

of applications and uses of Computer Imaging. For MNIST digit recognition, simple thresholding / defining explicit rules isn't enough for great results at distinguishing digits. We want to be able *learn* the rules from data!

- Images are always arrays of pixels (possible 3 channels, or 1 - grayscale). Input image $x = [x_1, x_2, \dots, x_n]^T$; x_i 's are the columns of the image.
- $\Theta^x = \Theta_1 \circ \Theta_2 \circ \Theta_3 \circ \dots \circ \Theta_n$; this has many parameters - 785 if we have 28 pixels in our image! ($28^2 \times 1$ for bias).

We can reduce number of input features in other ways - by capturing different features than raw pixel values, e.g. num white pixels, digit width, height, slant, thickness, length. Using all of them we can make a linear model: $x = [\text{length}, \text{thickness}, \text{height}]^T$; $\Theta^x = \Theta_0 + \Theta_1 x$.

Convolutional layer: applying filter K with weights to the input producing a feature map that's feed into the layers (with a bias term) and an activation function (for non-linearity). $Z = ReLU(Wx + b)$.

Convolutional layer: $(W \otimes x) = (W \otimes x)^T$, this has many parameters - 785 if we have 28 pixels in our image! ($28^2 \times 1$ for bias).

We can reduce number of input features in other ways - by capturing different features than raw pixel values, e.g. num white pixels, digit width, height, slant, thickness, length. Using all of them we can make a linear model: $x = [\text{length}, \text{thickness}, \text{height}]^T$; $\Theta^x = \Theta_0 + \Theta_1 x$.

Logistic regression: $P(x) = 1/(1 + e^{-(\Theta^T x)})$

Convolutional layer: applying filter K with weights to the input producing a feature map that's feed into the layers (with a bias term) and an activation function (for non-linearity). $Z = ReLU(Wx + b)$.

Convolutional layer: $(W \otimes x) = (W \otimes x)^T$, this has many parameters - 785 if we have 28 pixels in our image! ($28^2 \times 1$ for bias).

3 rows of $[1 \dots 1]$ is used for computing y gradient / edge detection, 3 columns of $[1 \dots 1]$ T is used for computing y gradient / edge detection.

Convolutional filters are translation equivariant! Key for images

Equivariant - changes the output same way it changes input.

Padding and Strides: For an Input Size (M) we have the following terms:

• Padding (P) - how many extra values are put around edges in the conv.

• Stride (S) - how much the kernel is moved across after each application

• Kernel Size (K). The below formula assumes square sized kernels.

• Output Size = $(M + 2P - K) / S + 1$

Pooling: Reduces the size of feature maps without doing learning. Allows for some local translation invariance. We define kernels specifically with stride = kernel width so there's no overlapping, e.g. **MaxPooling, AveragePooling**

AlexNet: $22 \times 22 \times 3$ input, then maxpools & convs, until we get $7 \times 7 \times 512$, then 2 fc layers to make 1000 outputs (as it's image classifier for 1000 categories)

Building Novel Detection Systems (Skin Cancer):

• Collect Data to try self supervised learning features OR consult experts to get idea of features to extract.

• Start from a model with initialised weights (a classifier) and train that - **Transfer Learning by finetuning**. Training just the classification head!

Semantic Segmentation

Segmented regions are also assigned a semantic meaning.

Different to segmentation based on 'pure' clustering of the image into coherent regions. For imaging, we care about Semantic Segmentation.

Useful for 1) quantitative analyses, e.g. measuring the volume of the ventricular cavity, 2) determining the precise location and extent of an organ or tissue type e.g. a tumour, for treatment such as radiation therapy.

3) creating 3D models used for simulation, e.g. generating a model of an abdominal aortic aneurysm for simulating stress/strain distributions.

Challenges of Segmentation: noise, partial volume effects, intensity inhom., anisotropic res, imaging artifacts, limited contrast, morphological variability

• **Partial Volume Effects** - b/c of coarse sampling the result im. shows partial volume effects at the boundaries; both colours (black and white) contribute to the intensity val or boundary pixels b/c large influence area of each pixel.

• **Intensity Inhomogeneities**: difference of an image with diff contrast (sun)

• **Anisotropic Resolution**: 3D images are often blurry (CT, MRI)

• **Shadows:** shadow of some axial slices of the left ventricle may have low intraslice resolution of 1.3mm and an inter-slice resolution of 8mm

• **Imaging Artifacts:** e.g. From metal teeth, disrupts the CT scan

• **Limited Contrast**: diff tissues have similar phys properties and thus similar intensities. Intensity based algo can fail or leak into adjacent tissue.

• **Morphological Variability**: hearts/livers vary in shape person to person!

• **Ground Truth:** The reference the method can be compared against. There's not really a ground truth segmentation (except physical phantom/simulation)

• **Gold Standard:** The best truth we can have - often manual segmentation by experts. But reg training, tedious, takes time, intra/inter observer variation.

Measuring Performance: We overlay the segmentation map (SM) & gold standard (GS). We then easily see regions of agreement and disagreement. SM & GS = TP, SM & GS = FP, SM & GS = FN, GS & TN.

We can then construct the **Confusion Matrix**: $\begin{bmatrix} \text{TP} & \text{FP} \\ \text{FN} & \text{TN} \end{bmatrix}$. FP = type 1 error, FN = type 2 error. From here we can compute performance metrics:

1. **Accuracy**: $(\text{TP} + \text{TN}) / (\text{P} + \text{N})$ **Precision**: $\text{tp} / (\text{tp} + \text{fp})$; TP = $(\text{TP} + \text{FP}) / \text{P}$

3. **Recall** (Sensitivity / Hit Rate / PTP): $\text{TP} / \text{P} = \text{TP} / (\text{TP} + \text{FN})$

4. **Specificity** (True Negative Rate): $\text{TN} / \text{N} = \text{TN} / (\text{TN} + \text{FP})$

5. **FI Score**: Hmean of Prec&Rec: $2(\text{PPV} \cdot \text{TPR}) / (\text{PPV} + \text{TPR}) = 2\text{TP} / (2\text{TP} + \text{FP} + \text{FN})$

6. **Jaccard Index**: $\text{A} \cap \text{B} / |\text{A} \cup \text{B}|$

7. **Dice Similarity Coefficient**: $DSC = 2(A \cap B) / (A + B)$. Dice can be worked out from Jaccard & DC & F1 are equal: $\text{DC} = \text{TP} / (\text{TP} + \text{FN})$; $\text{F1} = 2 \cdot \text{DC} \cdot \text{Rec} / (\text{DC} + \text{Rec})$; $|\text{A} \cap \text{B}| / |\text{A} \cup \text{B}| = 1 - \text{FP} / (\text{TP} + \text{FP}) = 1 - \text{FN} / (\text{TP} + \text{FN})$

8. **Volumetric Dice Coefficient**: $|\text{A} \cap \text{B}| / |\text{A} \cup \text{B}| = 1 - \text{FP} / (\text{TP} + \text{FP}) = 1 - \text{FN} / (\text{TP} + \text{FN})$

9. **Surface Distance Metric**: $SD = \sqrt{(\sum_{i=1}^n d_i^2) / n}$ where d_i is the distance between each pixel's value is its distance to the boundary.

• Swap the contours of the two images and walk along the contour edge, for Hausdorff take the maximum value on the edge. For ASD just mean

• Distance maps have the nice side effect of letting you easily generate the contour formed if the boundary changes efficiently

pitfalls in Segmentation Evaluation:

1. **Effect of Structure Size:** For the Dice Coefficient, a 1 pixel difference has a much bigger impact on a small than big image. This is problematic if our metric aggregates over many structures of different sizes. Dice is quite invariant to different structure shapes, but very sensitive to pixel offsets

2. **Effect of Annotation Noise:** The ground truth standard might have a label error - if there's a pixel marked somewhere far away, massively impacts Hausdorff

3. **Effect of 'empty' labelmaps:** Might have nothing labelled - then TP, FN = 0, and we get $NaN(0/0)$ for a lot of our metrics! Needs explicit handling.

4. **Effect of resolution:** Affects calculations like DSC, Hausdorff...

Segmentation Methods:

1. **Simple Thresholding:** Select a thresh to split intensity range into 2 classes. Can use multiple thresholds for multi-class, or upsample for one class that's within the interval and the other is outside of it. Simplest, fast.

2. **Region Growing:** Start from (user selected) seed point(s), and grow a region according to an intensity threshold. Adv: relatively fast, yields connected region (from a seed point). Disadv: regions must be homogeneous, leakages and 'rough' boundaries likely, needs user input for SPS.

3. **Atlas Based Segmentation:** Generate template of an example segm & scene those (in an "Atlas") used for future segm. Has geo info abt points, curves or surfaces, or label info about voxels (anatomical regions or function).

Usually constructed from examples: single subjects, populations of subjects, e.g. averaging to produce probabilistic atlases

a. We segment via **label propagation**: check atlas (alt. has expert segms) & align them to test images via **registration**. If you can establish spatial correspondence b/w the imgs you can prop info onto the test images. To aggregate info: label fusion from multi atlases (from multiple examples) propagation. i.e. each atlas "votes" for what it thinks the segmentation is, label fusion just picks the overall label via majority voting or other metrics.

Adv: Robust, acc. Makes plausible segm. Automatic. Disadv: expensive. Doesn't deal well with abnormalities. Doesn't work well for tumour segmentation.

4. **Segmentation via Dense Classification:** We can use trained CNNs (as we've seen) to perform dense (pixel by pixel) classification on the image. **Classical Approach:** Patchify img, do conv (e.g. dim 5 kernels); w/FC layer & 2 neurons at end. We patchify into $17 \times 17 \times 3$; 4.5x5 kern convs results in 1x1 outputs. **Problem:** Overlapping patches; the FC layer at the end is very inefficient as we must combine all patch results! **Better Approach:** LeNet. **Receipt:** Feed (F) a decision - set of pixels that influence the size they were designed for if we use FC layers. We can fix that by replacing FC layers with a Fully Convolutional Layer which have kernel size exactly equal to input size. Feeding a 32×32 into LeNet ends up with a 1x1 classification answer. Feeding 64×64 ends in a 9x9 classification map, which has both classification and localization info (where high activation is = where object is!)

```
class LeNet(nn.Module):
    def __init__(self):
        super(LeNet, self).__init__()
        self.conv1 = nn.Conv2d(1, 6, kernel_size=5)
        self.conv2 = nn.Conv2d(6, 16, kernel_size=5)
        self.fc1 = nn.Linear(16 * 5 * 5, 120)
        self.fc2 = nn.Linear(120, 84)
        self.fc3 = nn.Linear(84, 10)

    def forward(self, x):
        x = F.max_pool2d(F.relu(self.conv1(x)), 2)
        x = F.max_pool2d(F.relu(self.conv2(x)), 2)
        x = x.view(x.size(0), -1)
        x = F.relu(self.fc1(x))
        x = F.relu(self.fc2(x))
        x = self.fc3(x)
        return F.log_softmax(x, dim=1)
```

```
class FCLENet(nn.Module):
    def __init__(self):
        super(FCLENet, self).__init__()
        self.conv1 = nn.Conv2d(1, 6, kernel_size=5)
        self.conv2 = nn.Conv2d(6, 16, kernel_size=5)
        self.conv3 = nn.Conv2d(16, 64, kernel_size=5)
        self.fc1 = nn.Linear(64 * 5 * 5, 120)
        self.fc2 = nn.Linear(120, 84)
        self.fc3 = nn.Linear(84, 10)

    def forward(self, x):
        x = F.max_pool2d(F.relu(self.conv1(x)), 2)
        x = F.max_pool2d(F.relu(self.conv2(x)), 2)
        x = F.relu(self.conv3(x))
        x = self.conv5(x)
        return F.log_softmax(x, dim=1)
```

Upsampling: We can AverageUnpooling or MaxUnpooling - average unpooling has each pixel duplicated by unpool kernel size, top-left MaxUnpool has the pixel duplicated once into top-left of the kernel and the rest zeroed. **Unpool with spatial information:** you record which pixel indices were the maxpool donors. Then after some intermediate calcs, we put the intermediate results back into those specific indices. A form of **Encoder-Decoder** net.

Encoder-Decoder Networks: Have high res input, convs and downsampling, then we resample at the end with upscaling. VGG-16 instead uses maxpooling to downsample and unpooling to upsample. UNet has 'skip layers' that copy and connect the features from the downsampling stage to the upscaling stage before any convolutional, producing better mappings. **Joint Embedding Prediction:** combines the two prior approaches.

• **MAE** avoid some drawbacks of CL, but training a full res recon model is costly

• Not necessary to fully learn a perfect dec. model as unused after training

• The idea behind I-Jepa: combining strength of CL and MAE approaches

• We mask most of the image, leaving one region visible to the net. We encode this producing an embedding via the **context decoder**.

• We do the same for many parts of the image, putting it through the **target encoder**. You obtain one embedding for each patch

• Then we take the context of the predictor, and try to work out which embedding passed through the target encoder are part of that context

• Loss is MSE b/w the emb space: $\sum_i D(s_i, s_j)$ = $\sum_i D(s_i, s_j) / n$

• Loss is $\text{L1} / 2 \cdot \text{dist}(h, p)$ patch $\sum_i D(s_i, s_j) / n$

• Advantages: 1) less computation 2) more robust to noise

• Disadvantages: 1) need to learn patch embeddings 2) need to learn context encoder

• **Downsample**: Dilate by 2 means each pixel in the receptive field of the convolution has a gap of 0 to the next either horizontally or vertically.

• **Atrous Spatial Pyramid Pooling:** Stacks dilated convolutions of multiple sizes within a single layer to gain lots and lots of context at once!

• **Padding** obviously introduces artifacts.

Multi-Scale Processing: To help the network see more context, we can add upsample layers which process down-sampled images of a bigger patch, and concat them to the pathway which processed the normal sized patch.

Vision Transformers: uses Multi-Head Attention for images by patchifying. Each gets a separate embedding. Uses the patch for a linear projection before adding on the position embedding before encoding & classifying w/ MLP head. We also pass in a special token - the class token to the encoder that controls what summarises patch info. Learnable. Enc. involves: input: A > norm: B > multi-head: C > concat A&C: D > norm: E > MLP: F > concat D&F: Output

Transformers: patchify > linearly proj patches > transformer enc. > dec. > output

Segment Anything model produces segment masks for all detectable objects.

Self-Supervised Learning

Self-Supervised Learning: Goal - building pattern recognition into networks without supplying the labels. Why? - supervised classifiers require lots of data, and in domains like healthcare the labelling process is very expensive. We do however have lots of labelled images. We want to learn w/o labels!

• Supervised Learners from their encoders outputted image representation. h. We want to learn the repr h, that summarises the main information in the image without needing any external labels. We need to create "proxy" tasks that only require the images to create the supervision signal.

Contrastive Learning: Teaching the network to recog. meaningful im. pairs

• Main idea: a good repr. will put similar ims. close together in emb. space

• An artificially perturbed version of the same image, should have repr. closest to the original image compared to all other images. This is the key principle of Contrastive Learning: creating multiple versions of the same image and then having the network to choose the correct version.

• SIMCLR - implements this. Encoder produces h for downstream tasks. MLP produces a smaller repr z to compute the contrastive loss.

• Pipeline: Take image, augment it to produce copy. Encode both, MLP both, produce two zs - compute contrastive loss on zs - the two similar images need close z, but the further ones should have distant z. This is our training objective

• **Augmentations** must be sophisticated enough. They need to reflect what information the model should disregard and what it should focus on. Needs to be hard enough, otherwise trivial information is learned

• Rotate, crop, flip, colour distort, cutout (mask), blur, noise, sobel

• If the augmentation isn't hard the model will learn colour, bg, etc

• Keep in mind we don't want to learn b/w the two images that's it's a cat but that it is THIS cat. We push other cats further from our pair

• **NT-Xent loss** (normalised temperature contrastive loss):

• Similarity $s(u, v) = (u^\top v) / \|u\| \|v\|$

• For each positive pair (i, j) we have loss:

$l_{ij} = -\log(s(z_i, z_j) / \sum_k l_{ik} l_{jk}) = -\log(s(z_i, z_j) / \sum_k s(z_i, z_k))$

• To min loss we need to max the above -> max numerator, min denominator

• Maximizing numerator: pulling the positive pairs closer

• Minimizing denominator: push the negative pairs apart

• The final loss is computed across all positive pairs, both (i, j) and (j, i), in a minibatch, by averaging all l_{ij}

• **Applications:**

• Satellite Imaging (to reconcile / compare / combineimages of the same thing from different locations - satellite is always moving)

• Point Correspondences - Panoramic Stitching.

• Temperature τ controls how much to penalise hard negatives (negative pairs wrongly mapped close to each other). A low temperature penalise them more.

• There's other losses too: triplet loss:

$$\mathcal{L}_{\text{triplet}}(x, x^+, x^-) = \sum_{x \in \mathcal{X}} \max(0, \|f(x) - f(x^+)\|_2^2 - \|f(x) - f(x^-)\|_2^2 + \epsilon)$$

• x is one image, x^+ is the correpos pos and x^- a randomly selected neg image. ϵ is the margin parameter controlling how far the negatives should be compared to the positive pairs (like temperature).

Evaluating Self-Supervised Models:

• Check perf on d-stream classification task or use (smaller) labelled dataset.

• **Lined Plotting**: pretrain the encoder with self-supervised learning e.g. SimCLR & then finetune on classifier, fit τ to the weights of the encoder. Train classifier head on top of the frozen encoder with cross-entropy loss.

• For finetuning, we don't freeze the weights.

- Generative Models are in one of two ways:
 - Exogenous Prior: $p_{\text{A}}(z)$ - the latent space not touched by any of the observed variables directly
 - Conditional Prior: $p_{\text{A}}(z|x)$ - the latent space IS affected by observed variables.
- In any case we're stacking lots of residual blocks and have a very large latent space in order to try learn what we need
- Latent Model: The conditional prior (2) induces a latent mediator since z_A is no longer exogenous
- It still has a Markovian interpretation, so we can do causal mediation analysis
- CMA: the study of $p(\text{U}) = p(\text{U}_x) \prod_{i=1}^n p(\text{U}_{p_i}) \prod_{i=1}^n p(\text{U}_k)$ how a transformation effect is mediated by another variable, to help explain why or how an individual may respond to certain stimulus
- This enables estimation of Direct (DE), Indirect (IE) and Total (TE) causal effects: $\text{DE}_x(p_{\text{A}}) = \mathbb{E}[g(p_{\text{A}}, z_{1:L})] - g(p_{\text{A}}, z_{1:L-1})$
- As an example in Morpho MN: $\text{IE}_x(z_{1:L}) = \mathbb{E}[g(p_{\text{A}}, z_{1:L})] - g(p_{\text{A}}, z_{1:L-1})$ -IST, you get to see that $\text{TE}_x(p_{\text{A}}, z_{1:L}) = \mathbb{E}[g(p_{\text{A}}, z_{1:L})] - g(p_{\text{A}}, z_{1:L-1})$ observed digit thickness causes pixel intensity. Changing thickness changes the generated thickness but not vice versa!

Evaluating Counterfactuals:

- Soundness Theorem: the properties of composition, effectiveness, and reversibility are necessary in all causal models
- The completeness theorem states that these properties are sufficient
- We measure counterfactual soundness using these axiomatic properties:
 - Composition: If you intervene on a variable X by setting it to the value it would have had anyway, no other variables' values are changed.
 - Then there exists a null-intervention
 - Effectiveness: Intervening on a variable to have a specific value will indeed cause the variable to take on that value.
 - For all of these axioms, we say the model is debiased if it does the 3 Reversibility: If two variables, X and Y , influence each other in a way that creates a feedback loop, they must settle to a single, consistent pair of values rather than multiple possible solutions.
 - These properties have a trade off: While composition \rightarrow reversibility, you can't really satisfy all 3 at the same time optimally. Need to pick one which is most important and maximize it unbiasedly.

Overall, it's crucial to account for the data-generating process and causality in our modelling to avoid biased predictions.

Deep SCMs can generate plausible high-fidelity image counterfactuals
Latent mediator models enable estimation of direct, indirect and total causal effects for high-dimensional variables
Sadv: measuring/counterfactual effects needs separately trained classifiers.

Inverse Problems

- Inverse Problems:** Given a measurement $y = Ax + n$, we want to estimate x .
 • Note the noise term n . Applications: inpainting/de-blurring/recon.
 • The matrix A is known as the forward model
- The Classical Solution for Deblurring:**
- Formulating it as a Least Squares Problem: $\arg \min D(Ax, y)$: we have that $D(Ax, y) = 1/2 \|Ax - y\|^2$; the solution is $x^* = (A^T A)^{-1} A^T y$ (air know A)
 - This is an ill-posed problem though. Slightly perturbing the input image gives us a wildly different solution. There's multiple "correct" solutions...
 - To get a more correct solution we add a **Regularization** term to the above: $R(x, y) = R_1 \|x\|_2^2 + R_2 \|y\|_2^2$; i.e. the solution is $\arg \min D(Ax, y) + R(x)$
 - Least SQ with Tikhonov regularisation: $x = (A^T A + R)^{-1} A^T y$ - this is better conditioned and helps suppress noise.
 - This gives the general approach for images: solve the optimization prob: $\arg \min D(Ax, y) + R(x)$, $D(Ax, y)$ - Data Consistency term, $R(x)$ - Regularization term.
 - Common Regularization: $R(x) = \|\nabla x\|_2^2$ - Tikhonov Regularisation (first derivative of the image should be smooth).
 - $R(x) = \|\nabla x\|_2^2 = \sum_{i,j} ((\nabla x)_{i,j}^2)$ - Total Variation Regularization (don't have too much intensity fluctuations in image x - dimage)
 - $R(x) = \|\nabla x\|_1$ - Wavelet Regularisation (∇x = gradient of image x) (don't have too much components in the wavelet domain space)
 - Instead of choosing, we'd like to learn a good R from the training data: **Building a Forward/Inverse Model**:
 - Data: We either have lots of images and their perturbed version OR we have a forward model and we generate the perturbed versions ourselves
 - We have a specific inverse depending on what our forward is **Solving the problem of building the Inverse Model**:

1. **Model Agnostic Approach:** we can train a Neural Network that maps y to x ; just optimize loss based on how close the output is to the groundtruth.
- **Partly Model Agnostic:** We improve by doing preprocessing before feeding into inverse NN model: we know our foward (e.g. downsample), so we can upsample the image before feeding into NN, improves results

Decoupled Approach (first learn then reconstruct):

- a. Deep Proximal Gradient set $\tilde{x}^{(0)}$ and stepsize η
- $$\begin{aligned} & \tilde{x}^{(k+1)} = \tilde{x}^{(k)} + \eta A^T(y - A\tilde{x}^{(k)}) && \text{gradient descent (DC)} \\ & \tilde{x}^{(k+1)} = \arg \min_{\tilde{x}} \|\tilde{x}\|_2^2 + \eta \|A\tilde{x}\|_2^2 && \text{denoising} \end{aligned}$$
- b. We could also use a GAN: $H(x) = 0$ if x is on image manifold, \sim otherwise
- The goal is to find images similar to the natural images we expect
 - GANs have 2 components: 1) a generator G (forward model takes a noise vector) and converts to an image (encoder-like). 2) we also have a discriminator D which we supply the generated image and real image. The job of the discriminator is to see whether the generation looks real or not.
 - We train D to max the probability of assigning correct label. Train G to "fool" D . These correspond to the following two terms in Value Function: $\min_D \max_G V(G, D) = E_{x \sim p_{\text{data}}(x)} [\log D(x)] + E_{x \sim p_g(x)} [\log (1 - D(G(x)))]$

• This discriminator is useful as a regularization function!

• But then you can use the generator instead of gradient descent as it produces "solutions" to the optimization problem.

- Learn generator G that outputs $x \in \mathcal{D}$ given $x \in \mathcal{R}^d$ for $d < d'$
- $R(x) = 0$ for $x \in \mathcal{R}^d$, \sim otherwise
- We then choose $x \in \mathcal{R}^d$ that best fits the data:

$$x = \arg \min_{x \in \mathcal{R}^d} \text{gen}(G(x)) - \text{real}^2$$

3. Unrolled Optimization: assume that $R(x)$ is differentiable. Then:

$$\begin{aligned} & \text{set } \tilde{x}^{(1)} \text{ and stepsize } \eta \\ & \text{for } k = 1, 2, \dots \end{aligned}$$

$$\begin{aligned} & \tilde{x}^{(k+1)} = \tilde{x}^{(k)} + \eta A^T(y - A\tilde{x}^{(k)}) - \eta \nabla_R(\tilde{x}^{(k)}) \text{ replace with learned NN} \\ & (\tilde{x}^{(k+1)} - \tilde{x}^{(k)}) = (\tilde{x}^{(k)} - \tilde{x}^{(k)}) + \eta A^T(y - A\tilde{x}^{(k)}) \end{aligned}$$

Unrolled optimization framework, trained end-to-end to what we want to compute H derivative, and we expand out the grad desc term (do everything in one line). We often however don't know how to differentiate R so we use an NN in practice. We need an NN for each step - it becomes costly unless we ensure there isn't too many iter steps

We can apply this inverse problem solving framework to many applications: Deep Learning for Image Super-Resolution:

- Upsample low res im, to high res. Forward model: im degradation followed by downampling. Inverse model: interpolation based models
- **Post-Upsampling Super Resolution (Model Agnostic):**
 - Directly upsample LR to SR. Requires many learnable upsampling layers
 - Fast, low mem required. But net has to learn entire upsampling pipeline (no prior knowledge) & only performs one specific up-sampling factor
- **Pre-Upsampling Super Resolution:**
 - Use traditional upsampling algorithm (e.g. LERP) to obtain SR images. Then refine upsampling by a deep neural network (usually a CNN)
 - CNN only needs to correct smaller details. Can be applied to a range of upscaling ratios and image sizes. But operates on SR image, thus reducing more computation and memory

Progressive Upsampling Super-Resolution:

- Multi-stage: Use a cascade of CNNs to recorr higher and higher res. At each stage, the images are upsampled to higher res & refined by CNNs
- Decomposes complex task into simple tasks, reasonable efficiency. But it's sometimes difficult to train very deep models
- **Iterative Up-and-Down Sampling Super-Resolution (Sota if using CNNs):**
 - Alternate between upsampling and downsampling (back-projection) operations. We have up-up and down-down skip conns to later layers (for error prop), & we have mutually comp up & down-sampling stages
 - Superior perf as it allows error feedback. Easier training of deep nets
 - Loss: Pixel-wise loss function (either L_1 or L_2): $L = 1/M \sum_{i=1}^M \|x_i - \hat{x}_i\|^p$
 - Huber-Loss: $L = 1/M \sum_{i=1}^M \min\{\alpha |x_i - \hat{x}_i|, \alpha^2/2\}$ where $\alpha = 0.5$ for $|x_i| < 1$ else $|x_i| - 0.5$ this is a mix of L_1 and L_2 . L_2 is mostly used in practice.

Homomorphic Encryption:

- The entire process is encrypted: encrypted model is trained on encrypted data, the output is only decrypted as the final result is revealed. We assume we can do basic arithmetic on encrypted vals: $x[y] = [x][y] = [x][y] - [x][0]$; no intermed. decryp - useful if we don't trust the server & are scared of middlemen. $* \rightarrow *$ easily homomorphically encrptable, but trig and log are not; & time consuming

Secure Multi-Party Computation:

- individual contributions are 'shared' and only the individual shares are shared with the participants. Only when shards are combined is the final result revealed. So we have an "Add" characteristic, and for M parties, we have M shards and each shard contains both parts of (x, y) (e.g. $x[0], x[1], \dots, x[M-1]$, $y[0], y[1], \dots, y[M-1]$)

Total Variation:

- $L = 1/M \sum_{i=1}^M \|x[i] - x[i-1]\|^2$ is the abss value of gradient of the image is low, e.g. image is piecewise constant. For 2D im: $L = \frac{1}{n_1 n_2} \sum_{i,j} \sqrt{|x_{i,j} - x_{i,j-1}|^2 + |x_{i,j} - x_{i-1,j}|^2}$

GAN Discriminator Loss:

We can use the discriminator output of GANs (as explained earlier) $L = \log D(G(x))$

SoTA: Use $\text{GAN Loss} + L_2$ Loss + Perceptual-Loss

Methods of Upsampling 1) Nearest Neighbour 2) bilinear (add intermediate pixels LERP (the original pixels) 3) Learnable weights via transpose convs Parameterizing Images via Generative Networks:

- Consider a generator net ψ , with a random fixed noise input vector z_0 and pass it through an encoder-decoder net. We get an output image.
- The network params w can be thought as a parameterization of images
- We want to fit the net to ONE single image. Given an image x , its params w are recovered by solving the following optim problem: $\min \|x - \psi(z_0; w)\|^2$. The net now encodes x , and we can rebuild the by giving it z at each time

Deep Image Priors:

Optimizing this net is very fast - for most generative nets naturally looking images is easier/faster than fitting others

Natural looking images fit models easily. Unnatural don't fit! We could use that to evaluate our image outputs.

Concrete Attack Example: Gradient-based model inversion:

- If you're a MITM then you send updates sent by client to server. You can use an algorithm to reconstruct the input from the gradient updates and discover the original data!

The adversary randomly generates an image-gradient pair

- The adversary captures the gradient update made by the victim
- Using a suitable cost function (often cosine similarity), the adversary minimises the difference between the captured and the generated update by perturbing the image they control

The algorithm is repeated until the final iteration is reached

$$\arg \min_{x' \in \mathbb{R}^{n \times n}} \left\{ 1 - \frac{\|\nabla_Q(L(x, y), \nabla_Q(L(x', y), y)\|_2}{\|\nabla_Q(L(x, y), y\|_2 \|\nabla_Q(L(x', y), y)\|_2} + \alpha TV(x') \right\}$$

where x' is recon target, x is ground truth, y label, VL gradient wrt weights, ϵ inner prod. α is hyperparam scaling the variation penalty over the image.

Concrete Mitigation: Differentially private stochastic gradient descent (DP-SGD):

- Compute gradients for each individual sample (they represent independent clients)
- Clip the calculated gradients to obtain a known sensitivity
- Add the noise scaled by the sensitivity from step 2
- Perform the gradient descent step

Interpretability and Explainability:

We want interpretability to understand why a model made a decision. This does also help debug

Interpretability maps/T-SNE help us see why the model made a decision.

Helps debug classifications (or mortor learning) & detect bias

Common Misunderstandings - some models are actually more interpretable. (0) Interpretablity \rightarrow fairness/just/accuracy/interpretable \rightarrow interpretability. (1) Interpretability doesn't always matter \rightarrow can also lead to gaming the sys

We care about Interpretability in our final model - not before/during train

Methods of trying to find interpretability:

Ablation Test: How important is a data point/feature? Train without some data or features and observe/study the impact. Difficult and expensive.

Fit Functions: (use first derivatives) - sensitivity analysis or saliency maps

• Looks at the distrib, takes a sample, then we look at the gradient of that wrt features. i.e. if output says class 1, perturb features & check again

• Not necessarily easy for humans to interpret - gradient might be jumpy

i.e. take first deriv and look at image - might give a clear good edge map

Direct visualization of filters - Easy to implement; Limited practical value

• First layers are easy to interpret (mostly low-level features)

• Higher layers more difficult to interpret (non-interpretable features)

Occlusions: Mask out region in the input image and observe net output

• If masked out region causes a significant drop in confidence, the masked-out region is important

Salinity Maps: three similar approaches which do backprop differently

• DeconvNet: Given a trained network and an image, choose activation at one layer (set all others to zero), and try invert network

• Backprop looks the same as depth in net with fwd pass then invert.

• Guided Backprop: Improve results by "guiding" the backprop start.

Guided Backprop: Improve results by "guiding" the backprop start.

• Positive gradients = features the neuron is interested in

• Negative gradients = features the neuron is not interested in

DeepDream / Inceptionism: Attempt to understand inner workings of net

• Optimize with respect to image

• Arbitrary image or noise as input. Instead of adjusting network parameters, tweak image towards high "X" where "X" can be - Neuron/Activation map/Layer or Logits/Class probability. Search for images that are "interesting". Different layers enhance different features.

- For Federated SGD, we have data distributed to K different clients, k is the set of indices of the dataset of client k , $n_k = |\mathcal{D}_k|$. The loss f_k is then just summing over all clients and weighting their losses:
 - $L(k) = \sum_{i \in \mathcal{D}_k} (x_i - \hat{x}_i)^2$ (at each iteration we transmit a gradient - this could be huge (70 billion params = gradient / 70 billion entries)
 - $L(k)$ is just the loss of a specific client - that's what's sent to the server). With full assumption, loss of 1 client = expected loss of whole dataset.
- For Federated Averaging: The clients themselves do parameter updates - does this for N iterations before transmitting updated grad to server. We randomly choose some clients each round to update and avg.
- If the data is not IID, our gradient descent is not very good. If unbalanced or few training samples, this happens too. Communication is expensive, and no federal security/privacy provided; prone to adversarial influence

Overcoming Existing Privacy Preserving Methods:

- K-anonymity: if a participant's records are indistinguishable from at least k-1 other records which are at least k-1 other records which are indistinguishable.

• **Linkability:** the record should be unique to the person

• **Local Differential Privacy:** the noise term (which is proportional to the size of the record) is added to the record to make it look similar as possible by scaling, shearing and rotating them.

• **Fast Gradient Sign Method:**

- Perform gradient descent in order to maximize the loss (the goal of adversarial attack).
- Consider the input image x to be a trainable param and compute the gradient with respect to the input image to create a perturbation.
- $n = \epsilon$, sign($\nabla_L(x, 0, y)$)
- An adversarial example can be created as: $x' = x + \epsilon \cdot \text{sign}(\nabla_L(x, 0, y))$

• Defence: Generate adversarial examples, add the generated adversarial examples to the training set, retrain model using training set. There is however almost infinite adversarial examples...

Tutorials

Tutorial 0: Feature Learning - Handwriting the features can be hard - we might not know what's relevant. Instead we could try **learn** the relevant features.

• Logistic Regression has the advantage of letting us visualise our learned features (i.e. for MNIST what the model actually thinks the digits look like)

• It's done by taking the 784 coefficients and putting back into the image

• It hasn't actually learnt what the two digits look like, it's learnt the separation of the two (you can imagine a divisor in a hyperplane)

• A generative model understands the digits (and hence can make them)

Tutorial 1: Just a normal simple ML setup like the CW. Just use AlexNet. AlexNet still performs well on out of dist data (input images shift by 2 pixels)

Tutorial 2: Segmentation: Literally just implement U-Net.

Tutorial 3: Self Supervised Learning: Just implementing some stuff from the lecture like infoNCE.

Tutorial 4: Spatial Transformers: just implementing some spatial transformer architecture as in the lectures.

Tutorial 5: Looks like T-SNE - similar to CW

Tutorial 6: this tutorial actually looks important / different enough

Exam Questions & Formulae

Receptive Field = (previous receptive field) + (product of strides before curr layer) * (curr layer kernel size - 1) * (curr layer dilation)

Parameters: $com = (Kh \times Kw \times Channels \times Bias-1)$

$$H_{out} = H_{in} + 2 \cdot padding - dilation \cdot (kernel_size - 1) + 1$$

1) Param count of simple logistic regression model have where raw pixel values are used as inputs? Number of channels * raw pixels * 1*bias

2) Why do we split data into training, validation and test sets? train-test

3) Why do we use vgg16, vgg19, vgg16n, vgg19n and nvidia? vgg16, vgg19, vgg16n, vgg19n

4) Calculate the SSD, SAD, CC, and NMI for the following pair of images. Write down any intermediate values that you may calculate in the process.

$$im_1 = [1, 2, 2, 1, 2, 1, 3, 5, 4, 1], im_2 = [1, 2, 1, 1, 0, 1, 2, 1, 2, 1]$$

$$SSD = 1/V \sum (im_1 - im_2)^2 = 1/(2-1)^2 + 1/(2-1)^2 + ... + 1/(2-1)^2 = 113/9$$

$$SAD = 1/V \sum |im_1 - im_2| = 25/9$$

5) Design a simple enc-dec network with five conv layers and five convT layers that maps an img of size $16 \times 16 \times 64$ to an output of size $32 \times 32 \times 10$. enc-dec between the enc and dec of size $32 \times 32 \times 10$.

6) What is the size of the output for regularisation in a neural network? (check all that apply) periodicity in the weights, Dropout, Stochastic Gradient Descent, Batch Normalization, Adam Optimizer

7) A particular version of YOLO uses 11 x 11 cells for object detection across images. Furthermore, it uses four object bounding box predictions per cell and it can detect 10 classes of objects. What is the size of the output that the CNN has to predict? YOLO per bounding box predicts confidence(x,y,w,h) = 5 values as well as 10 conditional class prob \rightarrow 4x10-30 values per cell.

8) How to modify this particular version of YOLO to be able to detect 1000 object classes? (4x51000)*121

9) Briefly explain what steps you could take to avoid the computational difficulties that may arise from increasing the number of classes to 1000, separate obj detection from classification using a generic object detector (e.g. RCNN) and then a separate object classifier.

10) What would be a suitable dissimilarity measure that could be considered for intensity-based registration for each of the following three pairs of images? The CNN has to **strongly diagonal dots but with an isolated cluster** The histogram is obtained after registering a pre- and post-contrast scan of the same modality. The linear correlation of the majority of intensity pairs suggests mono-modality. The isolated cluster corresponds to a lesion that appears bright in the post-contrast image and dark in the pre-contrast one.

11) From 2020 exam: Design a simple spatial transformer network with two convolutional, two max pooling and two fully connected layers for estimating the params of a 2D rigid transformation for input of size 16×64 .

c1 = Conv2d(in=1, out=2, kernel_size=3, stride=1, pad=0)

c2 = Conv2d(in=2, out=4, kernel_size=3, stride=1, pad=0)

c3 = Conv2d(in=4, out=8, kernel_size=3, stride=1, pad=0)

c4 = Conv2d(in=8, out=16, kernel_size=3, stride=1, pad=0)

c5 = Conv2d(in=16, out=32, kernel_size=3, stride=1, pad=0)

c6 = ConvTranspose2d(in=32, out=10, kernel_size=3, stride=2, pad=0)

c7 = ConvTranspose2d(in=10, out=20, kernel_size=3, stride=2, pad=0)

c8 = ConvTranspose2d(in=20, out=40, kernel_size=3, stride=2, pad=0)

c9 = ConvTranspose2d(in=40, out=80, kernel_size=3, stride=2, pad=0)

c10 = ConvTranspose2d(in=80, out=160, kernel_size=3, stride=2, pad=0)

c11 = ConvTranspose2d(in=160, out=320, kernel_size=3, stride=2, pad=0)

c12 = ConvTranspose2d(in=320, out=640, kernel_size=3, stride=2, pad=0)

</