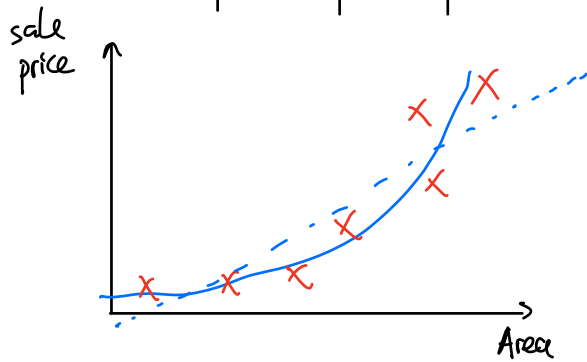


1. Features
2. Convexity
3. Closed-form solution

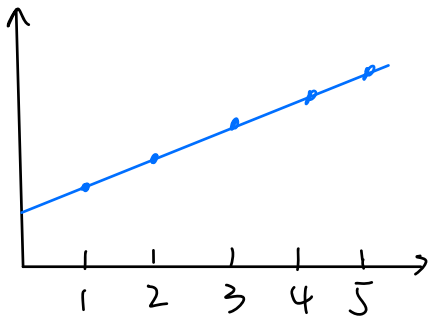
Featurization

(y)	Area	# bed	house type	Area ²	Area ³
sale price	500	1	condo	250000	
	1000	2	town house	1000000	



$$\hat{y} = w_1 \text{Area} + w_2 \text{Area}^2 + w_3 \text{Area}^3 + \dots$$

Linear regression is linear in the features.



indicator features

y	#bed = 1	bed = 2	= 3	≥ 4
	1	0	0	0
	0	1	0	0

$$I\{\text{true}\} = 1$$

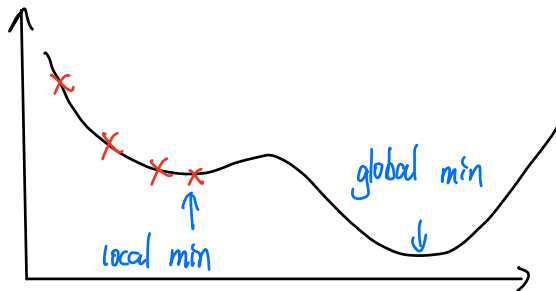
$$I\{\text{false}\} = 0$$

$$\hat{y} = w_1 I\{\text{bed} = 1\} + w_2 I\{\text{bed} = 2\} + \dots$$

zip code : naively 10^5 features

“Feature engineering”: art of choosing features to use
→ $I\{\text{zip code is in LA}\}$
allows sharing information between nearby zip codes

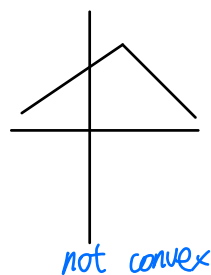
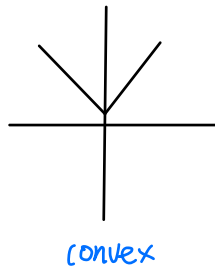
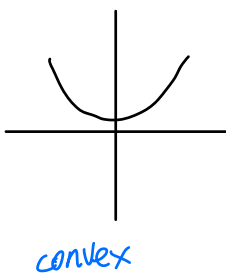
Convexity Why does gradient descent work?



① Linear regression is a convex problem, $L(w)$ is a convex function.

② For a convex function, any local minimum is a global minimum.

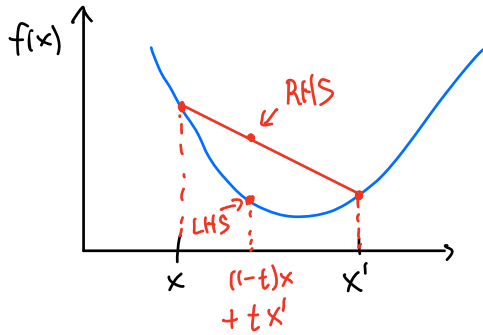
Def 1: $f(x)$ is convex $\Leftrightarrow f''(x) \geq 0$



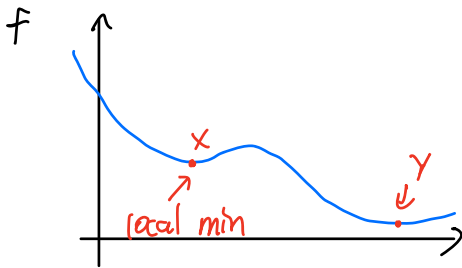
“Def 2”: convex function holds watter

“Def 3”: a function f is convex iff for all x, x' in domain of f
and $t \in [0, 1]$

$$f((1-t)x + tx') \leq (1-t)f(x) + tf(x')$$



If you draw a line connecting $(x, f(x))$ and $(x', f(x'))$, it must lie above the function itself.



Given that x is a local minimum of f , assume $\exists y, f(y) < f(x)$, show f cannot be convex. (by contradiction)

x is a local min of f iff $\exists \epsilon > 0$ s.t. $\forall z \in B_\epsilon(x), f(z) \geq f(x)$
 set of points where $\|x - z\| \leq \epsilon$

Choose some $t > 0$ s.t.

$$(1-t)x + tx' \in B_\epsilon(x)$$

(1) Because it's a local min, $f((1-t)x + tx') \geq f(x)$

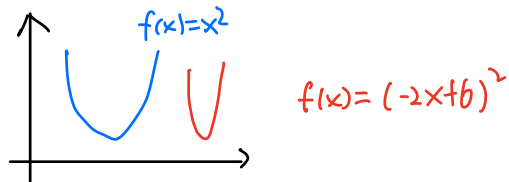
(2) Because of convexity,

$$\begin{aligned} f((1-t)x + tx') &\leq (1-t)f(x) + tf(x') \\ &< (1-t)f(x) + tf(x) \\ &= f(x) \end{aligned}$$

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

① If $f: \mathbb{R} \rightarrow \mathbb{R}$ and $f''(x) \geq 0$, then f is convex.

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



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

① $f(x) = x^2$ is convex (by ①)

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

③ $L(w)$ is convex (by ③)

Closed-form for Linear Regression ("Normal Equations")

$$\nabla_w L(w) = \frac{1}{n} \sum_{i=1}^n (w^T x^{(i)} - y^{(i)}) \cdot x^{(i)} = 0$$

$$\sum_{i=1}^n (w^T x^{(i)}) x^{(i)} = \sum_{i=1}^n x^{(i)} y^{(i)} \quad \leftarrow x^T y$$

$$X^T X w = \sum_{i=1}^n x^{(i)} y^{(i)} \quad \text{y = vector of } y^{(i)}_s$$

$$= \sum_{i=1}^n x^{(i)} (x^{(i)T} w)$$

$$= \left(\sum_{i=1}^n (x^{(i)} x^{(i)T}) \right) w \quad \begin{array}{l} \text{i-jth entry} \\ \sum_{k=1}^n x_i^{(k)} x_j^{(k)} \end{array}$$

$$\downarrow$$

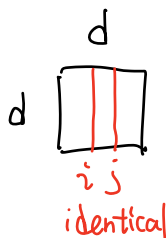
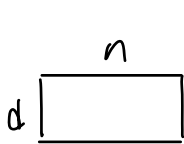
$$X^T X = \sum_{k=1}^n x_i^{(k)} x_j^{(k)}$$

$$X^T X w = X^T y$$

$$w = (X^T X)^{-1} X^T y \quad \boxed{\text{normal equations}}$$

Question: When would $X^T X$ is not invertible?

What if



suppose $x_i^{(k)} = x_j^{(k)} \quad \forall k$

$$\hat{y} = w_1 x_1 + \dots + \cancel{w_i x_i} + \dots + \cancel{w_j x_j} + \dots$$

$w_{i-100} \qquad w_{j+100}$

Result: answer is not unique anymore!

If i & j th column almost identical \Rightarrow leads to instability

In practice,

① Pseudoinverse A^+

– $A^+ = A^{-1}$ when A^{-1} exists

– for normal equations, A^+ always gives you an optimal solution

$$w = (X^T X)^+ X^T y$$

② Avoid highly-correlated features