# A whirlwind review/overview of deep learning (continued)

## CPSC 440/550: Advanced Machine Learning

`cs.ubc.ca/~dsuth/440/23w2`
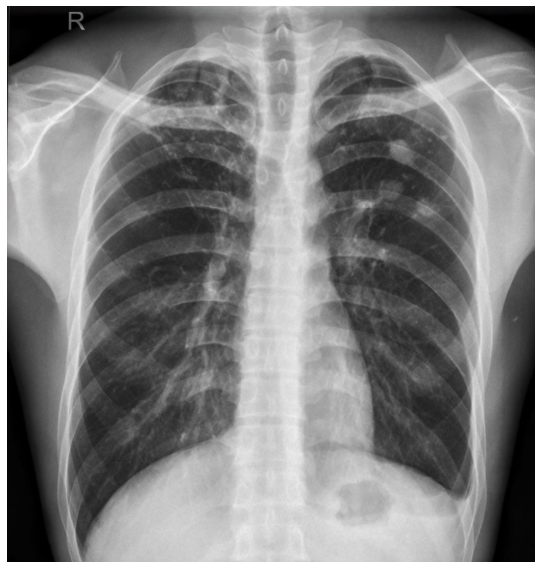
University of British Columbia, on unceded Musqueam land

2023-24 Winter Term 2 (Jan–Apr 2024)

# Next Topic: CNN review

# Motivation: X-ray abnormality detection

- Want to build a system that recognizes abnormalities in x-rays:



"Abnormality detected"
(binary classification)

- Applications:
  - Fast detection of tuberculosis, pneumonia, lung cancer, etc
- Deep learning has led to incredible progress on computer vision tasks
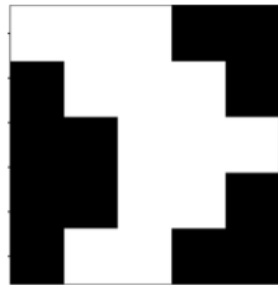  - Much of this progress driven by convolutional neural networks (CNNs)

# Convolutional Neural Network (CNN) motivation

- Consider training neural networks on 500 pixel by 500 pixel images
  - So the number of inputs $d$ to first layer is 250,000 (or 3x that, if colour)
- If first layer has $k$=10,000, then $W$ has <span style="color:red">2.5 billion parameters</span>
  - We want to avoid this huge number (due to storage and overfitting)

- <span style="color:blue">CNNs</span> <span style="color:green">drastically reduce the number of parameters</span> by:
  - Having <span style="color:green">activations only depend on a small number</span> of inputs
  - Using the <span style="color:green">same parameters on the connections</span> of many activations
- Done using layers that look like "<span style="color:blue">convolutions</span>" in signal processing

# Illustration of 2D Convolution

- ## 2D convolution:
  - Inputs: an "input" image $x$ and a "filter" image $w$
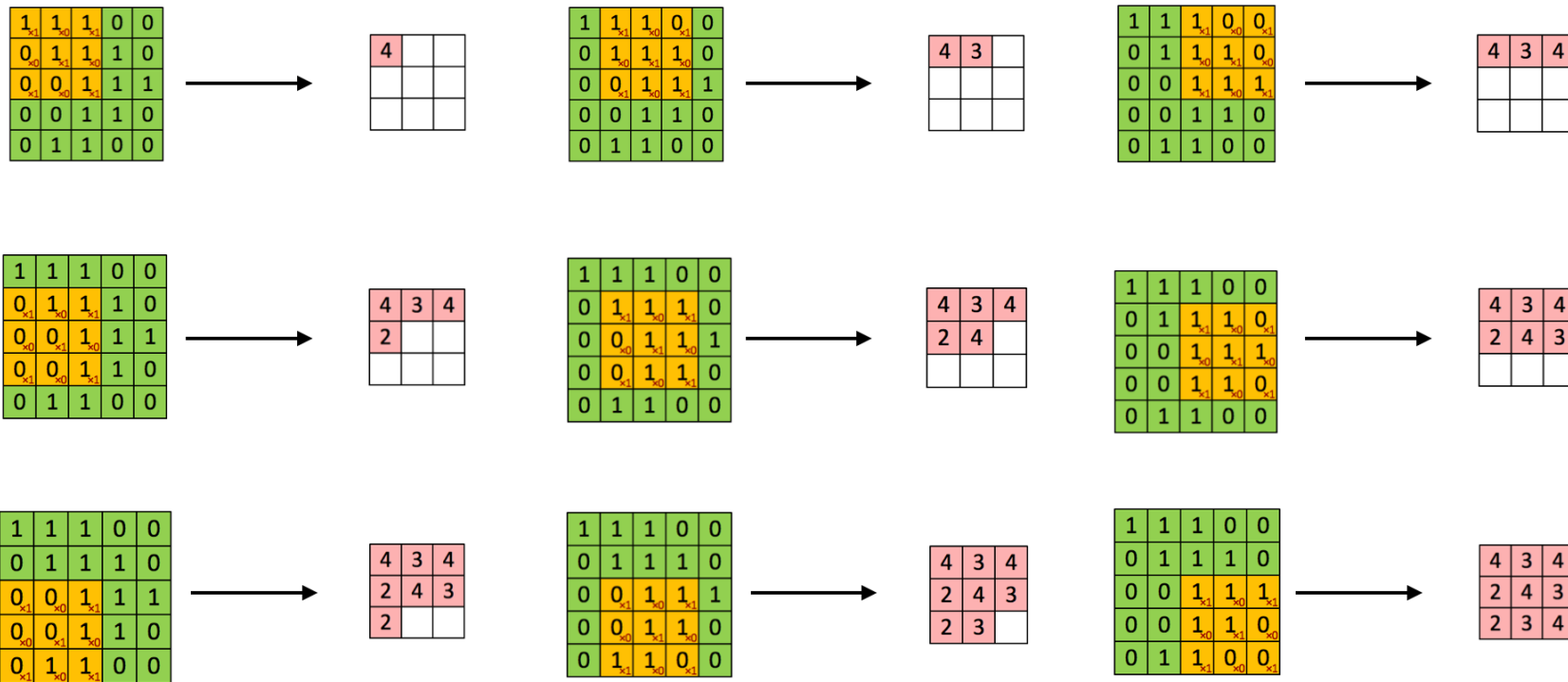  - Output: new image $z$ whose pixels are dot products of filter and image region)



**Input image**  **Filter image**  **Output image**

| 1 | 1 | 1 | 0 | 0 |
|---|---|---|---|---|
| 0 | 1 | 1 | 1 | 0 |
| 0 | 0 | 1 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

$x$

*

| 1 | 0 | 1 |
|---|---|---|
| 0 | 1 | 0 |
| 1 | 0 | 1 |

$w$

=

| 4 | 3 | 4 |
|---|---|---|
| 2 | 4 | 3 |
| 2 | 3 | 4 |

$z$

6

# Illustration of 2D Convolution

- ## 2D convolution:
  - Inputs: an "input" image *x* and a "filter" image *w*
  - Output: new image *z* whose pixels are dot products of filter and image region)
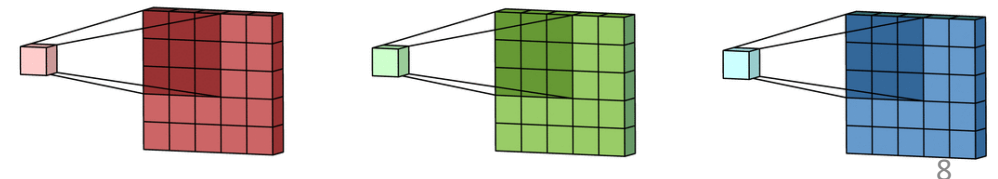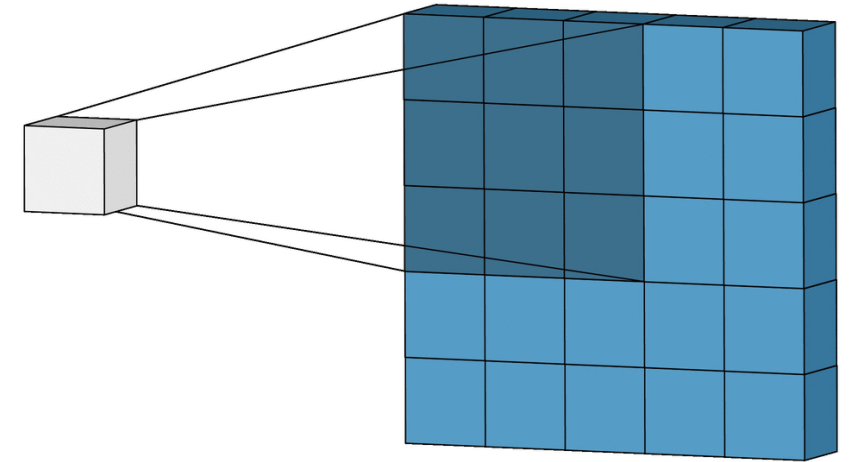
**Filter image**

# Illustration of 2D Convolution

- **2D convolution**:
  - Inputs: an "input" image *x* and a "filter" image *w*
  - Output: new image *z* whose pixels are dot products of filter and image region)
- As a formula:

$$z[i_1, i_2] = \sum_{j_1=-m}^{m} \sum_{j_2=-m}^{m} w[j_1, j_2] x[i_1 + j_1, i_2 + j_2]$$

  - Final image *z* can be written as usual *z=W' x*
    - *W'* will be sparse, with filter values in *W* repeated
- **3D convolution** (for colour images):
  - Weighted dot product across all three dimensions

# Formal Convolution Definition

- We have defined the convolution as:

$$z_i = \sum_{j=-m}^{m} w_j x_{i+j}$$

- In other classes you may see it defined as:

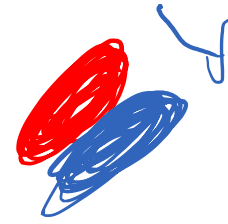$$z_i = \sum_{j=-m}^{m} w_j x_{i-j}$$

(reverses 'w')

$$z_i = \int_{-\infty}^{\infty} w_j x_{i-j} \, dj$$

(assumes signal + filter are continuous)

- For simplicity we're skipping the "reverse" step,
  and assuming $w$ and $x$ are sampled at discrete points (not functions)

- But keep this mind if you read about convolutions elsewhere
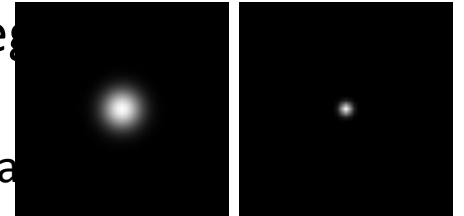
# Convolutions

- Pre-2012, people often designed the filters by hand
  - Filters can approximate "derivatives" or "integrals" of the image regions.
    - Derivative filters will up to 0, integral filters will add up to 1
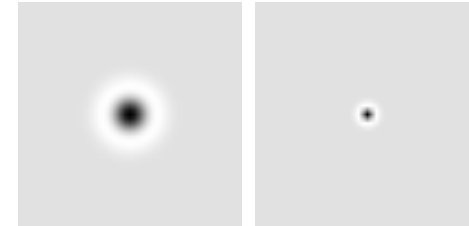  - Three of the most-common filters that people used:
    - Gaussian filters: integral filter, giving the average brightness in a reg
      - Variance of the Gaussian controls the amount of smoothness
      - This produces a pixel feature that is less sensitive to noise than pixel's raw va
    - Gabor filters: derivative filters, measuring changes in brightness along a direction
      - We typically compute these for different orientations and "frequencies"
      - This gives a set of features that is useful in describing edges in the image
    - Laplacian of Gaussian filter: total second-derivative filter
      - Complements Gabor filters: helps describe if change is due to an edge, line, or continuous change
  - Similar filters may be used early in the eyes visual processing
  - Results of convolutions kind of like the "bag of words" making up images

10

# Image Convolution Examples



Gaussian Convolution:

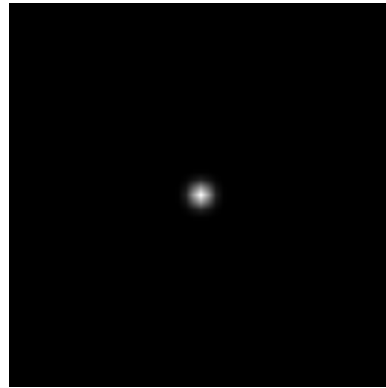blurs image to represent average (smoothing)

# Image Convolution Examples



Gaussian Convolution:

$*$

(smaller variance)
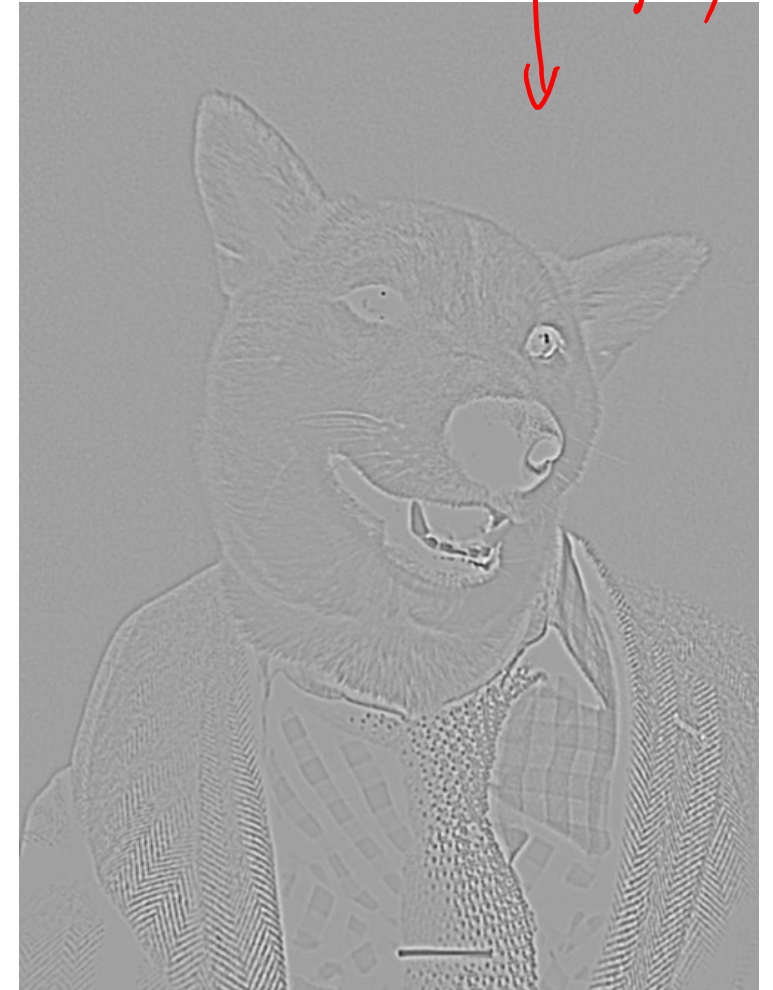
$=$

blurs image to represent average (smoothing)

# Image Convolution Examples



Laplacian of Gaussian

$*$     $=$

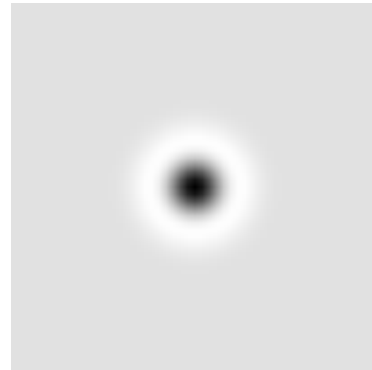"How much does it look like a black dot surrounded by white?"

"signed" image (gray is 0)

# Image Convolution Examples



Laplacian of Gaussian

(larger variance)

Similar preprocessing may be done in basal ganglia and LGN.
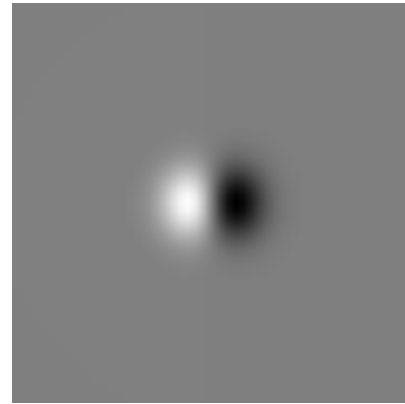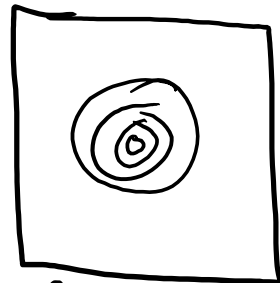
Black/white as sides of edge

# Image Convolution Examples



Gabor Filter
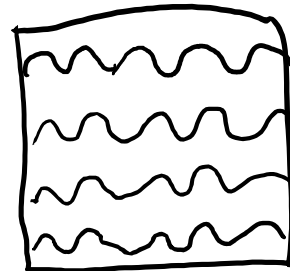(Gaussian multiplied by sine or cosine)

\*  =

||

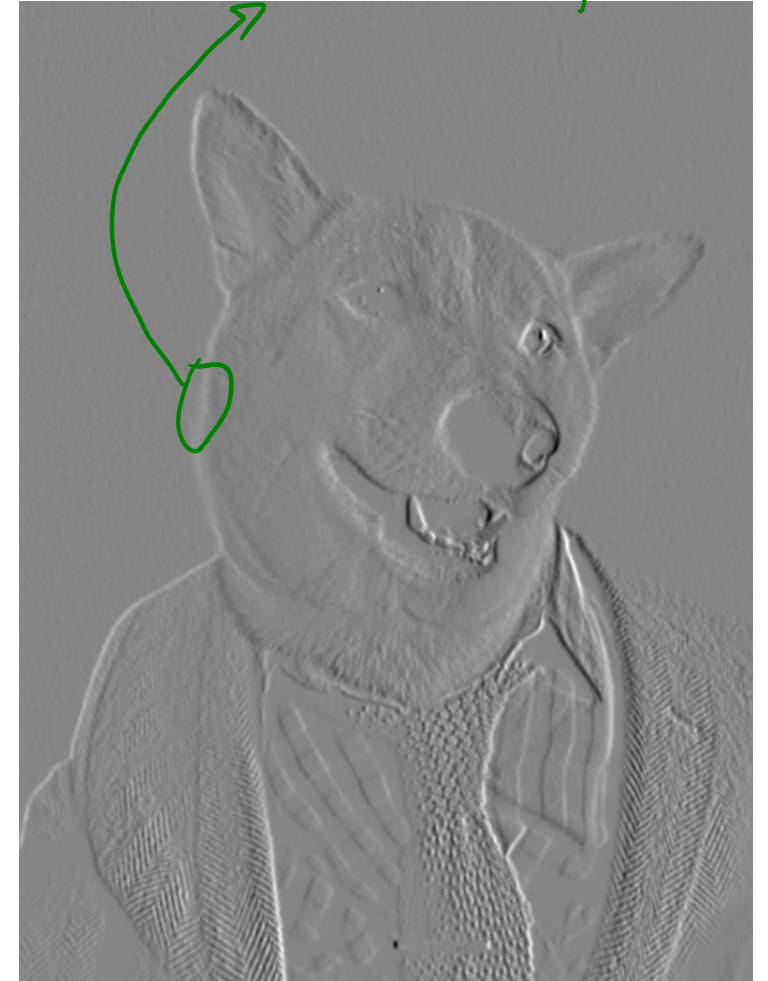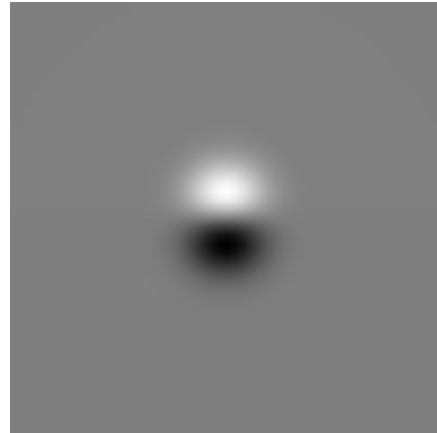Gaussian  .\*  Parallel Sine functions

Horizontal "bright to dark"

15

# Image Convolution Examples



Gabor Filter
(Gaussian multiplied by sine or cosine)

Different orientations of the sine/cosine let us detect changes with different orientations.

2d derivatives have a direction.

# Image Convolution Examples



Gabor Filter
(Gaussian multiplied by
sine or cosine)

*

(smaller variance)

=

# Image Convolution Examples



Gabor Filter
(Gaussian multiplied by
sine or cosine)

∗    =

(smaller variance)

Vertical orientation

—Can obtain other orientations by
rotating.

—May be similar to effect of V1 "simple cells."

# Unsupervised Learning of Filters for Image Patches

- Consider building an unsupervised model of image patches:



$X =$

$x_i$

Size of $X$: (image height) × (image width) by ((patch height) × (patch width) × 3)

# Unsupervised Learning of Filters for Image Patches

*bonus!*

- Some methods to do this generate Gaussian/LoG/Gabor filters:
  - These filters are motivated from both neuroscience and ML experiments.



*"colour opponency"*

(c) With whitening - gray.    (d) With whitening - RGB.

http://lear.inrialpes.fr/people/mairal/resources/pdf/review_sparse_arxiv.pdf

# Motivation for Convolutional Neural Networks

- Classic vision methods use <span style="color:red">fixed convolutions</span> as features:
  - Usually have <span style="color:green">different types/variances/orientations</span>
  - Can do subsampling or take <span style="color:green">maxes across locations/orientations/scales</span>

# Motivation for Convolutional Neural Networks

- Convolutional neural networks learn the convolutions:
  - Learning $W$ and $v$ automatically chooses types/variances/orientations
  - Don't pick from fixed convolutions, but learn the elements of the filters

# Motivation for Convolutional Neural Networks

- **Convolutional neural networks** learn the convolutions:
  - Learning $W$ and $v$ automatically chooses types/variances/orientations
  - Can do multiple layers of convolution to get deep hierarchical features

23

# Convolutional Neural Networks

- Classic architecture of a convolutional neural network:



- **Convolution layers**:
  - Apply convolution with several different filters
  - Sometimes these have a "stride": skip several pixels between applying filter
- **Pooling layers**:
  - Aggregate regions to create smaller images (usually "max pooling")
- **Fully-connected layers**: usual "multiplication by $W_i$" in layer

"Stride" of 2

# Max Pooling Example

- Max pooling:

$z_i$:



- Decreases size of hidden layer, so we need fewer parameters
  - Gives some local translation invariance:
    - The precise location of max is not important
- This is continuous and piecewise-linear but non-differentiable
  - Like ReLU, we can still optimize this type of objective with SGD

# LeNet Convolutional Neural Networks

- Classic convolutional neural network (LeNet):



- Visualizing the "activations":
  - http://scs.ryerson.ca/~aharley/vis/conv
  - http://cs231n.stanford.edu

26

# ImageNet Competition

- **ImageNet**: Millions of labeled images, 1000 object classes
  - Task is to classify images into one of the 1000 class labels.
    - We will discuss multi-class classification in Part 2 of the course.
  - Everyone submits their "best" model, winners announced.

https://www.youtube.com/watch?v=40riCqvRoMs

# AlexNet Convolutional Neural Network

*bonus!*

- Modern CNN era started with AlexNet (won 2012 competition):
  - 15.4% error vs. 26.2% for closest competitor
  - 5 convolutional layers
  - 3 fully-connected layers
  - SG with momentum
  - ReLU non-linear functions
  - Data translation/reflection/ cropping
  - L2-regularization + Dropout
  - 5-6 days on two GPUs
    - at the time – could run much faster with current hardware



Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input is 150,528-dimensional, and the number of neurons in the network's remaining layers is given by 253,440–186,624–64,896–64,896–43,264–4096–4096–1000.

# ImageNet Insights

- Filters and stride got smaller over time
  - Popular VGG approach uses 3x3 convolution layers with stride of 1
    - 3x3 followed by 3x3 simulates a 5x5, and another 3x3 simulates a 7x7, and so on
    - Speeds things up and reduces number of parameters
    - Also increases number of non-linear ReLU operations

# ImageNet Insights

- Filters and stride got smaller over time.
  - Popular VGG approach uses 3x3 convolution layers with stride of 1.
  - GoogLeNet used multiple filter sizes ("inception layer"), but not as popular.
- Eventual switch to "fully-convolutional" networks.
  - No fully connected layers.
- ResNets allow easier training of deep networks.
  - Won all 5 tasks in 2015, training 152 layers for 2-3 weeks on 8 GPUs.
- Ensembles help.
  - 2016 winner combined predictions of previous networks.
- Competition ended in 2017!

# Discussion of CNNs

- Convolutional layers reduce the number of parameters in two different ways:
  - Each hidden unit only depends on small number of inputs from previous layer
  - We use the same filters across the image
    - So we do not learn a different weight for each "connection" like in classic neural networks

- CNNs give some amount of translation invariance/equivariance:
  - Because the filters are used across the image, they can detect a pattern anywhere in the image
    - Even in image locations where the pattern has never been seen
  - The pooling layer can also give some local invariance, against small translations of the image

- CNNs are not only for images!
  - Can use CNNs for 1D sequences like sound or language
  - Can use CNNs for 3D objects like videos or medical image volumes
  - Can use CNNs for graphs
- But you do need some notion of "neighbourhood" for convolutions to make sense.

# Next Topic: Autoencoders

# Autoencoders



- Autoencoders try to make their output the same as the input
    - Usually have a bottleneck layer with dimension $k$ < input $d$
    - First layers "encode" the input into bottleneck
    - Last layers "decode" the bottleneck into a (hopefully valid) input

# Autoencoders



- This is an unsupervised learning method.
  - There are no labels $y$
- Relationship to principal component analysis (PCA):
  - With squared error and linear network, equivalent to PCA
    - Size of bottleneck layer gives number of latent factors $k$ in PCA
  - With non-linear transforms: a non-linear/deep generalization of PCA

# Encoder as Learning a Representation

- Consider the encoder part of the network:
  - Takes features $x^i$ and makes low-dimensional $z^i$

- Ways you could use the encoder:
  - Use $z^i$ as compressed input (reduce memory needed)
  - Set bottleneck size to 2, and plot the $z^{(i)}$ to visualize the data
  - Try to interpret what the bottleneck features $z^{(i)}$ mean
  - Use the $z^{(i)}$ as features for supervised learning
    - For the special case of PCA and regression with L2 loss, this is called "partial least squares"
  - You could add a supervised $y^i$ to final layer of trained autoencoder + fit with SGD
    - This is called "unsupervised pre-training"
    - If you use unlabeled data to do this initialization, an example of "self-supervised" learning
      - Usually it is easier to get a lot of unlabeled data than it is to get labeled data.

# PCA vs. Deep Autoencoder (Document Data)

(these days t-SNE is the usual way to make visualizations like this; see these guidelines)

# Unsupervised pre-training

- One version of this: semi-supervised learning
    - 1. Train an autoencoder on all of Instagram (e.g.)
    - 2. Use its features for softmax regression on your (small) labeled dataset
- Could also "fine-tune" the whole network;
  this would be one version of unsupervised pre-training

# Some clarification about unsupervised pre-training *bonus!*

- What "unsupervised pre-training" used to mean
- Old scheme for deep networks: **stacked** denoising autoencoders
  - Train a two-layer denoising autoencoder
  - Freeze the encoder layer as first layer of your deep net
  - Train a denoising autoencoder on activations from that layer
  - Freeze its encoder as the second layer of your deep net
  - Repeat
  - Fine-tune with SGD at the end

- People don't do this anymore: we can do end-to-end SGD now

# Decoder as Generative Model

- Consider the decoder part of the network:
  - Takes low-dimensional $z^{(i)}$ and makes features $\hat{x}^{(i)}$

- Can be used for outlier detection:
  - Check distance to original features to detect outliers

- Can be used to generate "new data":
  - If the decoder is good, new values of $z$ that "look like real $z$" should decode into $\hat{x}$ that "look like real $x$"
  - To do this "properly," need to estimate the distribution p(z)
    - This is what "Stable Diffusion" does

# Font Manifold

- Going from encoding to decoding for different fonts:



Please drag the black and white circle around the heat map to explore the 2D font manifold.

Unlikely    Probability    Likely

Select Character: t

- Demo [here](#).
  - The above was generated by a Gaussian process and not an autoencoder.
  - But the decoder part of autoencoders is trying to do something like this.

# Latent Space Interpolation



- Encode both ends; decode various points on a line between

# Neural Networks with Multiple Outputs

- Previous neural networks we have seen only have 1 output $y$.
- In autoencoders, we have $d$ outputs (one for each feature).

$$\hat{x}_1 = v_1^T h(W^3 h(W^2 h(W^1 x)))$$
$$\hat{x}_2 = v_2^T h(W^3 h(W^2 h(W^1 x)))$$
$$\hat{x}_d = v_d^T h(W^3 h(W^2 h(W^1 x)))$$

$$\hat{x} = V h(W^3 h(W^2 h(W^1 x)))$$

- For training, we add up the loss across all $j$:

$$= v_j^T h(W^3 h(W^2 h(W^1 x^i)))$$

$$f(W^1, W^2, V) = \sum_{i=1}^{n} \sum_{j=1}^{d} \left(\hat{x}_j^i - x_j^i\right)^2$$

squared error for continuous $x_j$

$$f(W^1, W^2, V) = \sum_{i=1}^{n} \sum_{j=1}^{d} \log\left(1 + \exp\left(-\hat{x}_j^i x_j^i\right)\right)$$

logistic loss for binary $x_j^i \in \{-1, +1\}$

- Fit with SGD (sampling random $i$), and usual deep learning tricks can be used
  - Even though network has multiple outputs, $f$ is a scalar so autodiff works as before
  - For images, may want to use convolution layers

# Denoising Autoencoders

- A common variation on autoencoders is denoising autoencoders:
  - Use "corrupted" inputs, and learn to reconstruct uncorrupted originals



  - "Learn a model that removes the noise". Easy to get lots of training data.
    - You can apply the model to denoise new images.
    - Do not necessarily need a "bottleneck" layer.

# What Denoising Autoencoders Learn

**Theorem 1** *Let $p$ be the probability density function of the data. If we train a DAE using the expected quadratic loss and corruption noise $N(x) = x + \epsilon$ with*

$$\epsilon \sim \mathcal{N}\left(0, \sigma^2 I\right),$$

*then the optimal reconstruction function $r^*(x)$ will be given by*

$$r^*(x) = \frac{\mathbb{E}_\epsilon\left[p(x - \epsilon)(x - \epsilon)\right]}{\mathbb{E}_\epsilon\left[p(x - \epsilon)\right]} \tag{3}$$

*for values of $x$ where $p(x) \neq 0$.*

   *Moreover, if we consider how the optimal reconstruction function $r_\sigma^*(x)$ behaves asymptotically as $\sigma \to 0$, we get that*

$$r_\sigma^*(x) = x + \sigma^2 \frac{\partial \log p(x)}{\partial x} + o(\sigma^2) \quad as \quad \sigma \to 0. \tag{4}$$

Alain and Bengio (2012)

- Can use to estimate "Hyvärinen score" $\frac{(r_\sigma^*(x) - x)}{\sigma^2} \approx \nabla_x \log p(x)$
- Closely related to diffusion models (later in the course!)

# Image Colourization

Colorado National Park, 1941     Textile Mill, June 1937     Berry Field, June 1909     Hamilton, 1936

- Gallery: http://iizuka.cs.tsukuba.ac.jp/projects/colorization/extra.html
- Video: https://www.youtube.com/watch?v=ys5nMO4Q0iY

http://iizuka.cs.tsukuba.ac.jp/projects/colorization/en/

# Image Colourization

- Instead of noisy inputs, you use de-coloured inputs:



- Another application is super-resolution:
  - Learn to output a high-resolution image based on low-resolution images.

# Next Topic: Multi-Label Classification

# Motivation: Multi-Label Classification

- Consider multi-label classification:

$$X = \begin{bmatrix} \phantom{xxxxxxxxxxxx} \end{bmatrix} \Big\} n$$

$$\underbrace{\phantom{xxxxxxxxxxxx}}_{d}$$

cat dog person chair mouse

$$Y = \begin{bmatrix} 1 & \sim1 & 1 & 1 & 1 \\ \sim1 & \sim1 & 1 & -1 & \sim1 \\ \sim1 & -1 & \sim1 & 1 & -1 \\ 1 & 1 & \sim1 & 1 & \sim1 \\ -1 & \sim1 & 1 & 1 & \sim1 \\ 1 & \sim1 & 1 & -1 & 1 \end{bmatrix} \Big\} n$$

$$\underbrace{\phantom{xxxxxxxxxxxx}}_{k}$$

- Which of the *k* objects are in this image?
  - There may be more than one "correct" class label.

# Independent Classifier Approach

- One way to build a multi-label classifier:
  - Train a <span style="color:green">classifier for each label</span>
    - Train a neural network that predicts +1 if the image contains a dog, and -1 otherwise
    - Train a neural network that predicts +1 if the image contains a cat, and -1 otherwise
    - …
  - To make predictions for the $k$ classes, <span style="color:green">concatenate predictions</span> of the $k$ models

- Can think of this as a "product of independent classifiers"

- Drawbacks:
  - <span style="color:red">Lots of parameters</span>: $k*$(number of parameters for base classifier).
  - Each classifier needs to <span style="color:red">"relearn from scratch"</span>
    - Each classifier needs to learn its own Gabor filters, how corners and light works, and so on
    - A lot of visual <span style="color:green">features for "dog" might also help us predict "cat"</span>

# Encoding-Decoding for Multi-Label Classification

- Multi-label classification with an encoding-decoding approach:
  - Input is connected to a hidden layer.
  - Hidden layer is connected to multiple output units.



$$\hat{y}_1 = v_1^T h(W_x)$$

$$\hat{y}_2 = v_2^T h(W_x)$$

$$\hat{y}_3 = v_3^T h(W_x)$$

- Prediction: compute hidden layer, compute activations, compute output:

$$\hat{y} = V h(W_x)$$

- Number of parameters and cost is O($dm + mk$) for $k$ classes and $m$ hidden units.
  - If we trained a separate network for each class, number of parameters and cost would be O($kdm$) (for 'W' for each class)
- Might have multiple layers, convolution layers, and so on.
- No need to have a "bottleneck" layer – "encoder"/"decoder" is just terminology

# Encoding-Decoding for Multi-Label Classification



- We usually assume that the classes are independent given last layer:

$$p\left(y_1, y_2, \ldots, y_k \mid x_1, x_2, \ldots, x_d, W, V\right) = p\left(y_1 \mid x_1, x_2, \ldots, x_d, W, v_1\right) p\left(y_2 \mid x_1, x_2, \ldots, x_d, W, v_2\right) \cdots p\left(y_k \mid x_1, x_2, \ldots, x_d, W, v_k\right)$$

with: $\quad p(y_1 = 1 \mid x, W, v) = \dfrac{1}{1 + \exp\left(-v_1^T h(Wx)\right)} \quad \underbrace{\phantom{xxxx}}_{\Theta_1} \quad p(y_2 = 1 \mid x, W, v) = \dfrac{1}{1 + \exp\left(-v_2^T h(Wx)\right)} \quad \underbrace{\phantom{xxxx}}_{\Theta_2} \quad \cdots$

  - Conditioned on features/parameters, this is ultimately a fancy product of Bernoullis model:
    - $p(y_1, y_2, \ldots, y_k \mid x, W, V) = p(y_1 \mid x, W, V)p(y_2 \mid x, W, V) \cdots p(y_k \mid x, W, V)$, where $p(y_c = 1 \mid x, W, V) = \theta_c$.
    - This makes decoding and other inference problems easy: you do inference on each $y_c$ independently.

# Encoding-Decoding for Multi-Label Classification

- The negative log-likelihood we optimize for MLE:

$$f(W, V) = \sum_{i=1}^{n} \sum_{c=1}^{x} \log\left(1 + \exp\left(-y_c^i \, v_c^\top h(W x^i)\right)\right)$$

- Use backpropagation or AD to compute gradient, train by SGD.
  - You randomly sample a training example $i$ and compute gradient for all labels
  - The updates of $W$ lead to features that are useful across classes
  - The updates of $V$ focus on getting the class labels right given the features

- Important:
  - We assumed independence of labels given the last layer
  - But the last layer can reflect dependencies
    - If "dog" and "human" are frequently together, this should be reflected in the hidden layer
      - For example, $\theta_{human}$ might be higher when the features give a high value for $\theta_{dog}$

# Pre-Training for Multi-Label Classification

- Consider a scenario where we get a new class label
  - For example, we get new images that contain horses (not seen in training)



- Instead of training from scratch, we could:
  - Add an extra set of weights $v_{k+1}$ to the final layer for the new class
  - Train these weights with the encoding weights $W$ fixed
    - This is a simple/convex logistic regression problem
    - If we already have "features" that are good for many classes,
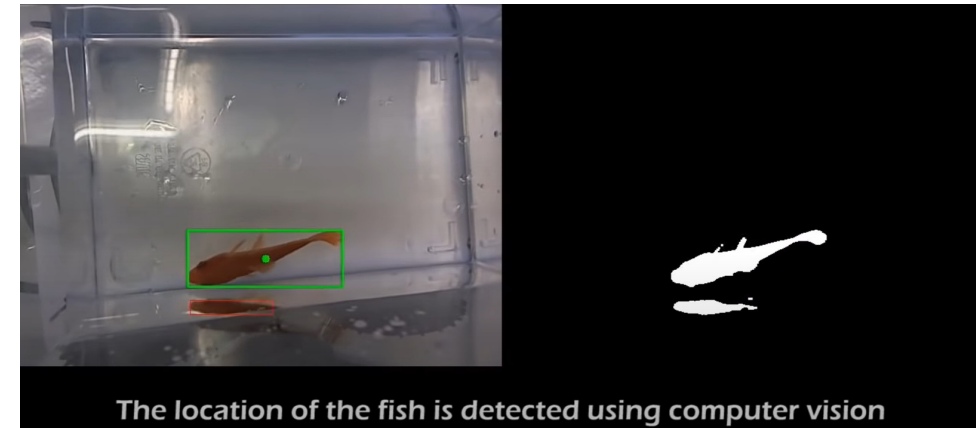      we may be able to learn a new class with very few training examples!

# Pre-Training for Multi-Label Classification

- Using an existing network for new problems is called "pre-training"
  - Typically, we start with a network trained on a large dataset
  - We use this network to give us features to fit a smaller dataset
    - "Few-shot learning"

- Depending the setup, you may also update $W$ and the other $v_c$
  - Useful if you have a lot of data on the new class
  - In this case, would typically mix in new examples with old ones

- Increasing trend in vision and language to using pre-training a lot
  - No need to learn everything about language for every language task!

# Next Topic: Fully-Connected Networks

# Motivation: Pixel Classification

- Suppose we want to assign a binary label to each pixel in an image:
  - Tumour vs. non-tumour, pedestrian vs. non-pedestrian, and so on



The location of the fish is detected using computer vision

- How can we use CNNs for this problem?

# Naïve Approach: Sliding Window Classifier

- Train a CNN that predicts pixel label given its neighbourhood



- To label all pixels, apply the CNN to each pixel
  - Advantages:
    - Turns pixel labeling into image classification
    - Can be applied to images of different sizes
  - Disadvantage: this is slow
    - (Cost of applying CNN) * (number of pixels in the image)

# Encoding-Decoding for Pixel Classification

- Similar to multi-label, could use CNN to generate an image encoding
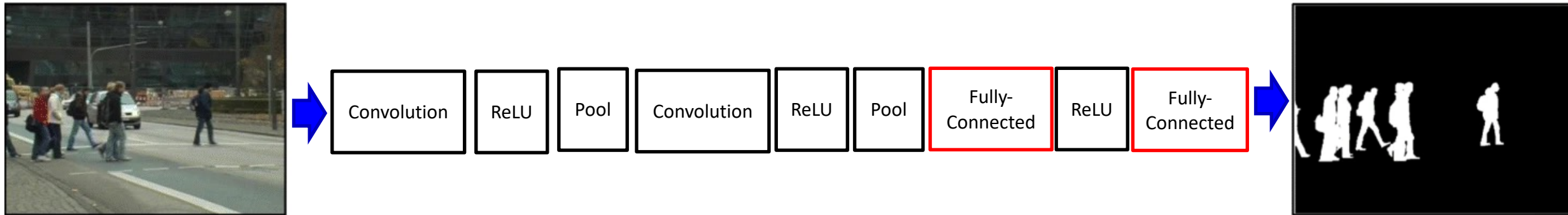  - With output layer making a prediction at each pixel



- Much-faster classification
  - Small number of "global" convolutions, instead of repeated "local" convolutions
- But, the encoding has mixed up all the space information
  - Fully-connected layers throw all that information out
  - Fully-connected layer needs to learn "how to put the image together"
    - And images must be the same size

# Fully-Convolutional Networks

- Fully-convolutional networks (FCNs):
  - CNNs with no fully-connected layers (only convolutional and pooling)



- Maintains fast classification of the encoding-decoding approach

- Same parameters used across space at all layers
  - This allows using the network on inputs of different sizes
  - Needs upsampling layer(s) to get back to image size

- FCNs quickly achieved state of the art results on many tasks

# Traditional Upsampling Methods

- In upsampling, we want to go from a small image to a bigger image
  - This requires interpolation: guessing "what is in between the pixels"
- Classic upsampling operator is nearest neighbours interpolation:
  - But this creates blocky/pixelated images



**Nearest Neighbor**

Input: 2 x 2  →  Output: 4 x 4

# Traditional Upsampling Methods

- Another classic method is bilinear interpolation:
  - Weighted combination of corners:
  - More smooth methods include "bicubic" and "splines"
- In FCNs, we learn the upsampling/interpolation operator

# Upsampling with Transposed Convolution

- FCN Upsampling layer is implemented with a transposed convolution
  - Sometimes called "deconvolution" in ML or "fractionally-strided convolution"
    - But not related to deconvolution in signal processing



Figure 3. Illustration of deconvolution and unpooling operations.

- Convolution generates 1 pixel by taking weighted combination of several pixels
  - And we learn the weights
- Transposed convolution generates several pixels by weighting 1 pixel
  - And we learn the weights
  - This generates overlapping regions, which get added together to make final image

# Upsampling with Transposed Convolution



Convolution

Each output (green) is a weighted combination of one 3x3 region

Transposed convolution

Each input (blue) gives a 3x3 region. These regions overlap, so we take a weighted combination at each pixel to give final result.

- Animations [here](#) and [here](#).

# Why is it called "transposed" convolution?

- We can write the convolution operator as a matrix multiplication, z = W' x.



- In transposed convolution, non-zero pattern of 'W' is transposed from convolution
  - You can implement transposed convolution as a convolution



- In this example the filter is the same, but does not need to be:
  - Transposed convolution is not the "reverse" of convolution (it only "reverses" the size)

# Increasing Resolution: FCN Skip Connections

*bonus!*

- Convolutions and pooling <span style="color:red">lose a lot of information</span>
- Original FCN paper considered adding <span style="color:green">skip connections</span> to help upsampling:



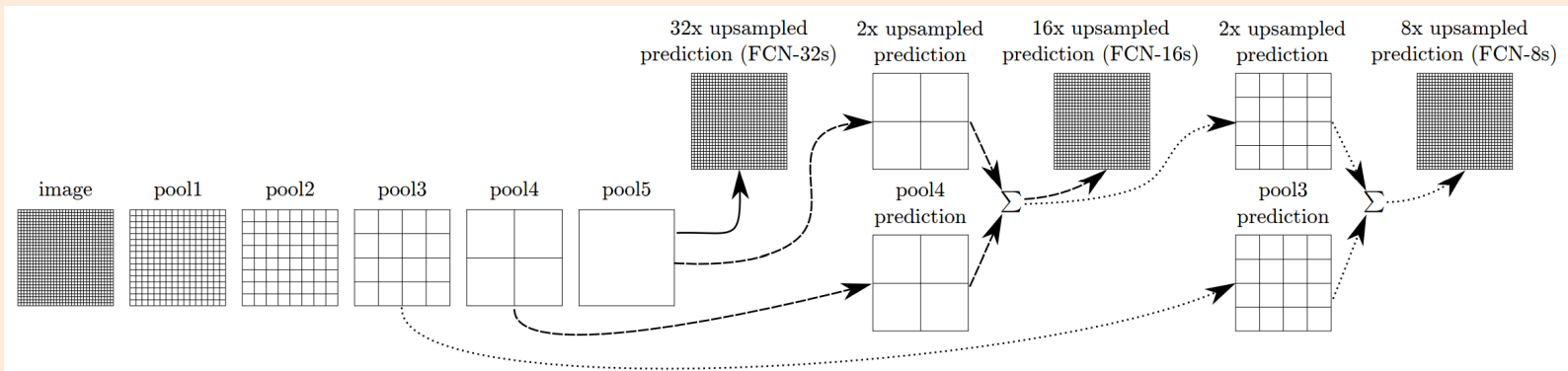Figure 3. Our DAG nets learn to combine coarse, high layer information with fine, low layer information. Layers are shown as grids that reveal relative spatial coarseness. Only pooling and prediction layers are shown; intermediate convolution layers (including our converted fully connected layers) are omitted. Solid line (FCN-32s): Our single-stream net, described in Section 4.1, upsamples stride 32 predictions back to pixels in a single step. Dashed line (FCN-16s): Combining predictions from both the final layer and the pool4 layer, at stride 16, lets our net predict finer details, while retaining high-level semantic information. Dotted line (FCN-8s): Additional predictions from pool3, at stride 8, provide further precision.

- Allows using <span style="color:green">high-resolution information from earlier</span> layers
- They first trained the low-resolution FCN-32, then FCN-16, then FCN-8
  - "First learn to encode at a low resolution", then slowly increase resolution
  - Parameters of transposed convolutions initialized to simulate "bilinear interpolation"

https://arxiv.org/pdf/1411.4038.pdf

# Increasing Resolution: Deconvolution Networks

*bonus!*

- Alternate resolution-increasing method is deconvolution networks:



- Includes transposed convolution layers and unpooling layers
  - Store the max pooling argmax values
  - Restores "where" activation happened
    - Still loses the "non-argmax" information



https://www.cv-foundation.org/openaccess/content_iccv_2015/papers/Noh_Learning_Deconvolution_Network_ICCV_2015_paper.pdf
https://towardsdatascience.com/transposed-convolution-demystified-84ca81b4baba

# Increasing Resolution: U-Nets

*bonus!*

- Another popular variant is u-nets:



- Decoding has connections to same-resolution encoding step

# Source of Labels

- Labeling every pixel takes a <span style="color:red">long time</span>.

- Where we get the labels from?
  - Domain expert (medical images)
  - Grad students or paid labelers (ImageNet)
  - Simulated environments
    - <span style="color:green">High number</span> of <span style="color:red">lower-quality</span> examples
    - Often a net gain with fine-tuning on real images
    - Can get data at night, in fog, or dangerous situations

Video game

Google street view

# Source of Labels

- Recent works recognize you <span style="color:red">do not need to label every pixel</span>
  - You can evaluate loss/gradient on a subset of labeled pixels
  - Could have labeler click on a few pixels inside objects, and a few outside
    - Many variations are possible, that let you label a lot of images in a short time
- Penguin counting based "single pixel" labels in training data:
  - And some tricks to separate objects and remove false positives:

# End of Part 1 ("Binary Variables"): Key Concepts

- We discussed binary density estimation
  - Model the proportion of times a binary event happens
- We discussed the Bernoulli parameterization
- We discussed various inference tasks, given the parameter:
  - Compute probabilities, find decoding, generate samples
- We discuss different learning strategies, given data:
  - Maximum likelihood estimation (MLE), maximum a posteriori (MAP)
  - Beta distribution as a prior gives a beta distribution as posterior ("$\propto$")
- We discussed modeling binary variables conditioned on features:
  - Tabular parameterization is flexible but has too many parameters
  - Logistic regression is limited but has a linear number of parameters

# End of Part 1 ("Binary Variables"): Key Concepts

- We discussed multivariate binary density estimation
  - Refined inference tasks when we have more than one random variable:
    - Joint probability, marginal probability, and conditional probability
  - Product of Bernoullis assumes variables are independent
    - Fast inference/learning but a strong assumption
- We discussed generative classifiers:
  - Build a model of the joint probability of features and labels
    - Compared to usual discriminative classifiers that model labels given features
  - Naïve Bayes assumes features are independent given label
- We discussed neural networks:
  - Model that learns the features and classifier simultaneously
  - Alternate between linear and non-linear transformations (universal approximator)
  - Training is a non-convex problem, but SGD often works better than expected:
    - For large-enough networks we often find global, and SGD seems to have implicit regularization

# End of Part 1 ("Binary Variables"): Key Concepts

- We discussed deep learning with multiple hidden layers
  - Biological motivations and efficient representation of some functions
  - Vanishing gradient problem and modern solutions:
    - ReLU, skip connections, ResNets
- We discussed automatic differentiation to generate gradient code
  - Code that generates gradient code for you (using chain rule)
- We discussed convolutional neural networks (CNNs):
  - Include convolution layers that measure image features
  - Include max pooling layers that highlight top features across space
  - Reduces number of parameters and gives some spatial invariance

# End of Part 1 ("Binary Variables"): Key Concepts

- We discussed autoencoders:
  - Networks where the output is the input
  - Encodes input into a bottleneck layer, then decodes back to input
  - Non-linear dimensionality reduction
  - Denoising autoencoders learn to enhance images
- We discussed multi-label classification:
  - Where each training examples can have 0-k correct labels
  - We discussed an encoding approach where the classes shares hidden layers
    - Reduces parameters and captures dependencies between labels
  - We discussed pre-training to learn new tasks with fewer labeled examples
- We discussed pixel labeling:
  - Fully-convolutional networks maintain spatial information at all layers
    - Requires upsampling to original image size
    - Can label images of different sizes