# CSCE 5218 & 4930 Deep Learning

## Neural Network Training
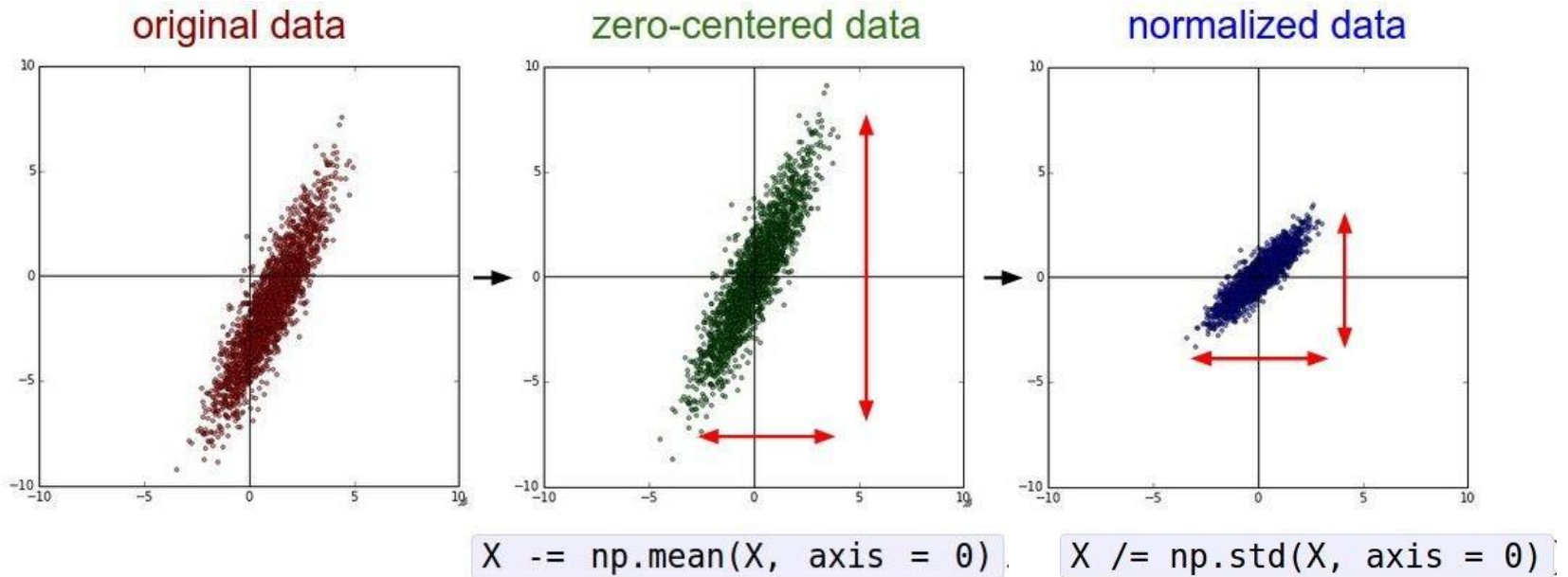
# Plan for this lecture

- Tricks of the trade
  - Preprocessing, initialization, normalization
  - Dealing with limited data
- Convergence of gradient descent
  - How long will it take?
  - Will it work at all?
- Different optimization strategies
  - Alternatives to SGD
  - Learning rates
  - Choosing hyperparameters
- How to do the computation
  - Computation graphs
  - Vector notation (Jacobians)

# Tricks of the trade

# Practical matters

- Getting started: Preprocessing, initialization, normalization, choosing activation functions
- Improving performance and dealing with sparse data: regularization, augmentation, transfer learning
- Hardware and software
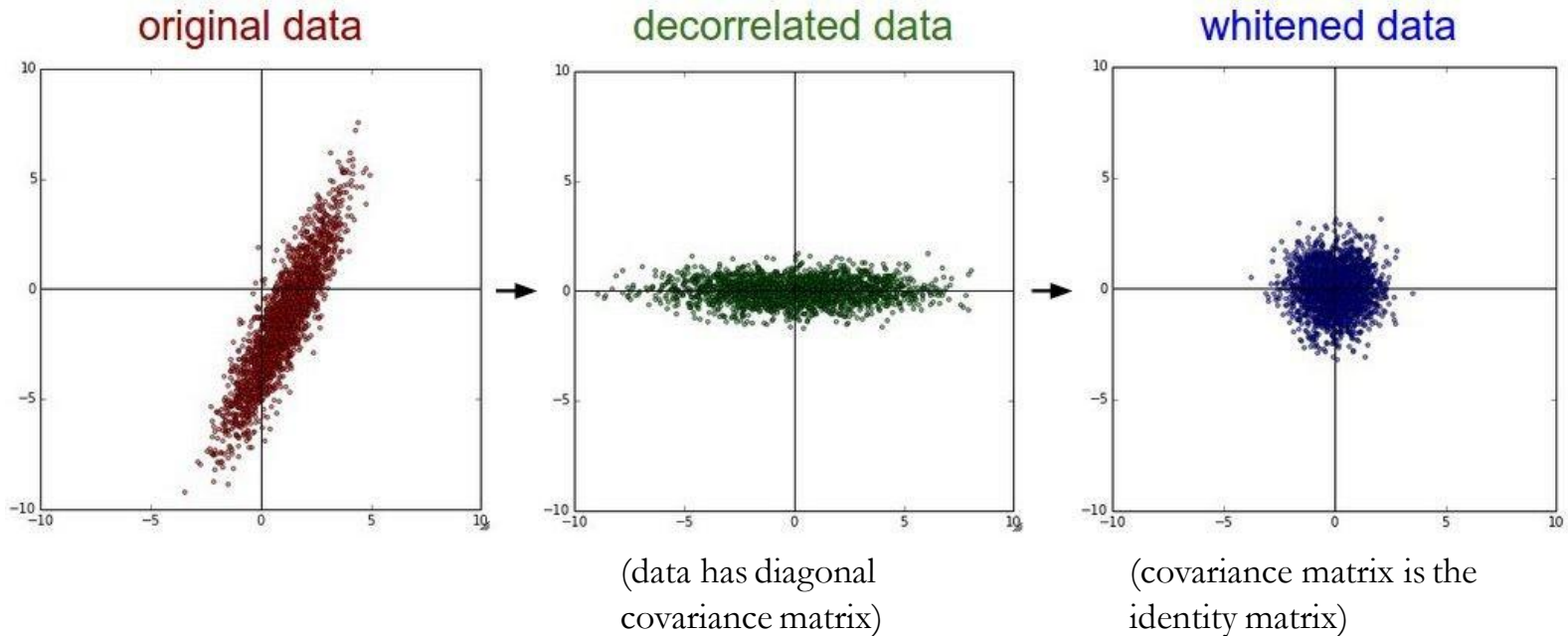- Extra reading/visualization resources
  - https://www.deeplearning.ai/ai-notes/initialization/
  - https://www.deeplearning.ai/ai-notes/optimization/

# Preprocessing the Data



original data      zero-centered data      normalized data

```
X -= np.mean(X, axis = 0)
```
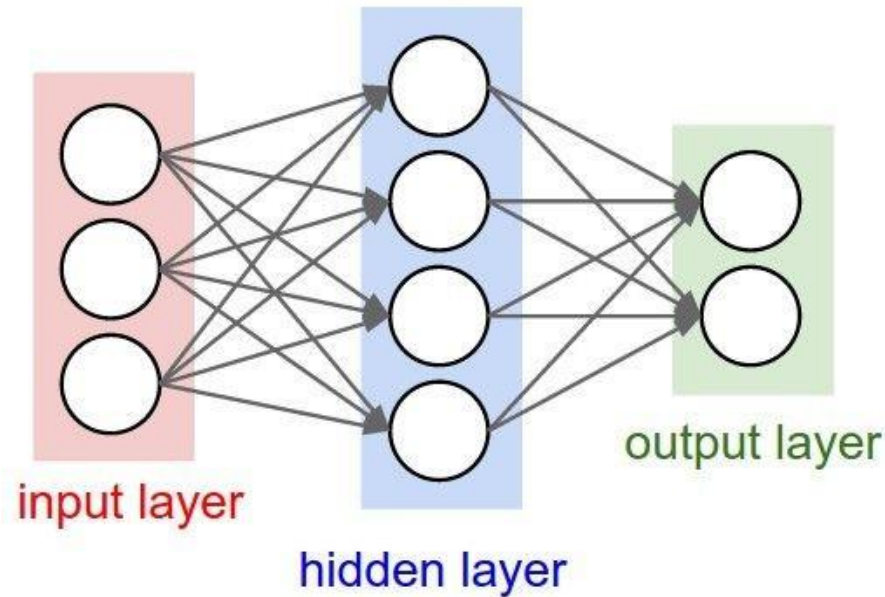
```
X /= np.std(X, axis = 0)
```

(Assume X [NxD] is data matrix,
each example in a row)

# Preprocessing the Data

In practice, you may also see **PCA** and **Whitening** of the data



| original data | decorrelated data | whitened data |
| --- | --- | --- |
| | (data has diagonal covariance matrix) | (covariance matrix is the identity matrix) |

# Weight Initialization



input layer

hidden layer

output layer

- Q: what happens when W=constant init is used?

# Weight Initialization

- Another idea: **Small random numbers**

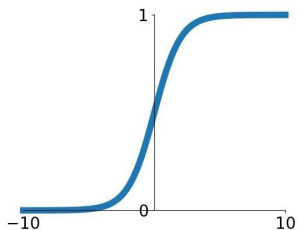(gaussian with zero mean and 1e-2 standard deviation)

```
W = 0.01* np.random.randn(D,H)
```

Works ~okay for small networks, but problems with deeper networks.
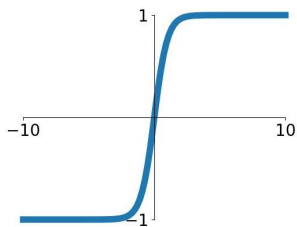
# Activation Functions
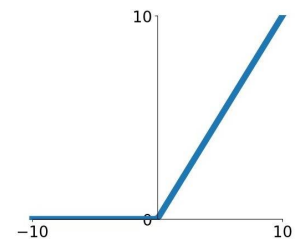
**Sigmoid**

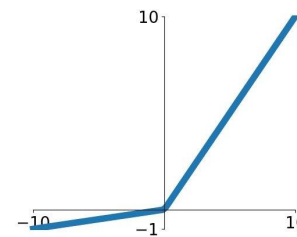$\sigma(x) = \frac{1}{1+e^{-x}}$

**tanh**

$\tanh(x)$

**ReLU**

$\max(0, x)$

**Leaky ReLU**

$\max(0.1x, x)$

**Maxout**

$\max(w_1^T x + b_1, w_2^T x + b_2)$

**EL**

$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$

Fei-Fei Li, Andrej Karpathy, Justin Johnson, Serena Yeung

# Activation Functions

$$\sigma(x) = 1/(1 + e^{-x})$$



**Sigmoid**

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating "firing rate" of a neuron
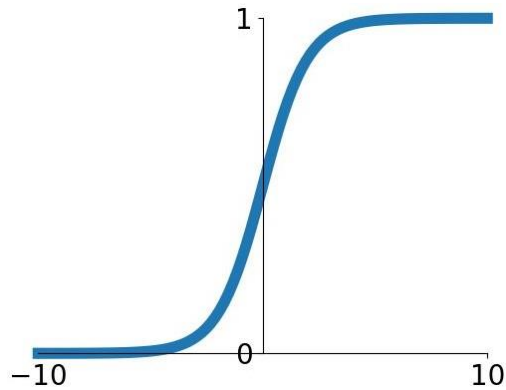
# Activation Functions

$$\sigma(x) = 1/(1 + e^{-x})$$



**Sigmoid**

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating "firing rate" of a neuron

• 3 problems:

1. Saturated neurons "kill" the gradients

x

$$\sigma(x) = 1/(1 + e^{-x})$$

$$\frac{\partial \sigma}{\partial x}$$ sigmoid gate

$$\frac{\partial L}{\partial x} = \frac{\partial \sigma}{\partial x} \frac{\partial L}{\partial \sigma}$$

$$\frac{\partial L}{\partial \sigma}$$

What happens when x = -10?
What happens when x = 0?
What happens when x = 10?

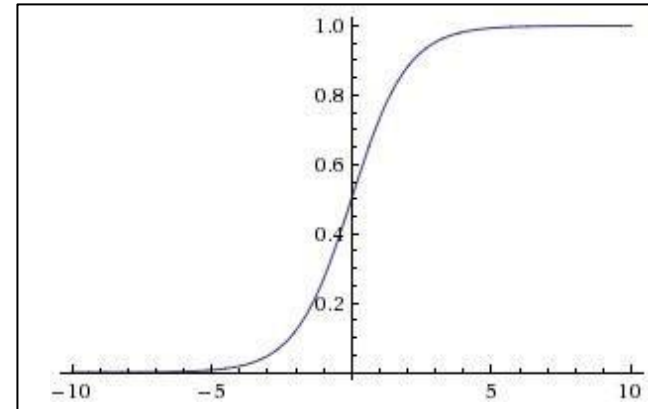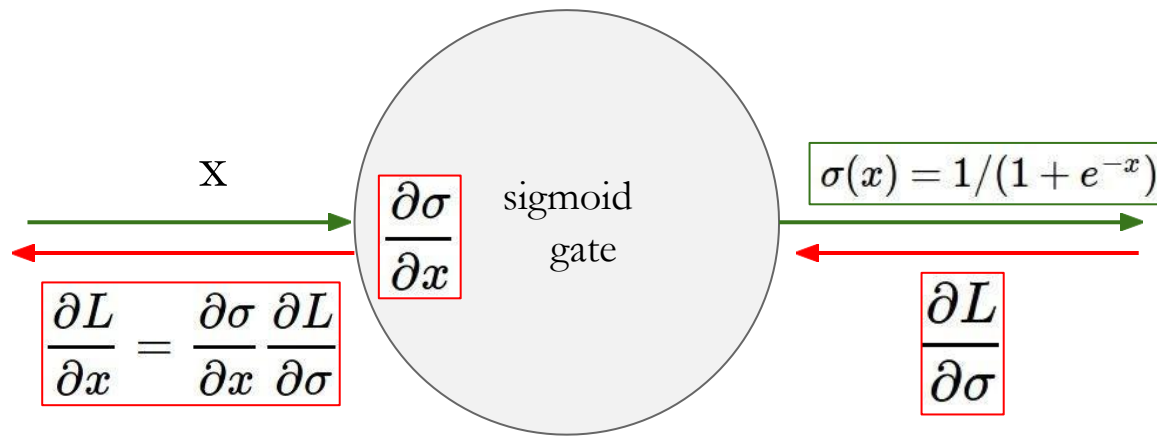# Activation Functions

$$\sigma(x) = 1/(1 + e^{-x})$$

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating "firing rate" of a neuron

- 3 problems:

  1. Saturated neurons "kill" the gradients
  2. Sigmoid outputs are not zero-centered

**Sigmoid**

# Activation Functions

$$\sigma(x) = 1/(1 + e^{-x})$$

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating "firing rate" of a neuron

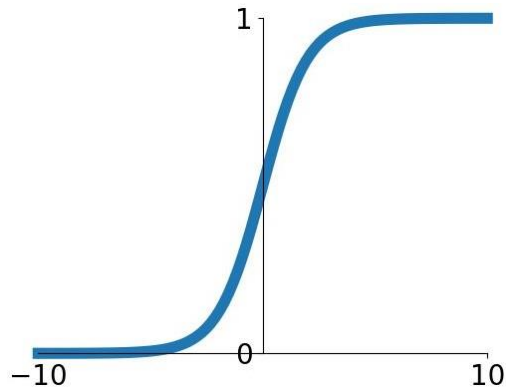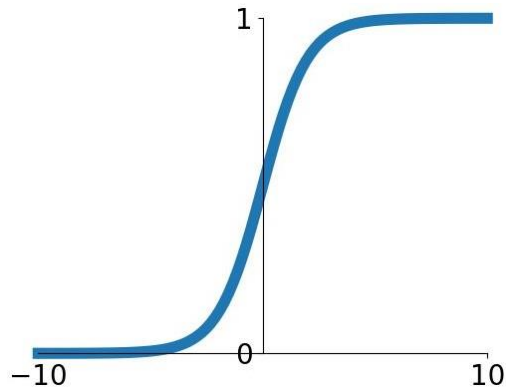- 3 problems:

  1. Saturated neurons "kill" the gradients
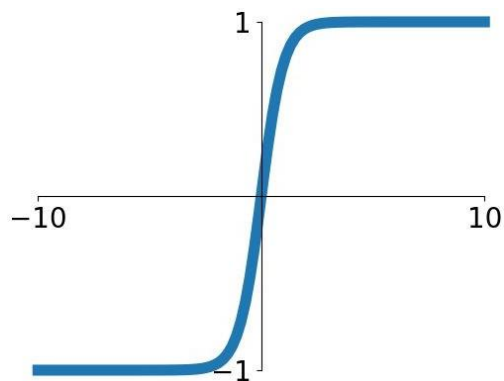  2. Sigmoid outputs are not zero-centered
  3. exp() is a bit compute expensive



**Sigmoid**

# Activation Functions



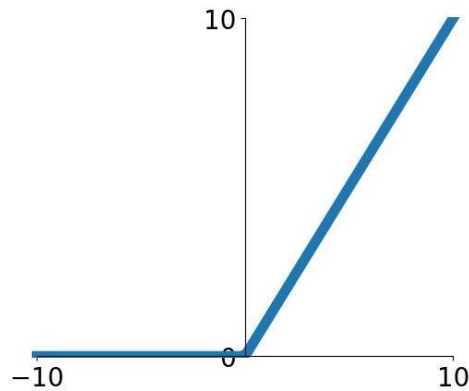**tanh(x)**

- Squashes numbers to range [-1,1]
- zero centered (nice)
- still kills gradients when saturated :(

[LeCun et al., 1991]

# Activation Functions



**ReLU**
(Rectified Linear Unit)

- Computes f(x) = max(0,x)

- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)

[Krizhevsky et al., 2012]

# Activation Functions

- Computes $f(x) = \max(0, x)$



**ReLU**
(Rectified Linear Unit)

- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)

- Not zero-centered output

# Activation Functions

- Computes $f(x) = \max(0,x)$



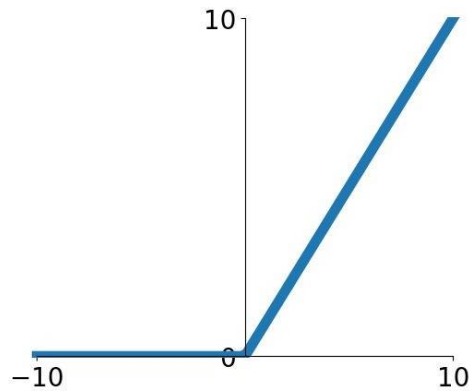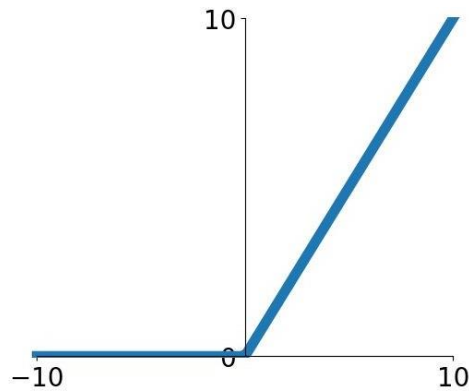**ReLU**
(Rectified Linear Unit)

- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)

- Not zero-centered output
- An annoyance:

hint: what is the gradient when x < 0?

x

$$\frac{\partial \sigma}{\partial x}$$

ReLU gate

$$\sigma(x) = \max(0, x)$$

$$\frac{\partial L}{\partial x} = \frac{\partial \sigma}{\partial x} \frac{\partial L}{\partial \sigma}$$

$$\frac{\partial L}{\partial \sigma}$$

What happens when x = -10?
What happens when x = 0?
What happens when x = 10?

# Activation Functions

- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not "die".**

**Leaky ReLU**

$$f(x) = \max(0.01x, x)$$

# Activation Functions

[Mass et al., 2013]
[He et al., 2015]



- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not "die".**

## Leaky ReLU

$$f(x) = \max(0.01x, x)$$

## Parametric Rectifier (PReLU)

$$f(x) = \max(\alpha x, x)$$

backprop into alpha (parameter)

# Activation Functions

## **Exponential Linear Units (ELU)**



- All benefits of ReLU
- Closer to zero mean outputs
- Negative saturation regime compared with Leaky ReLU  adds some robustness to noise

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha \left( \exp(x) - 1 \right) & \text{if } x \le 0 \end{cases}$$

- Computation requires exp()

Fei-Fei Li, Andrej Karpathy, Justin Johnson, Serena Yeung

Maxout "Neuron"  [Goodfellow et al., 2013]

- Does not have the basic form of dot product -> nonlinearity
- Generalizes ReLU and Leaky ReLU
- Linear Regime! Does not saturate! Does not die!

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

Problem: doubles the number of parameters/neuron :(

TLDR: In practice:

- Use ReLU. Be careful with your learning rates
- Try out Leaky ReLU / Maxout / ELU / PReLU
- Try out tanh but don't expect much
- Don't use sigmoid

# Batch Normalization

"you want zero-mean unit-variance activations? just make them so."

consider a batch of activations at some layer. To make each dimension zero-mean unit-variance, apply:

$$\widehat{x}^{(k)} = \frac{x^{(k)} - \mathrm{E}[x^{(k)}]}{\sqrt{\mathrm{Var}[x^{(k)}]}}$$

# Batch Normalization

"you want zero-mean unit-variance activations? just make them so."

N

D

1. compute the empirical mean and variance independently for each dimension.

2. Normalize

$$\widehat{x}^{(k)} = \frac{x^{(k)} - \mathrm{E}[x^{(k)}]}{\sqrt{\mathrm{Var}[x^{(k)}]}}$$

Fei-Fei Li, Andrej Karpathy, Justin Johnson, Serena Yeung

# Batch Normalization

Normalize:

$$\widehat{x}^{(k)} = \frac{x^{(k)} - \mathrm{E}[x^{(k)}]}{\sqrt{\mathrm{Var}[x^{(k)}]}}$$

And then allow the network to squash the range if it wants to:

$$y^{(k)} = \gamma^{(k)}\widehat{x}^{(k)} + \beta^{(k)}$$

Note, the network can learn:

$$\gamma^{(k)} = \sqrt{\mathrm{Var}[x^{(k)}]}$$

$$\beta^{(k)} = \mathrm{E}[x^{(k)}]$$

to recover the identity mapping.

# Batch Normalization

**Input:** Values of $x$ over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$;
Parameters to be learned: $\gamma, \beta$
**Output:** $\{y_i = BN_{\gamma,\beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^{m} x_i \qquad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^{m} (x_i - \mu_{\mathcal{B}})^2 \qquad // \text{ mini-batch variance}$$

$$\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \qquad // \text{ normalize}$$

$$y_i \leftarrow \gamma \widehat{x}_i + \beta \equiv BN_{\gamma,\beta}(x_i) \qquad // \text{ scale and shift}$$

- Improves gradient flow through the network
- Allows higher learning rates
- Reduces the strong dependence on initialization
- Acts as a form of regularization

# Batch Normalization

**Input:** Values of $x$ over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$;
Parameters to be learned: $\gamma, \beta$

**Output:** $\{y_i = \text{BN}_{\gamma,\beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^{m} x_i \qquad \text{// mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^{m} (x_i - \mu_{\mathcal{B}})^2 \qquad \text{// mini-batch variance}$$

$$\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \qquad \text{// normalize}$$

$$y_i \leftarrow \gamma \widehat{x}_i + \beta \equiv \text{BN}_{\gamma,\beta}(x_i) \qquad \text{// scale and shift}$$

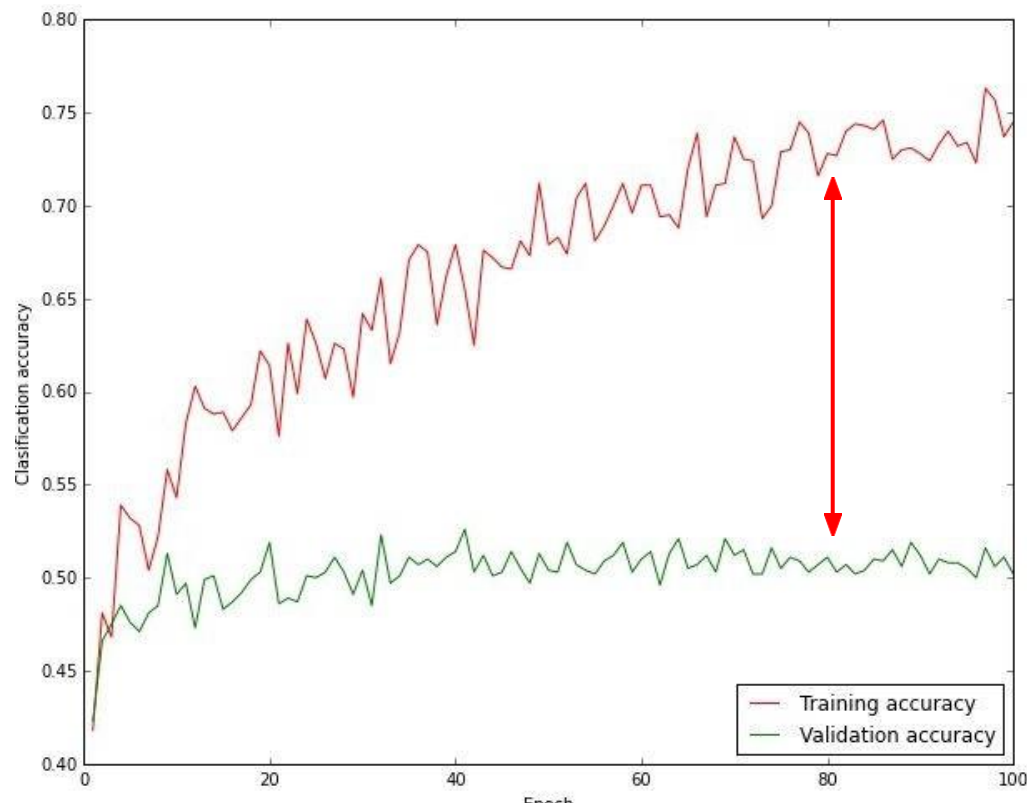**Note: at test time BatchNorm layer functions differently:**

The mean/std are not computed based on the batch. Instead, a single fixed empirical mean of activations during training is used.

(e.g. can be estimated during training with running averages)

# Babysitting the Learning Process

- Preprocess data

- Choose architecture

- Initialize and check initial loss with no regularization

- Increase regularization, loss should increase

- Then train – try small portion of data, check you can overfit

- Add regularization, and find learning rate that can make the loss go down

- Check learning rates in range [1e-3 … 1e-5]

- Coarse-to-fine search for hyperparameters (e.g. learning rate, regularization)

# Monitor and Visualize Accuracy



big gap = overfitting
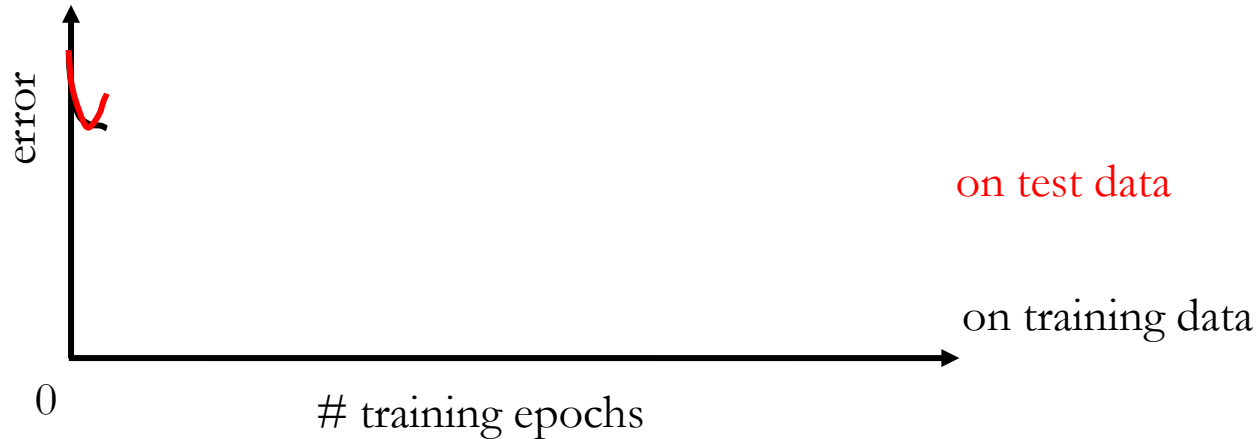=> increase regularization strength?

no gap
=> increase model capacity?

# Dealing with sparse data

- Deep neural networks require lots of data, and can overfit easily

- The more weights you need to learn, the more data you need

- That's why with a deeper network, you need more data for training than for a shallower network

- Ways to prevent overfitting include:
  - Using a validation set to stop training or pick parameters
  - Regularization
  - Data augmentation
  - Transfer learning

# Over-training prevention

- Running too many epochs can result in over-fitting.

error

0

\# training epochs

on test data

on training data

- Keep a hold-out validation set and test accuracy on it after every epoch. Stop training when additional epochs actually increase validation error.
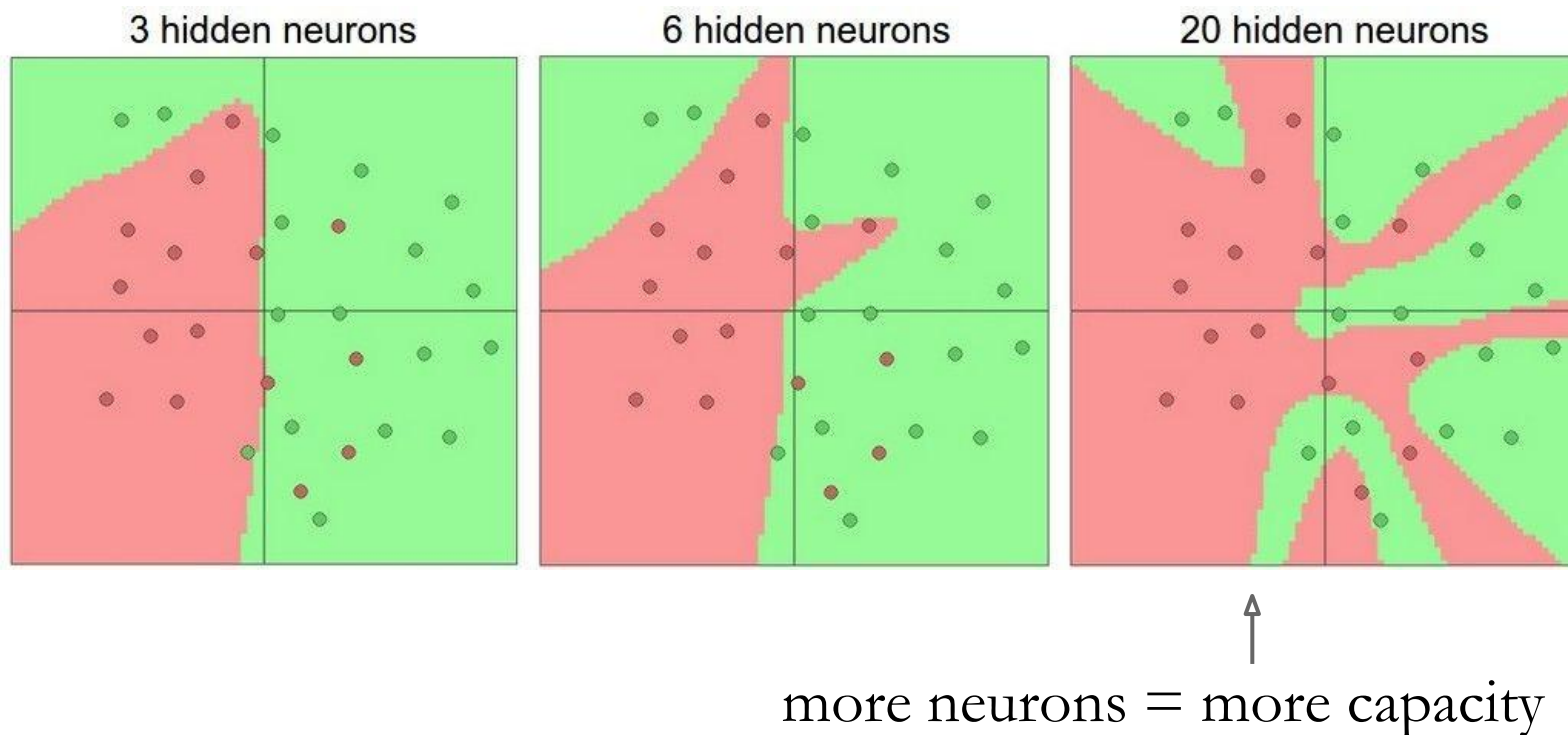
# Determining best number of hidden units

- Too few hidden units prevent the network from adequately fitting the data.

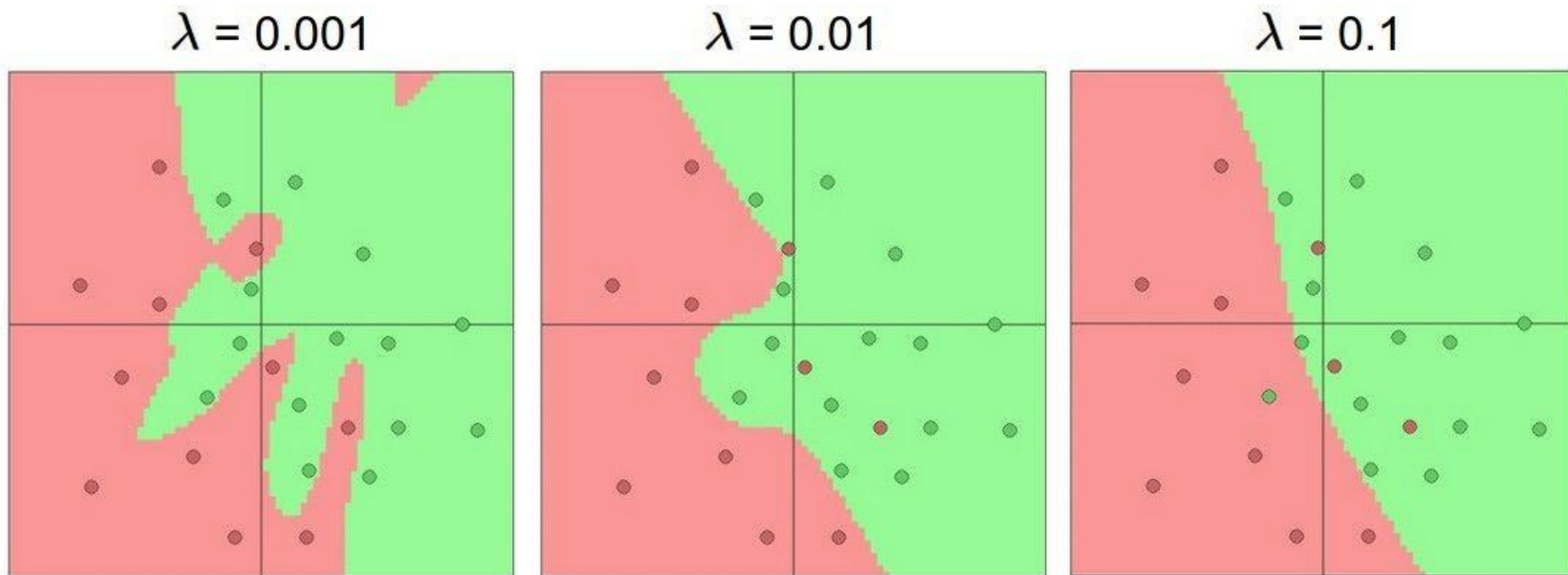- Too many hidden units can result in over-fitting.

error | # hidden units

on test data

on training data

0

- Use internal cross-validation to empirically determine an optimal number of hidden units.

# Effect of number of neurons



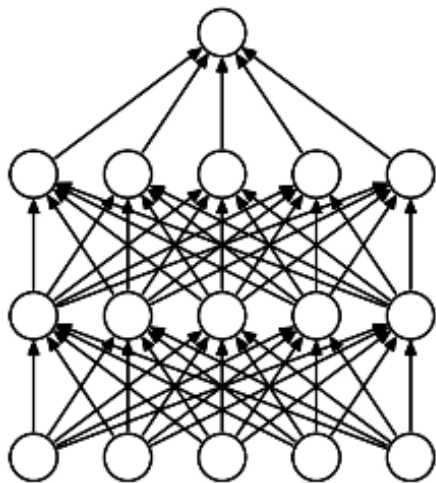3 hidden neurons      6 hidden neurons      20 hidden neurons

more neurons = more capacity

# Effect of regularization

Do not use size of neural network as a regularizer. Use stronger regularization instead:



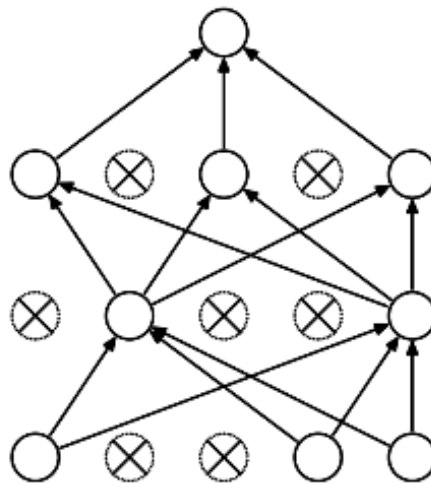$\lambda = 0.001$      $\lambda = 0.01$      $\lambda = 0.1$

(you can play with this demo over at ConvNetJS: http://cs.stanford.
edu/people/karpathy/convnetjs/demo/classify2d.html)
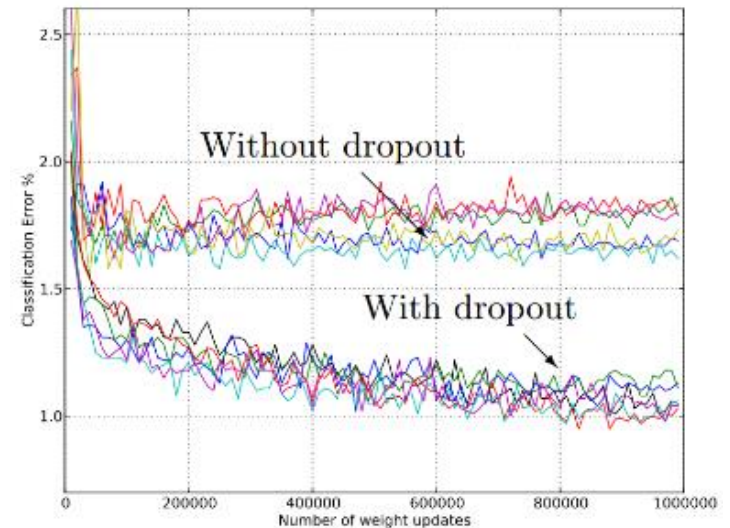
# Regularization

- L1, L2 regularization (*weight decay*)

- Dropout
  - Randomly turn off some neurons
  - Allows individual neurons to independently be responsible for performance
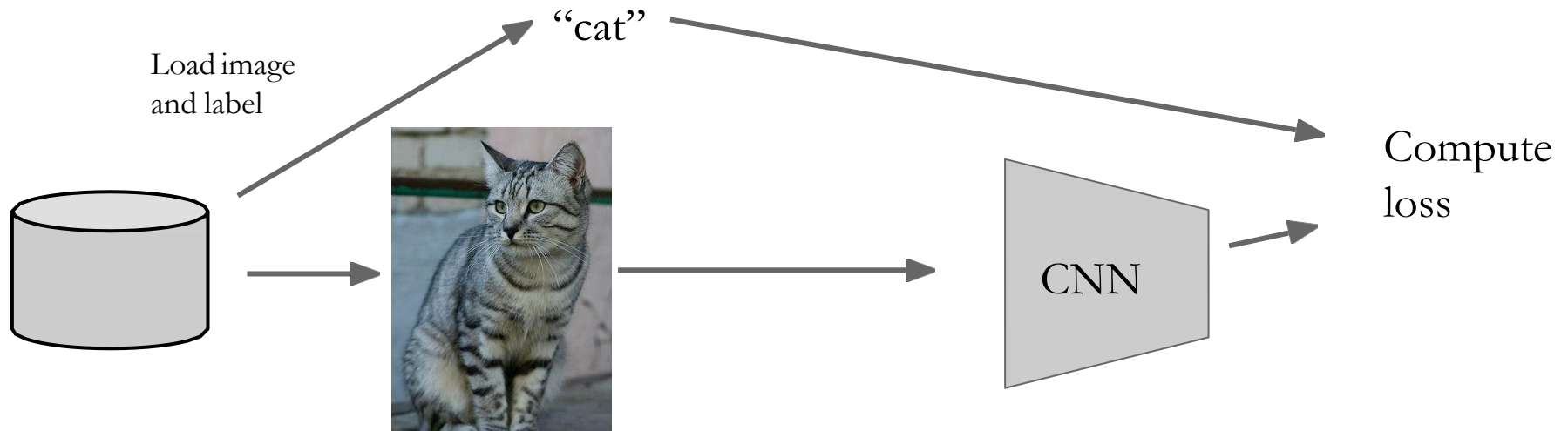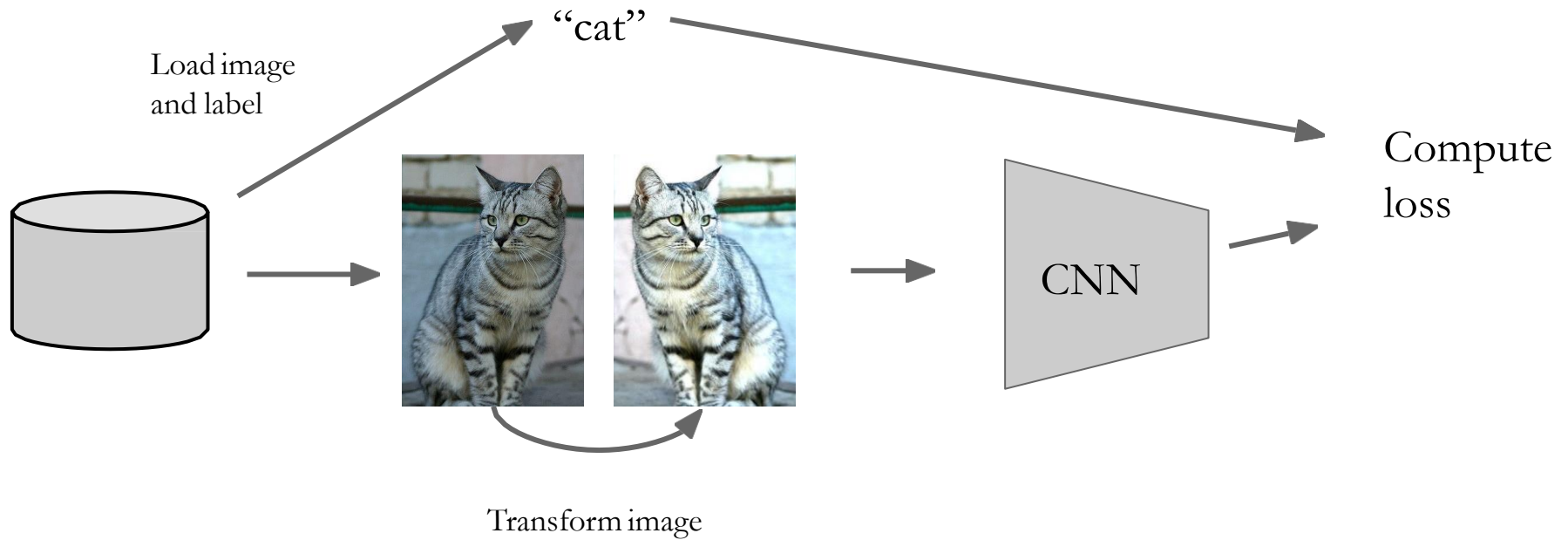


(a) Standard Neural Net  (b) After applying dropout.

Dropout: A simple way to prevent neural networks from overfitting [Srivastava JMLR 2014]

# Data Augmentation

Load image
and label

"cat"

CNN

Compute
loss

# Data Augmentation



Load image and label

"cat"

Transform image

CNN

Compute loss

# Data Augmentation

## Horizontal Flips

# Data Augmentation

## Random crops and scales

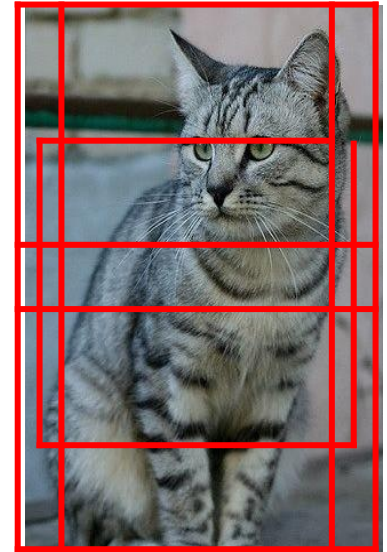**Training**: sample random crops / scales

ResNet:

1. Pick random L in range [256, 480]
2. Resize training image, short side = L
3. Sample random 224 x 224 patch

**Testing**: average a fixed set of crops

ResNet:

1. Resize image at 5 scales: {224, 256, 384, 480, 640}
2. For each size, use 10 224 x 224 crops: 4 corners + center, + flips

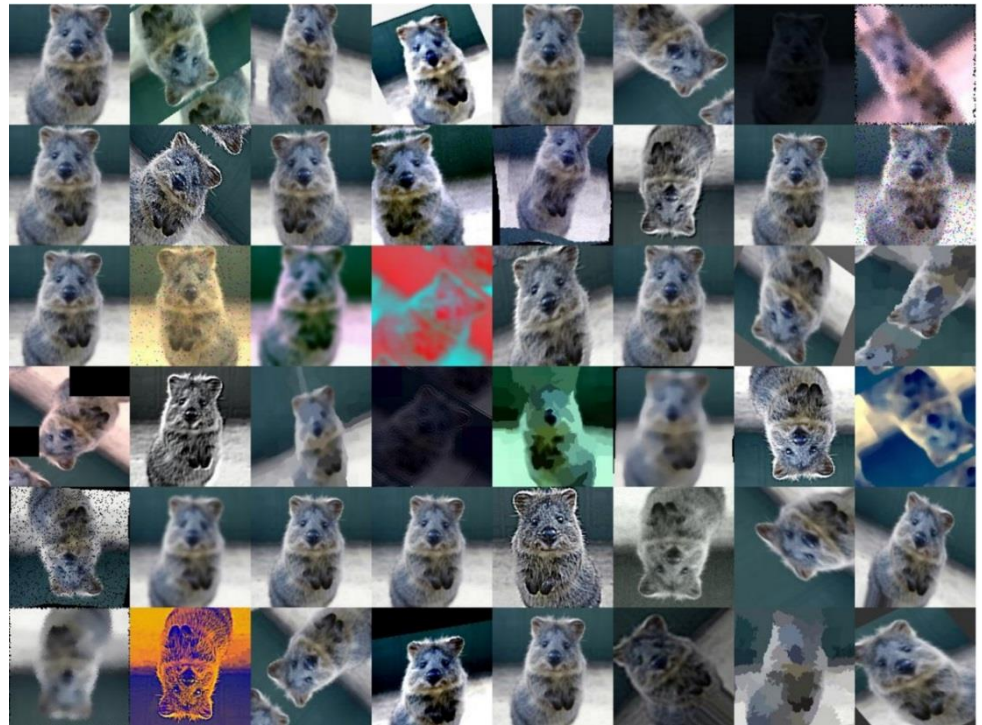# Data Augmentation

Get creative for your problem!

Random mix/combinations of :

- translation

- rotation

- stretching

- shearing,

- lens distortions

- …

# Transfer learning

- If you have sparse data in your domain of interest *(target),* but have rich data in a disjoint yet related domain *(source)*,

- You can train the early layers on the *source* domain, and only the last few layers on the *target domain:*



Set these to the already learned weights from another network

Learn these on your own task

# Transfer learning

Source: classify 20 animal classes          Target: 10 car classes

1. Train on source
(large dataset)

2. Small dataset:

3. Medium dataset:
**finetuning**

more data = retrain more of
the network (or all of it)

Freeze these

Freeze these

Train this

Train this

**Another option: use network as feature extractor,
train SVM/LR on extracted features for target task**

# Mini-batch gradient descent

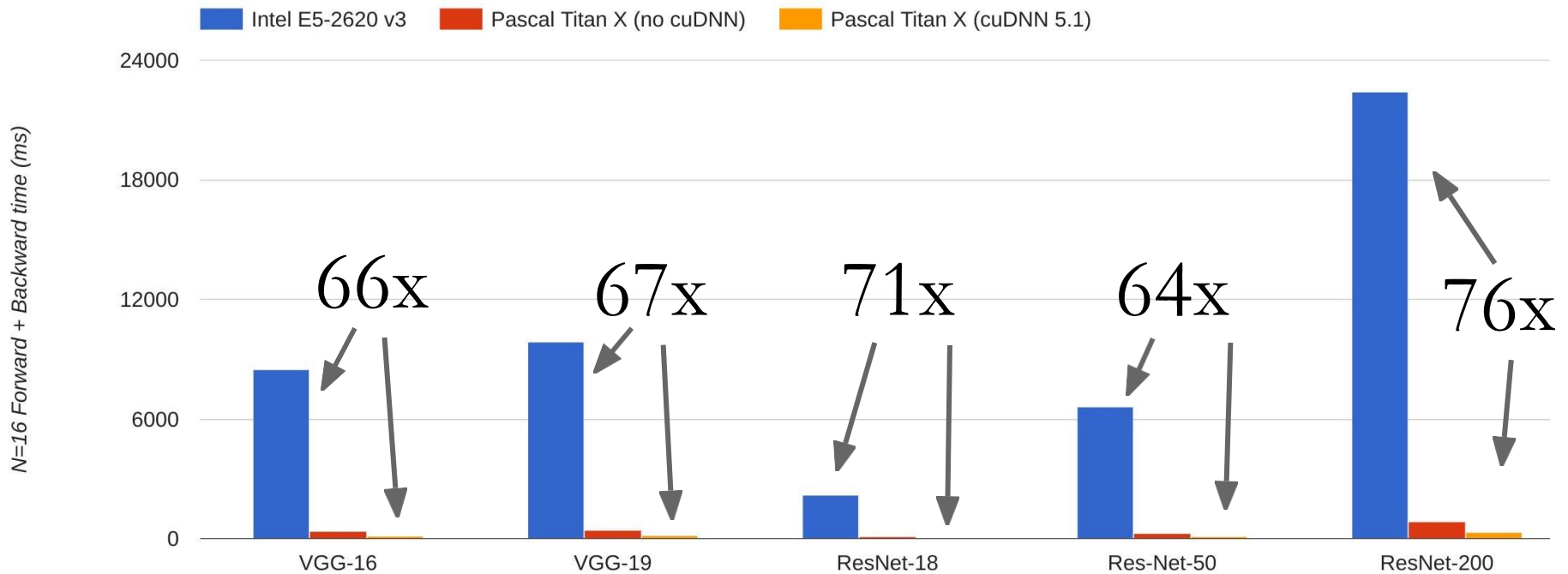- In classic gradient descent, we compute the gradient from the loss for all training examples

- Could also only use *some* of the data for each gradient update

- We cycle through all the training examples multiple times

- Each time we've cycled through all of them once is called an 'epoch'

- Allows faster training (e.g. on GPUs), parallelization

# Training: Best practices

- Center (subtract mean from) your data

- Use careful initialization for weights

- Use RELU or leaky RELU or ELU or PReLU

- Use batch normalization

- Use data augmentation

- Use regularization

- Use mini-batch

- Learning rate: too high? Too low?

- Use cross-validation for hyperparameters

# CPU vs GPU in practice

(CPU performance not
well-optimized, a little unfair)



Data from https://github.com/jcjohnson/cnn-benchmarks

Fei-Fei Li, Andrej Karpathy, Justin Johnson, Serena Yeung

# Software: A zoo of frameworks!

PaddlePaddle
(Baidu)

Chainer

Caffe
(UC Berkeley)

→ Caffe2
(Facebook)

MXNet
(Amazon)

CNTK
(Microsoft)

Torch
(NYU / Facebook)

→ PyTorch
(Facebook)

JAX
(Google)

Theano
(U Montreal)

→ TensorFlow
(Google)

And others...