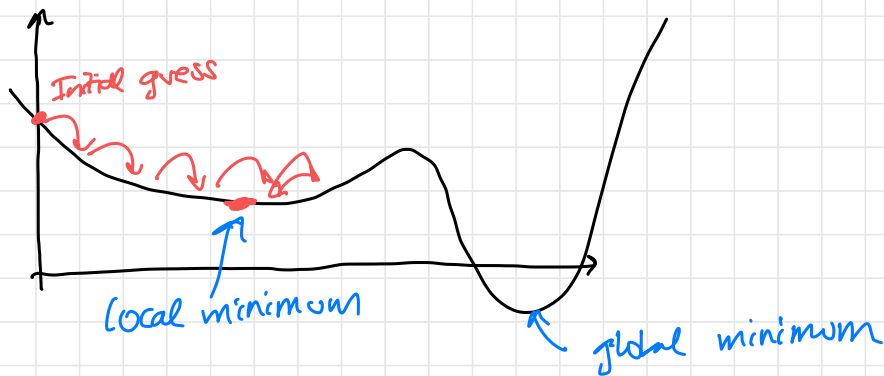
$$\text{prediction} = w_1 \cdot \mathbb{1}[\#bed=1] + w_2 \cdot \mathbb{1}[\#bed=2] + \dots$$

For categorical features, use some idea, e.g. $\mathbb{1}[\text{house type} = \text{"condo"}]$

"Feature Engineering": Process of choosing features to use

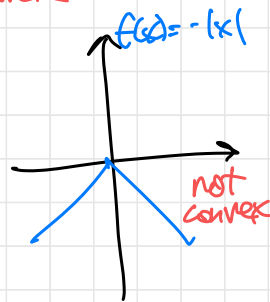
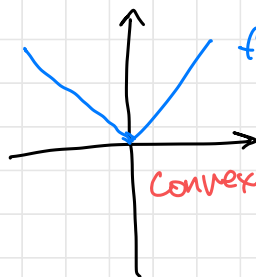
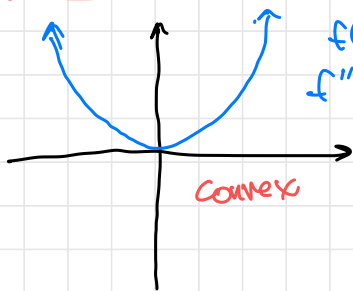
Why does gradient descent work?

Answer: Linear regression is convex



- ① Linear regression loss function $L(w)$ is convex
- ② For all convex functions, every local minimum is a global minimum

Def 1: $f(x)$ is convex $\Leftrightarrow f''(x) \geq 0$ everywhere
But this only holds when f'' exists everywhere



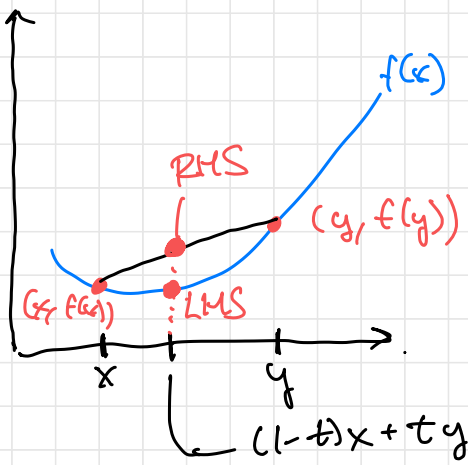
Def 2 (informally): convex functions "hold water"

Def 3 (formal): A function f is convex iff
For every x, y in its domain
and every $t \in [0, 1]$,

$$f((1-t)x + ty) \leq (1-t)f(x) + tf(y)$$

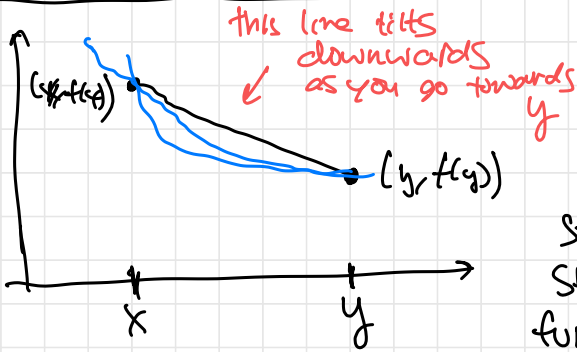
"go $t\%$ of the way from x to y "

"go $t\%$ of way from $f(x)$ to $f(y)$ "



TLDAR:
 If you draw line
 between $(x, f(x))$
 and $(y, f(y))$,
 the line must be
 above the function

All local minima of convex functions are global minima



Let y be global min
 x be any other point
 where $f(x) > f(y)$

Starting from x , if you
 step in direction of y ,
 function must go down.

therefore, x is not local minimum

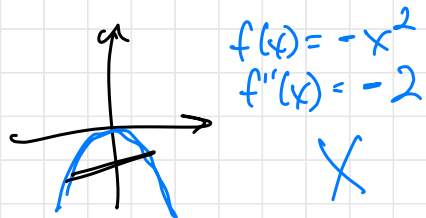
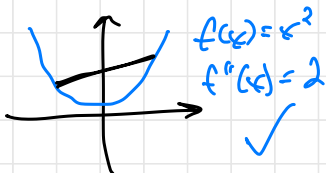
Linear regression is convex

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

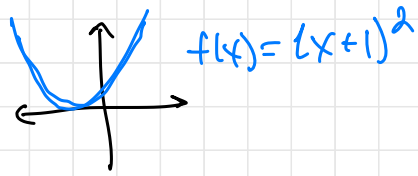
Rules:

① If $f: \mathbb{R} \rightarrow \mathbb{R}$ and $f''(x)$ exists everywhere
 and $f''(x) \geq 0$ everywhere

then $f(x)$ is convex



② If f is convex, then $g(x) = f(Ax + b)$ is convex for any A, b constants



③ If $f(x)$ and $g(x)$ are convex then so is $f(x) + g(x)$

④ If f is convex and C is a constant ≥ 0 then $C \cdot f(x)$ is convex

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

• $f(x) = x^2$ is convex by ①

• $(w^T x^{(i)} - y^{(i)})^2$ is convex by ②

↑ parameter ↑ constants

• $\sum_{i=1}^n (w^T x^{(i)} - y^{(i)})^2$ is convex by ③

• $\frac{1}{n} \sum_{i=1}^n (w^T x^{(i)} - y^{(i)})^2$ is convex by ④

Maximum Likelihood Estimation (MLE)

- Poit probabilistic process that generated the data
- Choose parameters to make observed data most likely

E.g. coin flips

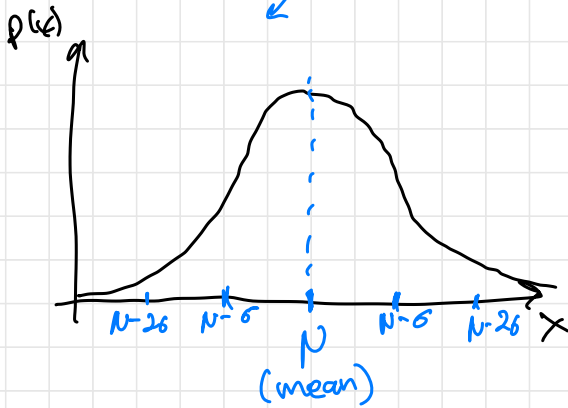
Observe data = [H, T, H, H, H]

observed data

unknown p = probability of heads ← unknown parameter

Goal: choose p that makes data most likely ← "learning"

Linear Regression: Assume $y^{(i)}$ is drawn from
Gaussian distribution w/ mean $w^T x^{(i)}$, variance σ^2
 Constant



determined by
 "true" value of w

$$p(x; \mu, \sigma^2) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Likelihood of data:

$$\begin{aligned} \mathcal{L}(w) &= \prod_{i=1}^n p(y^{(i)} | x^{(i)}; w) \\ &= \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(y^{(i)} - w^T x^{(i)})^2}{2\sigma^2}\right) \end{aligned}$$

"parameterized by w "

Trick: Take the log (increasing function)

$$\begin{aligned} \log \mathcal{L}(w) &= \sum_{i=1}^n \underbrace{\log\left(\frac{1}{\sigma \sqrt{2\pi}}\right)}_{\text{constant}} + \underbrace{-\frac{(y^{(i)} - w^T x^{(i)})^2}{2\sigma^2}}_{\text{constant}} \\ &= \text{constant} - \frac{1}{2\sigma^2} \cdot \sum_{i=1}^n (y^{(i)} - w^T x^{(i)})^2 \end{aligned}$$

maximizing $\log \mathcal{L}(w)$ is same as
 minimizing original $L(w)$ from linear regression