

People telling me AI is going
to destroy the world

My neural network



Organizational

- Homework: **You nearly did it!**
 - Homework 7 to be released next Wednesday (you will have until Dec 6th)
- No more extensions

Midterm Wednesday 11/29 4.30 to 5.45

- Instructions were released
- Proposals due Nov 8th
- Get to work!
 - Check boxes for grading first
 - Then, try to have a little fun =)



Interpretability, Generalization and Domain Gaps

Bag of Features - BagNet

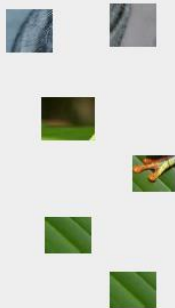


Bag-of-Features?

Step 1: Build codebook



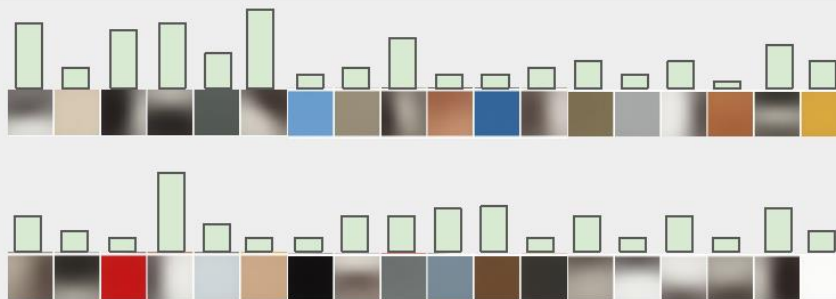
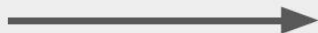
Extract random
patches



Cluster patches to
form “codebook”
of “visual words”



Step 2: Encode images

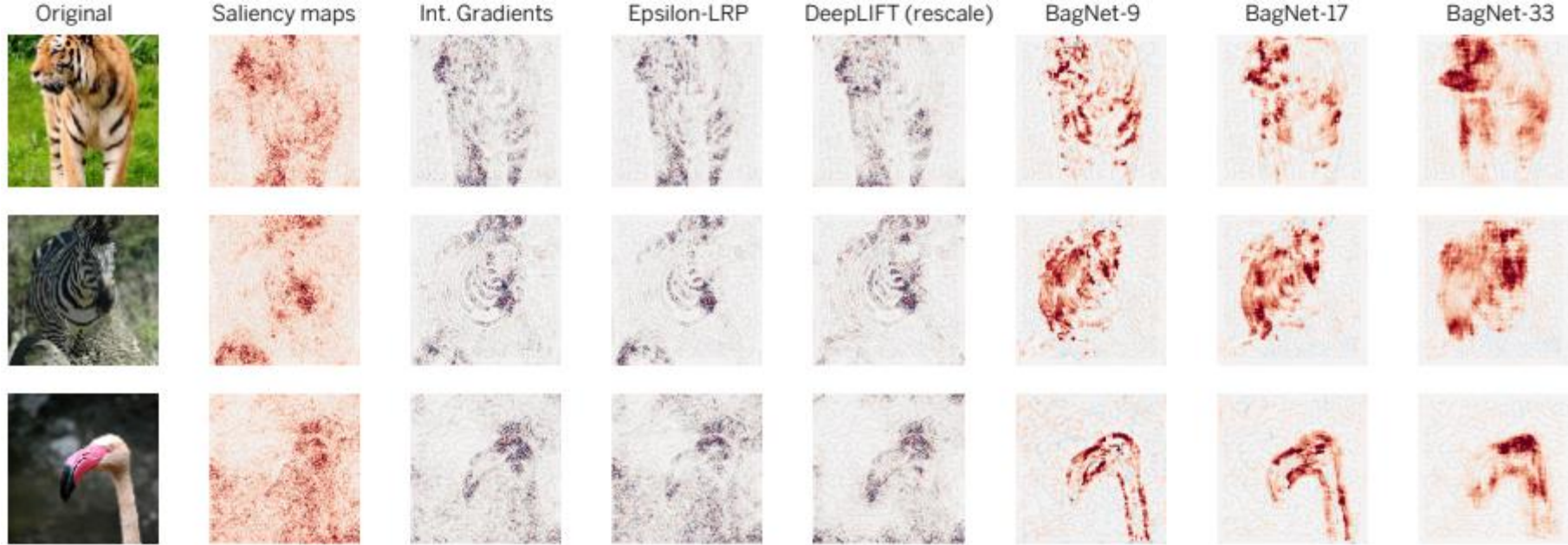


Bag-of-Features?

Deep Features

- Input $Q \times Q$ px image patch
- BagNet generates 2048-dim feature vector
- Linear classifier on this 2048-dim vector
- After linear classification: #classes heatmaps of “class evidence”
- Average class evidence and pass into softmax → Probability

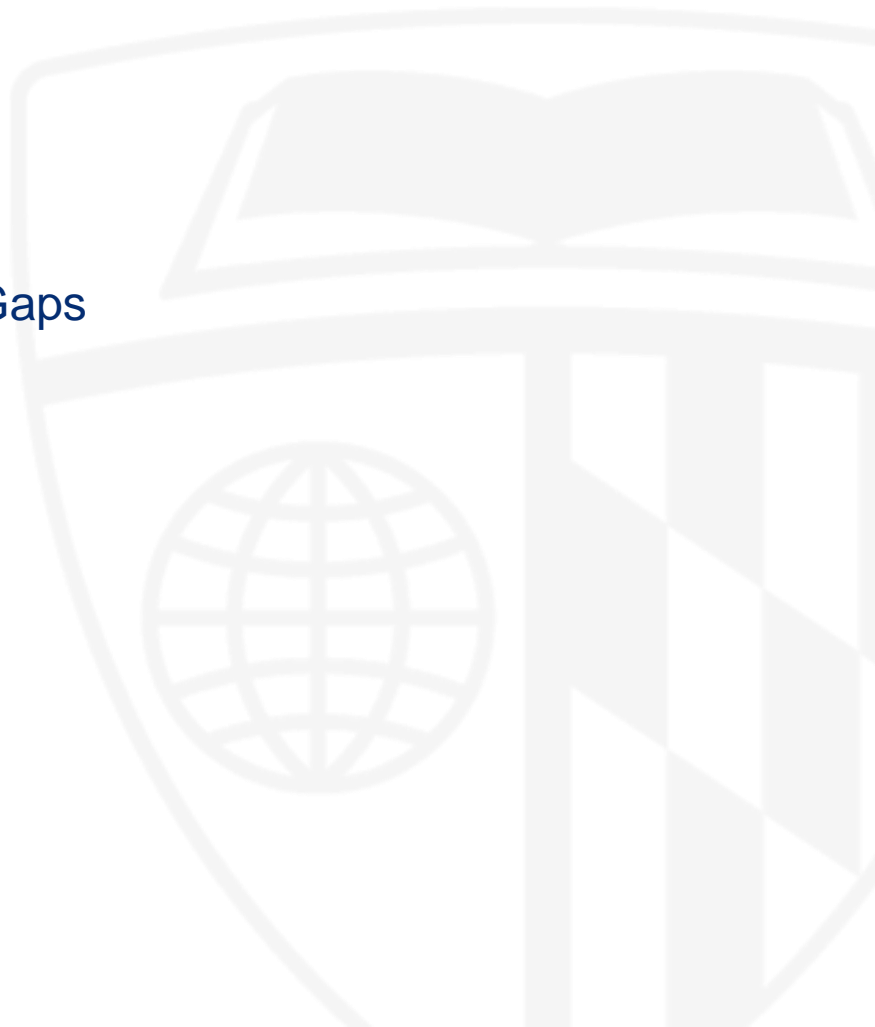
[Brendel, W., & Bethge, M. \(2018\). Approximating CNNs with Bag-of-local-Features models works surprisingly well on ImageNet. ICLR.](#)



- Similarities in “decision making” (saliency) to SOTA models suggest that current network architectures base their decisions on a large number of weak and local statistical regularities
- One way forward is to define novel tasks that cannot be solved using local statistical regularities

Interpretability, Generalization and Domain Gaps

Deep Image Prior



Deep Image Prior

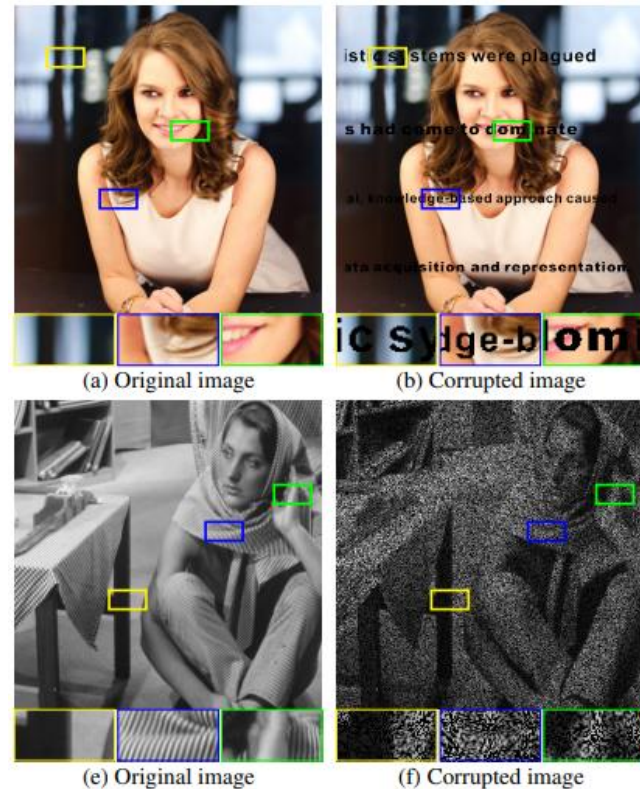
- Image restoration $x^* = \min_x E(x, x_0) + R(x)$, where

E is a task-dependent data term (similarity),
 x_0 is a corrupted observation (noise, occlusion, etc), and
 R is a regularizer (e.g. total variation, L2, wavelets,...)

- Regularizer is replaced by a CNN

$$\theta^* = \min_{\theta} E(f_{\theta}(z), x_0), \quad x^* = f_{\theta^*}(z)$$

Minimizer is optimized using a random z , starting from random initialization of θ



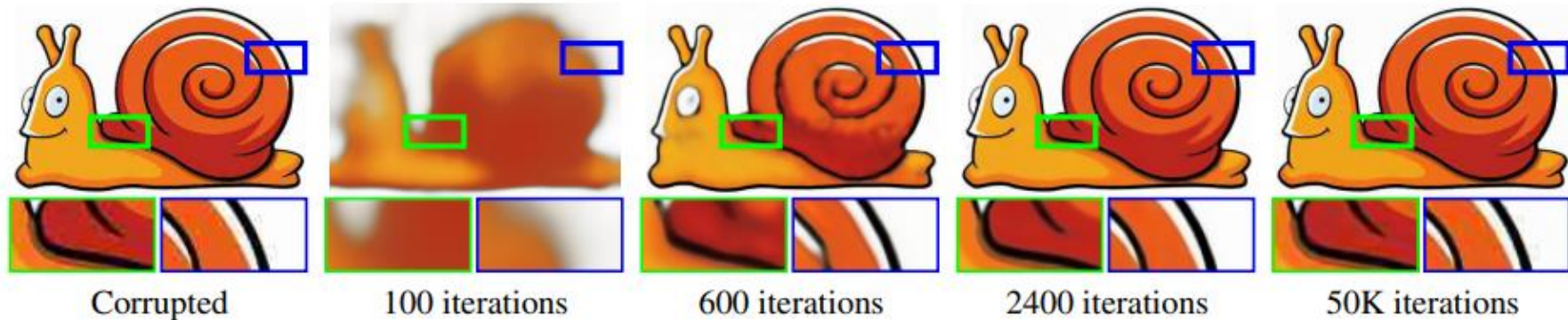


Figure 3: **Blind restoration of a JPEG-compressed image.** (*electronic zoom-in recommended*) Our approach can restore an image with a complex degradation (JPEG compression in this case). As the optimization process progresses, the deep image prior allows to recover most of the signal while getting rid of halos and blockiness (after 2400 iterations) before eventually overfitting to the input (at 50K iterations).

- Optimization is fast for “natural images” without corruption
- CNN can fit corruption, but very reluctantly
- Parametrization (CNN architecture) has high impedance to noise and low impedance to signal!

[Ulyanov, D., Vedaldi, A., & Lempitsky, V. \(2018\). Deep image prior. CVPR \(pp. 9446-9454\).](#)

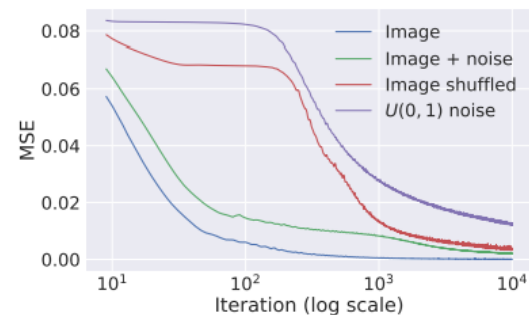
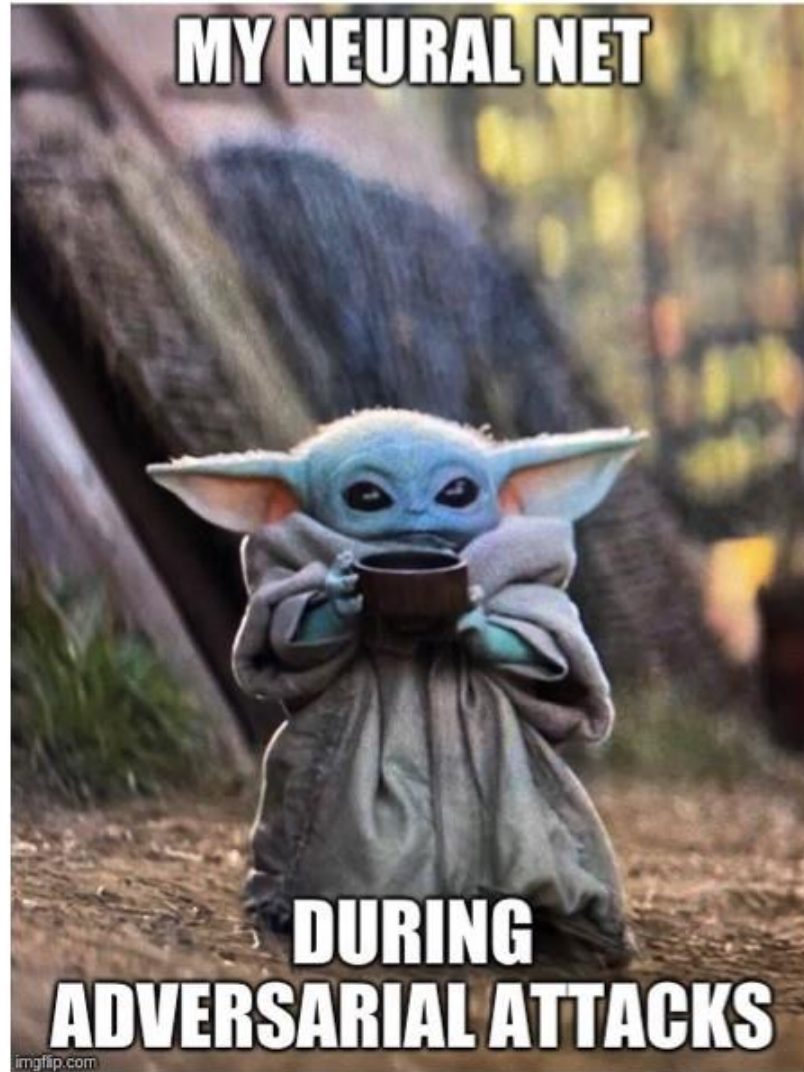


Figure 2: Learning curves for the reconstruction task using: a natural image, the same plus i.i.d. noise, the same randomly scrambled, and white noise. Naturally-looking images result in much faster convergence, whereas noise is rejected.



Adversarial Examples: What are they?



x

“panda”

57.7% confidence

$+ .007 \times$



$\text{sign}(\nabla_x J(\theta, x, y))$

“nematode”

8.2% confidence

$=$



$x +$

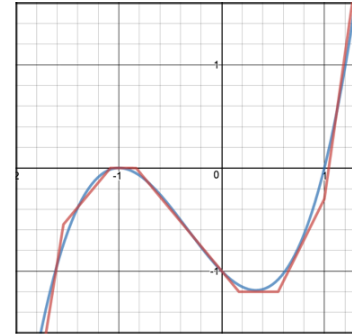
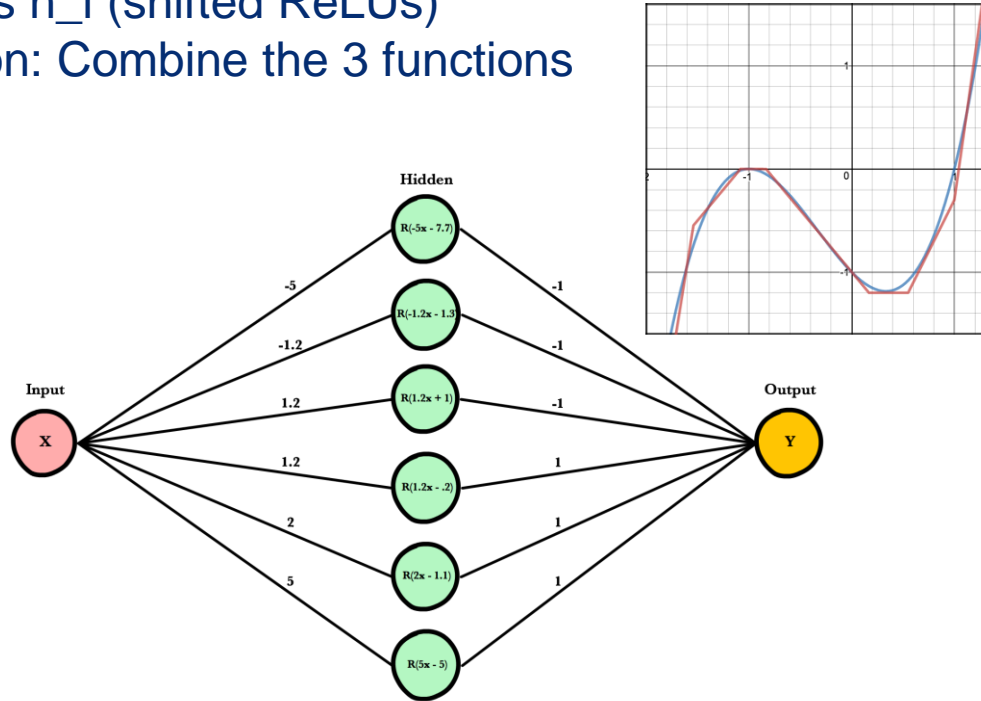
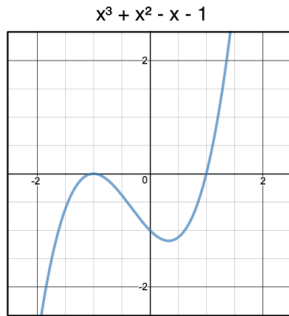
$\epsilon \text{sign}(\nabla_x J(\theta, x, y))$

“gibbon”

99.3 % confidence

Understanding Universal Approximation

- $f(W,x)$: ReLU activation; W in $\mathbb{R}(3 \times 1)$, b in $\mathbb{R}(3 \times 1)$, and x in $\mathbb{R}(1 \times 1)$
 \rightarrow 3 functions n_i (shifted ReLUs)
- Output neuron: Combine the 3 functions



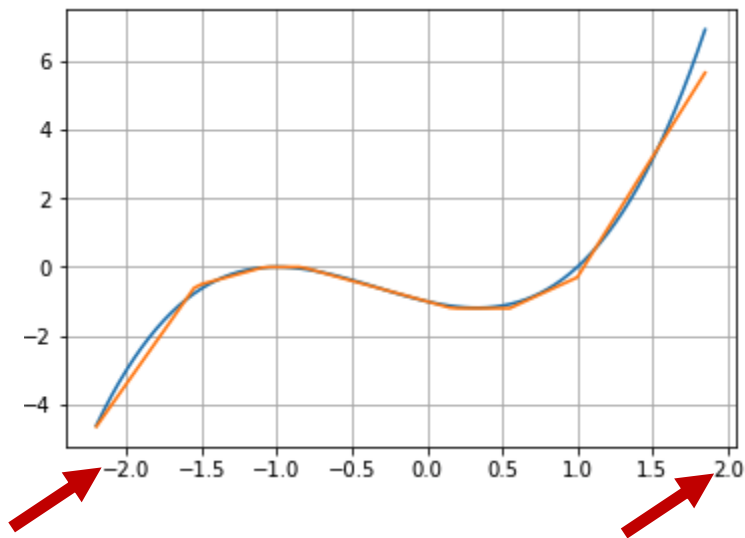
$$\begin{aligned} n_1(x) &= \text{Relu}(-5x - 7.7) \\ n_2(x) &= \text{Relu}(-1.2x - 1.3) \\ n_3(x) &= \text{Relu}(1.2x + 1) \\ n_4(x) &= \text{Relu}(1.2x - .2) \\ n_5(x) &= \text{Relu}(2x - 1.1) \\ n_6(x) &= \text{Relu}(5x - 5) \end{aligned}$$

$$\begin{aligned} Z(x) &= -n_1(x) - n_2(x) - n_3(x) \\ &\quad + n_4(x) + n_5(x) + n_6(x) \end{aligned}$$

Example taken from a blog post: [link](#).

Compact Domains

Another example



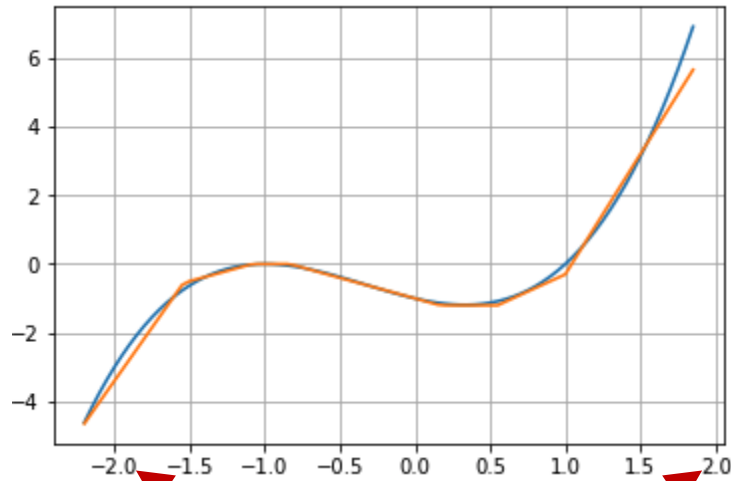
Note the range of x !

```
def feed_forward(input):  
    Zh = np.dot(input, Wh)  
    Ah = relu(Zh + Bh)  
    Zo = np.dot(Ah, Wo)  
    Ao = Zo + Bo  
    return Ao
```

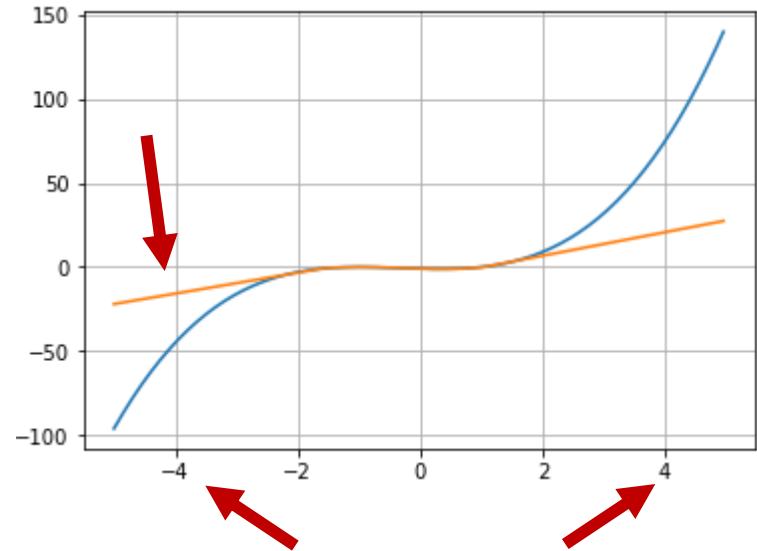
Example taken from a blog post: [link](#).

Compact Domains

Another example



Note the range of x !



What happened?

Example taken from a blog post: [link](#).