# Linear Regression I  (8/29/2023)
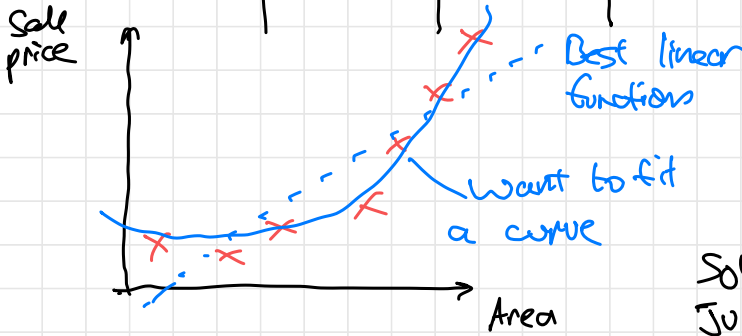
① How do we learn more complex functions?
② Why does gradient descent work for linear reg?
③ Why use squared error?

(y)

| Sale price | Area | #bed | house type | | Area$^2$ | Area$^3$ |
|---|---|---|---|---|---|---|
| 500k | 1200 | 2 | condo | | 1440000 | ~~ ... |

real number (Area) · integer (#bed) · categorical (house type) · Area$^x$



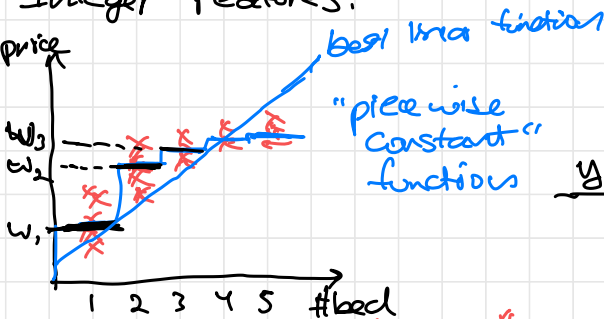Best linear function

Want to fit a curve

"predicted"

Want to learn:
$$\hat{y} = W_1 \cdot \text{Area}$$
$$+ W_2 \cdot \text{Area}^2$$
$$+ W_3 \cdot \text{Area}^3$$

Solution:
Just add more features!

Linear regression is linear in the input features
(that we can choose)

## Integer Features?



best linear function

"piecewise constant" functions

Add indicator features
some boolean function

| y | #bed = 1 | #bed = 2 | #bed 3 |
|---|---|---|---|
| 1 | 0 | 0 | |
| 0 | 1 | 0 | |
| 0 | 1 | 0 | |
| 0 | 0 | 1 | |

"indicator function"
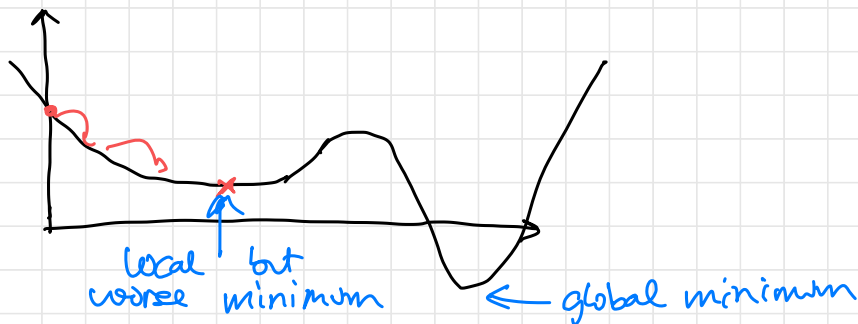$\mathbb{1}[\text{true}] = 1$, $\mathbb{1}[\text{false}] = 0$

Now: $\hat{y} = W_1 \cdot \mathbb{1}[\#bed = 1] + W_2 \cdot \mathbb{1}[\#bed = 2] + \cdots$

parameters

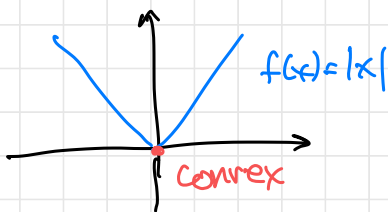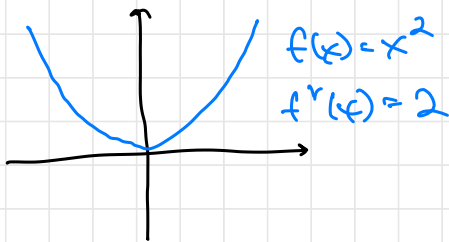"Feature Engineering": Process of choosing what features to use

Categorical features: add $\mathbb{1}[\text{is\_condo}]$, $\mathbb{1}[\text{is\_townhouse}]$

Convexity (why does gradient descent work?)



local but
worse minimum     ← global minimum

① Linear regression loss function $L(w)$ is convex
② For any convex function,
   all local minima are global minima

Def 1: $f(x)$ is convex $\iff$ $f''(x) \geq 0$ everywhere
                              assumes $f''$ exists everywhere

$f(x) = x^2$
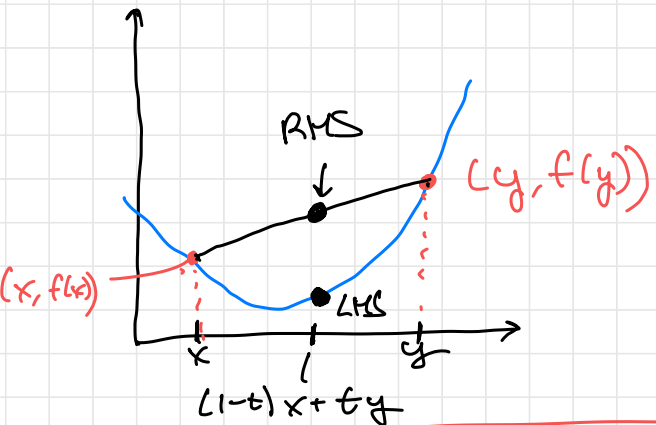$f''(x) = 2$

$f(x) = |x|$
convex

not convex

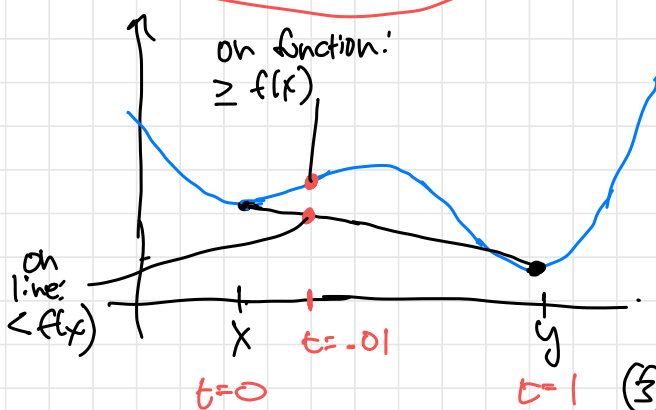Def 2 (informal): Convex function "holds water"

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

$$f\left((1-t)x + ty\right) \leq (1-t)f(x) + t\,f(y)$$

RHS

$(y, f(y))$

$(x, f(x))$

LHS

x

$(1-t)x + ty$

y

TL;DR:
If you draw line
connecting $(x, f(x))$
& $(y, f(y))$,

it must be above
the function

① All local minima of convex function are global minima



on function:
$\geq f(x)$

on line:
$< f(x)$

x

t=.01

y

t=0

t=1

① because x is local
min, there's some
small t
such that

$$f((1-t)x + ty) \geq f(x)$$

② Line from $(x, f(x))$
slopes downward

③ So $f((1-t)x + ty)$ must
be above the line between
x & y

④ Hence, f is not convex.

② linear regression is convex

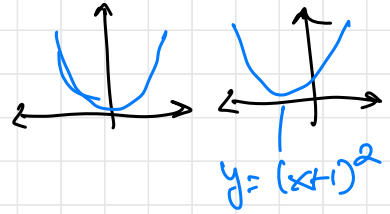→ $L(w) = \frac{1}{n} \sum_{i=1}^{n} \left( w^T x^{(i)} - y^{(i)} \right)^2$

Rules for convexity:
   ①. If $f: \mathbb{R} \to \mathbb{R}$ and $f''(x) \geq 0$ everywhere & exists everywhere

   then f is convex

② If $f$ is convex, then
$$g(x) = f(Ax + b) \underset{\text{is convex}}{\wedge} \text{ for any constants } A, b$$

③ If $f(x)$ and $g(x)$ are convex
$$f(x) + g(x) \text{ is convex}$$


$$y = (x+1)^2$$

④ If $f(x)$ is convex, and $c > 0$
$$c f(x) \text{ is convex}$$

$$L(w) = \frac{1}{n} \sum_{i=1}^{n} (w^T x^{(i)} - y^{(i)})^{\boxed{2}} \quad \leftarrow \text{why not 4?}$$
$$\text{absolute value?}$$

① $f(x) = x^2$ is convex by ①
② $(w^T x^{(i)} - y^{(i)})^2$ is convex ②
$$\underset{\text{parameter}}{\uparrow} \quad \underset{\text{constants}}{\uparrow}$$

③ $\frac{1}{n} \sum_{i=1}^{n} (w^T x^{(i)} - y^{(i)})^2$ is convex by ③ & ④

---

## Why square?. <u>Maximum Likelihood Estimation</u>

→ posit probabilistic process that generated data
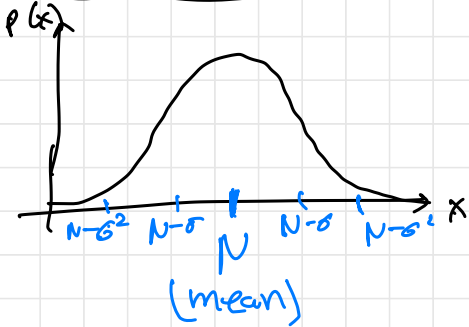→ choose parameters make observed data most likely?

E.J. Coin flips
observe $[H, T, H, H, H]$ ← observed data
unknown $p =$ prob. of heads ← parameter
Goal: Choose $p$ that makes data
most likely ← "learning"

Linear Regression: Assume $y^{(i)}$ drawn from Gaussian
w/ mean $w^T x^{(i)}$ & variance $\sigma^2$

determined by
"true" value of
parameter $w$

constant

independently

## Recall Gaussian



$p(x)$

$N-\sigma^2 \quad N-\sigma \quad N \quad N-\sigma \quad N-\sigma^2 \quad X$

$N$
(mean)

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

Likelihood of data (probability of data as a function of $w$)

$$\mathcal{L}(w) = \prod_{i=1}^{n} P(y^{(i)} \mid x^{(i)} ; w)$$

"parameterized by"

$$= \prod_{i=1}^{n} \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(y^{(i)} - w^T x^{(i)})^2}{2\sigma^2}\right)$$

Trick: take the log (monotonically increasing)

$$\log \mathcal{L}(w) = \sum_{i=1}^{n} \log\left(\frac{1}{\sigma \sqrt{2\pi}}\right) + \left(-\frac{(y^{(i)} - w^T x^{(i)})^2}{2\sigma^2}\right)$$

$$= \text{constant} + \frac{-1}{2\sigma^2} \sum_{i=1}^{n} (y^{(i)} - w^T x^{(i)})^2$$

Maximizing $\log \mathcal{L}(w)$ equivalent to
minimizing old $L(w)$