

ARGMIN AND MONOTONIC FUNCTIONS

① We've by now seen a few situations in which we want to compute $\operatorname{argmin}_x f(x)$ or $\operatorname{argmax}_x f(x)$:

e.g. $\operatorname{argmin}_\theta L(\theta)$ where L is a loss function

$\operatorname{argmax}_w P(w) \prod_{n=1}^N P(y^{(n)} | w, x^{(n)})$ for regression problems

In all of these situations, $f(x)$ has been nonnegative (i.e. $f: \mathbb{R} \rightarrow \mathbb{R}^+ \cup \{0\}$). Let's take a closer look at some strategies for computing $\operatorname{argmin}_x f(x)$ [or $\operatorname{argmax}_x f(x)$], given that f has a nonnegative range.

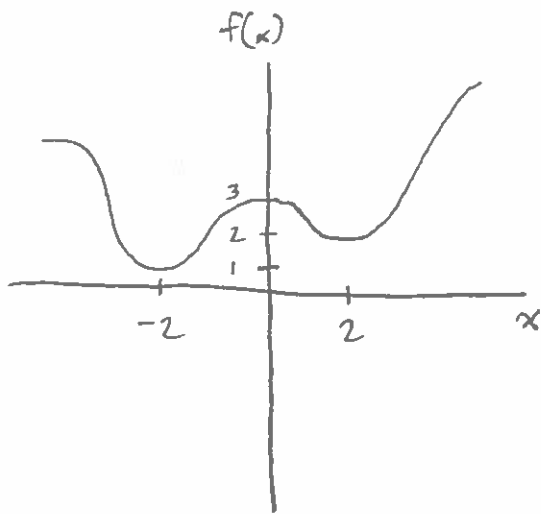
② One property we've already used (see REGRESSION PROBLEMS ①) is:

$$\operatorname{argmin}_x \frac{f(x)}{K} = \operatorname{argmin}_x f(x)$$

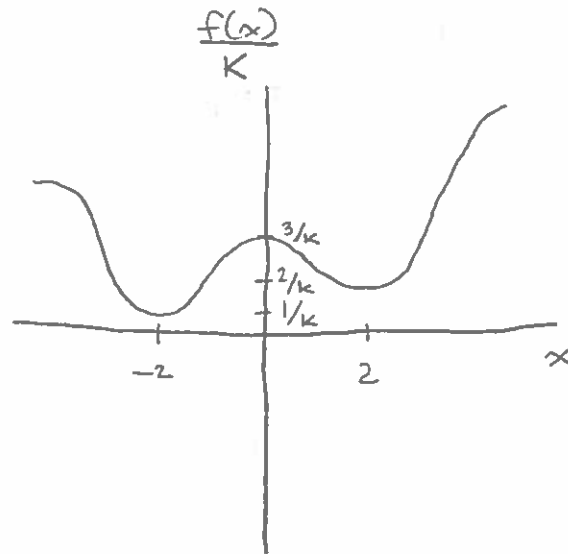
if K is a positive constant.

ARGMIN AND MONOTONIC FUNCTIONS

③ This is easy enough to justify with a picture



$$\operatorname{argmin}_x f(x) = -2$$

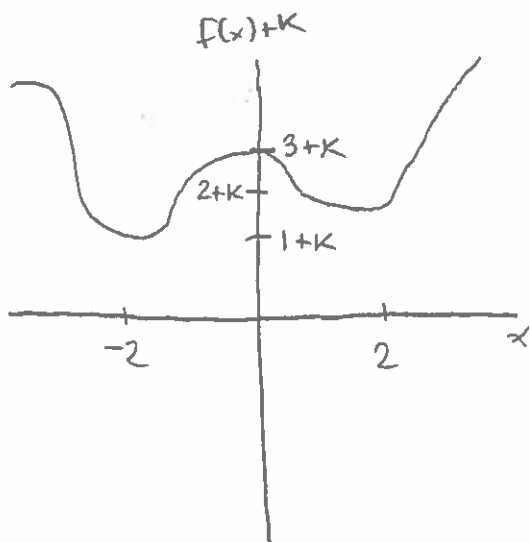


$$\operatorname{argmin}_x \frac{f(x)}{K} = -2$$

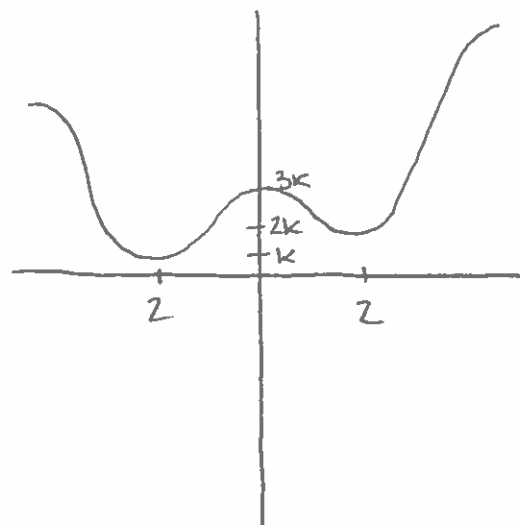
④ The same trick works in several other varieties:

$$\operatorname{argmin}_x f(x) + K = \operatorname{argmin}_x f(x) \quad \text{for any } K \in \mathbb{R}$$

$$\operatorname{argmin}_x K f(x) = \operatorname{argmin}_x f(x) \quad \text{for positive } K$$



$$\operatorname{argmin}_x f(x) + K = -2$$



$$\operatorname{argmin}_x K \cdot f(x) = -2$$

ARGMIN AND MONOTONIC FUNCTIONS

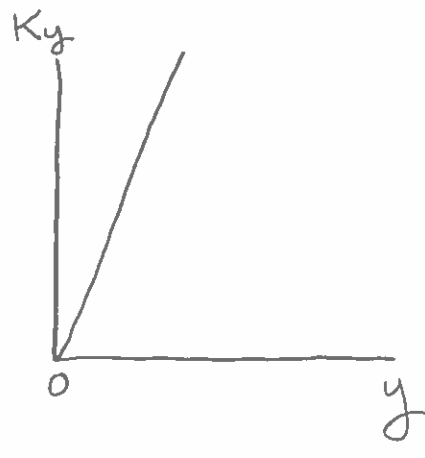
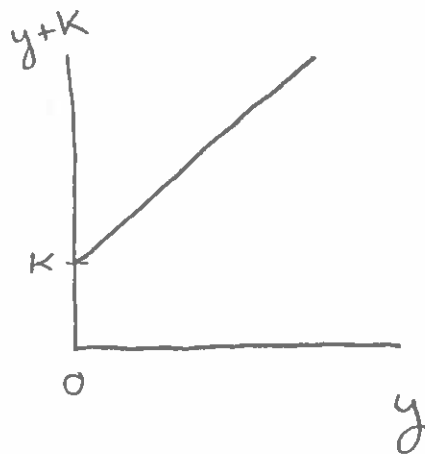
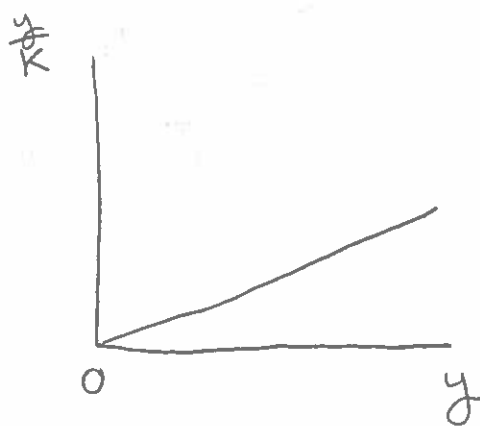
⑤ How can we generalize this trick? Well, let's look at these transformations as functions of $f(x)$, i.e.

$$g(y) = \frac{y}{K} \Rightarrow g(f(x)) = \frac{f(x)}{K}$$

$$g(y) \Rightarrow y + K \Rightarrow g(f(x)) = f(x) + K$$

$$g(y) \Rightarrow Ky \Rightarrow g(f(x)) = Kf(x)$$

We've assumed $f(x) \geq 0$ for all x , so let's look at $g(y)$ over the domain $[0, \infty)$:



⑥ All are monotonically increasing, i.e.

$$y_1 > y_2 \Leftrightarrow g(y_1) > g(y_2) \quad \text{for } y_1, y_2 \in \mathbb{R}^+ \cup \{0\}$$

This means:

$$f(x_1) > f(x_2) \Leftrightarrow g(f(x_1)) > g(f(x_2)) \quad \text{for } x_1, x_2 \in \mathbb{R}$$

which implies

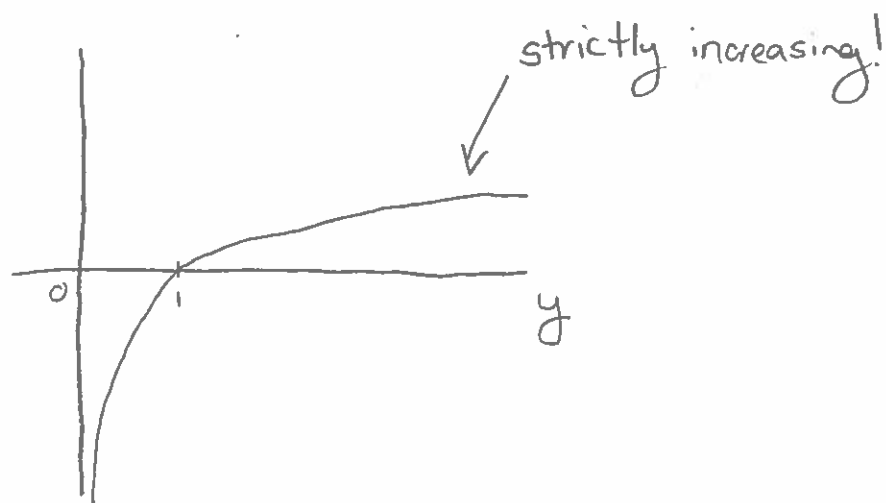
$$\operatorname{argmin}_x f(x) = \operatorname{argmin}_x g(f(x))$$

[Exercise: prove by contradiction]

ARGMIN AND MONOTONIC FUNCTIONS

- ⑦ The uncontested monarch of monotonic functions (over \mathbb{R}^+) is the logarithm, i.e.

$$g(y) = \log y$$



Because it is monotonically increasing, we know from ⑥ that:

$$\operatorname{argmin}_x f(x) = \operatorname{argmin}_x \log f(x)$$

- ⑧ This can be amazingly convenient when $f(x)$ is the product of complicated functions, e.g. suppose $f(x) = (2 + \sin x)(2 + \cos x)(x^2 + 1)$. First we observe that the range of $f(x)$ is positive. Thus:

$$\operatorname{argmin}_x f(x) = \operatorname{argmin}_x \log f(x)$$

$$= \operatorname{argmin}_x \log (2 + \sin x)(2 + \cos x)(x^2 + 1)$$

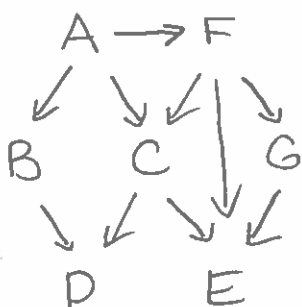
$$= \operatorname{argmin}_x \log (2 + \sin x) + \log (2 + \cos x) + \log (x^2 + 1)$$

ARGMIN AND MONOTONIC FUNCTIONS

- ⑨ The product has become a sum! Thus we can easily take the derivative:

$$\frac{d}{dx} \log f(x) = \frac{\cos x}{2 + \sin x} - \frac{\sin x}{2 + \cos x} + \frac{2x}{x^2 + 1}$$

- ⑩ This trick can be particularly useful when dealing with joint probability distributions. Consider a distribution that factors according to the following Bayesian network:



$$\text{i.e. } P(a, b, c, d, e, f, g) = P(a)P(b|a)P(c|a, f)P(d|b, c) \\ \cdot P(e|c, f, g)P(f|a)P(g|f)$$

Suppose we wanted to know the value of F that was most likely, given the other variables:

$$\hat{f} = \underset{f}{\operatorname{argmax}} P(F|a, b, c, d, e, g)$$

$$= \underset{f}{\operatorname{argmax}} \frac{P(a, b, c, d, e, f, g)}{\sum P(a, b, c, d, e, g)}$$

[total prob]

ARGMIN AND MONOTONIC FUNCTIONS

- ⑪ First thing to notice is that $P(a, b, c, d, e, g)$ is the same no matter what value we choose for f , so we can say:

$$\begin{aligned}\hat{f} &= \underset{f}{\operatorname{argmax}} \underbrace{P(a, b, c, d, e, f, g)}_K \\ &= \underset{f}{\operatorname{argmax}} P(a, b, c, d, e, f, g) \quad [\text{assuming } K > 0]\end{aligned}$$

- ⑫ We can then use the Bayesian network to simplify further:

$$\hat{f} = \underset{f}{\operatorname{argmax}} P(a)P(b|a)P(c|a, f)P(d|b, c)P(e|c, f, g)P(f|a)P(g|f)$$

There are two issues with this formula:

- products are tough to differentiate
- the overall quantity gets very small, very fast. Suppose each conditional probability term is .01 (i.e. a 1% probability). Then the overall expression evaluates to $(.01)^7 = .0000000000000001$. Pretty quickly this gets too small to represent in memory. (underflow).

ARGMIN AND MONOTONIC FUNCTIONS

⑬ But since $\log y$ is monotonically increasing over \mathbb{R}^+ :

$$\begin{aligned}\hat{f} &= \operatorname{argmax}_f \log P(a)P(b|a)P(c|a,f)P(d|b,c)P(e|c,f,g)P(f|a)P(g|f) \\ &= \operatorname{argmax}_f \log P(a) + \log P(b|a) + \log P(c|a,f) + \log P(d|b,c) \\ &\quad + \log P(e|c,f,g) + \log P(f|a) + \log P(g|f)\end{aligned}$$

Now, not only do we have an easy sum to differentiate, the quantity is also easy to represent in memory:

$$\begin{aligned}7 \cdot \log(.01) &= 7 \cdot (-2) \\ &= -14\end{aligned}$$