

# Deep Learning

- Linear/Logistic regression gradient descent
- Layers:
  - Linear
  - ReLU, Sigmoid
  - Softmax
  - $X$ -ent,  $X$ -ent w/ logits
  - Square loss
  - add, mul, sum, max
  - Conv, max-pool
    - padding, stride, dilation
  - batch norm & residual connections
- CNN:
  - general architecture
  - gradients problem & solutions
  - applications
- RNN:
  - simple RNN,  $\frac{\partial L}{\partial b}, \frac{\partial L}{\partial w}$
  - LSTM & thresholding
  - Applications
- Pretraining
- Optimisation:
  - SGD, why zig-zag, momentum & Nesterov
  - Newton, problem, Gauss Newton
  - Curvature Matrix approximations
- Attention & Memory:
  - Implicit; Jacobian
  - Explicit; glimpse, soft vs hard
- DL for NLP:
  - word embedding
  - LBL/NLM/k-gram/RNN-LM
  - seq-to-seq

# DL

## • Linear Regression:

$$L(\underline{w}, S) = \frac{1}{2} \sum (y_i - \underline{w}^T \underline{x}_i)^2$$

$$\nabla_{\underline{w}} L = \sum (y_i - \underline{w}^T \underline{x}_i) \underline{x}_i$$

## • Logistic Regression:

$$L(\underline{w}, S) = - \sum (y_i \log \sigma(\underline{w}^T \underline{x}_i) + (1-y_i) \log(1 - \sigma(\underline{w}^T \underline{x}_i)))$$

$$\nabla_{\underline{w}} L = - \sum (y_i (1 - \sigma(\underline{w}^T \underline{x}_i)) + (1-y_i) \sigma(\underline{w}^T \underline{x}_i)) \underline{x}_i$$

$$= \sum (\sigma(\underline{w}^T \underline{x}_i) - y_i) \underline{x}_i$$

Layer	$y$	$\frac{\partial L}{\partial x}$	$\frac{\partial L}{\partial \theta}$
Linear	$\underline{w}^T \underline{x} + b$	$\frac{\partial L}{\partial y} \underline{w}$	$\frac{\partial L}{\partial \underline{w}} = (\frac{\partial L}{\partial y})^T \underline{x}^T$
ReLU	$y_i = \max(x_i, 0)$	$\frac{\partial L}{\partial x_i} = (y_i > 0)$	$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial y}$
Sigmoid			
Softmax	$y_i = \frac{e^{x_i}}{\sum e^{x_i}}$	$\frac{\partial L}{\partial x_i} = \sum_j S_j \left( S_j = \frac{\partial L}{\partial y_i} \right) y_i$	
X-ent	$-\sum p_i \log x_i$	$\frac{\partial L}{\partial x_i} = -\frac{p_i}{x_i}$	
(w/ logits)	$-\sum p_i \log \left( \frac{e^{x_i}}{\sum e^{x_i}} \right)$	$\frac{\partial L}{\partial x_i} = y_i - p_i$	
Square Loss	$\  \underline{t} - \underline{x} \ ^2$	$\frac{\partial L}{\partial x_i} = -2(\underline{t} - \underline{x})$	
Add	$\underline{a} + \underline{b}$	$\frac{\partial L}{\partial a} = \frac{\partial L}{\partial y}$	
Mul	$\underline{a} \odot \underline{b}$	$\frac{\partial L}{\partial a} = \underline{b} \odot \frac{\partial L}{\partial y}$	
Sum	$\sum x_i$	$\frac{\partial L}{\partial x_i} = \frac{\partial L}{\partial y} \cdot 1^T$	
Max	$\max x_i$	$\frac{\partial L}{\partial x_i} = \frac{\partial L}{\partial y}$ if $x_i = \max \{\underline{x}\}$ else 0	

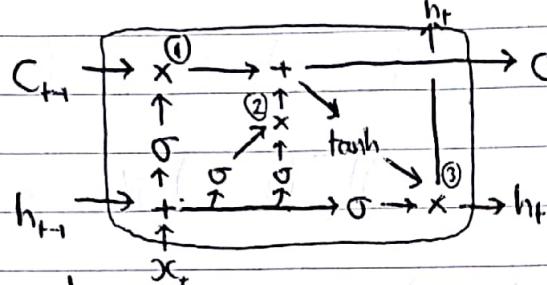
- CNN:
 
$$\text{output} \rightarrow M = \left[ \frac{N + 2p - R'}{S} \right] + 1, \quad R' = R + (R-1)(r-1)$$
  - SAME padding:  $p = (R-1)/2$
  - $[(\text{conv} \times m) \rightarrow \text{pool}] \times n \rightarrow \text{linear} \times k \rightarrow \text{softmax}$ 
    - Main workhorse:  $3 \times 3$  conv
  - Vanishing / exploding gradients:
    - add batch norm & residual connections
    - careful weight initialisation

• RNN:

$$h_0 \rightarrow h_1 \rightarrow h_2 \rightarrow \hat{p}_t = \text{softmax}(W_h h_t + U_b + b)$$

$$L(W, U, B) = -\sum_t \log p_t / (W_{i:t-1})$$

$$\nabla_B L = \sum_t (p_t - \text{label}_t) h_t^T, \quad \nabla_W L = \sum_t \sum_{k=1}^t \frac{\partial L_t}{\partial p_t} \frac{\partial p_t}{\partial h_k} \frac{\partial h_k}{\partial W}$$



• helps long term dependencies  
& prevents vanishing/exploding gradients  
(can also threshold gradients)

• Optimisation:

• GD:  $\Theta_{k+1} = \Theta_k - \alpha_k \nabla_{\Theta_k} L |_{\Theta_k}$

• implicitly approximates  $H$  by  $I$ ; overestimate curvature; zig-zags

• Momentum:  $\Theta_{k+1} = \Theta_k + \alpha_k V_{k+1}, \quad V_{k+1} = \eta_k V_k - \alpha_k \nabla_{\Theta_k} L |_{\Theta_k}$

• fixes zig-zagging; Nesterov better theoretical bound, similar in practice

• 2nd order:

• Newton:  $\Theta_{k+1} = \Theta_k - \alpha_k H^{-1} \nabla_{\Theta_k} L |_{\Theta_k}$

•  $H$  under-estimates curvature outside small radius

• can limit to  $r$ ; replace  $H$  with  $(H + \lambda I)$

• or use other matrix, e.g.  $G = \sum J_i^T H_i J_i, \quad J_i = \frac{\partial f(x, \Theta)}{\partial \Theta} \text{ wrt } \Theta$

•  $G$  works better in practice  
(it's always PSD)

$H$ : hessian of  $f(y, \hat{y})$   
wrt  $\hat{y}$

• Curvature matrix approximations:

•  $H: D \times D, \quad H^{-1} O(D^3) \quad (\& O(D^2) \text{ storage})$

• too expensive; use approximations to speed up

• diagonal & empirical Fisher matrix (RMSProp, Adam)

• low rank approx (L-BFGS)

• No asymptotic improvement from stochastic 2nd order, but  
improvement before asymptote  $\Rightarrow$  can still be extended to minibatch

• Attention & Memory (RNN):

• Implicitly given by Jacobian;  $J_{ij} = \frac{\partial y_i}{\partial x_j}$

• Explicitly done by attention model using network output + extra  
data to make a glimpse; then fed back into network as input

• Hard glimpse not differentiable; trains by RL. Soft-glimpse is; GD works  
good!

# COMP613 Advanced Topics in Machine Learning

## Deep Learning

### Supervised Learning

$$\underline{x} \rightarrow \boxed{f} \rightarrow y \quad \begin{cases} y \in \{-1, +1\}, \text{ classification} \\ y \in \mathbb{R}, \text{ regression} \end{cases}$$

#### • Generalisation:

- performance on unseen data
- empirical risk minimization: minimize training error
- Validation: hold out data to get unbiased ~~error~~ estimator
 
$$\text{validation error} = \text{training error} + \frac{\text{validation error} - \text{training error}}{\text{generalisation error}}$$

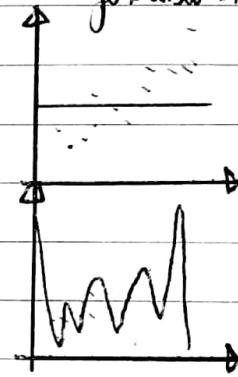
what we care about

biased!

eg cross validation

#### • Underfitting:

- error driven by approximation
- high bias, low variance



#### • Overfitting:

- error driven by generalisation
- low bias, high variance

#### • Early stopping:

- control degree of overfitting in gradient iterative optimisation methods

#### • Regularisation:

- impose a ~~scale~~ preference for simpler hypothesis
- weight decay (L1 or L2), dropouts..

loss function to minimize

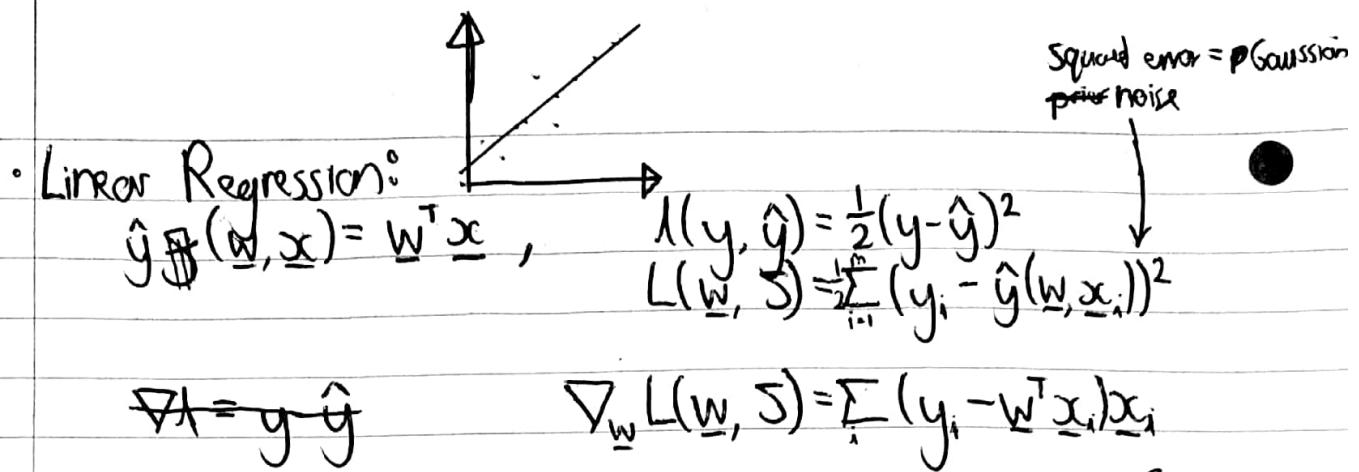
#### • Gradient descent:

$$\underline{w} \leftarrow \underline{w} - \eta \nabla L(\underline{w})$$

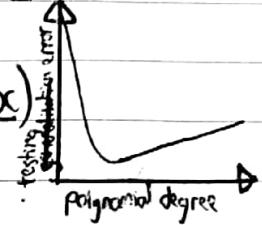
- useful when no closed form solution / dataset too large / high dimensional to practically solve  $S$
- often practical to use a mini-batch to approximate  $\nabla L(\underline{w})$ , where  $L(\underline{w}) = \sum_i f_i(\underline{w})$

$$\nabla L(\underline{w}) \approx \sum_{i \in S_j} \nabla f_i(\underline{w})$$

- minibatch size trades off computational cost & variance



- regularization →
- weight decay:  $L(w, S) = L(w, S) + \lambda \|w\|^2$  ← Gaussian prior on weights
  - Non-linear basis functions:
    - $\Phi(x) = (1, x, x^2)$ ,  $\hat{y}(w, x) = w^T \Phi(x)$
    - provides richer hypothesis space



- Logistic Regression:  $y \in \{0, 1\}$

$$z(w, x) = w^T x, \quad \hat{p}(z) = \frac{1}{1 + e^{-z}}$$

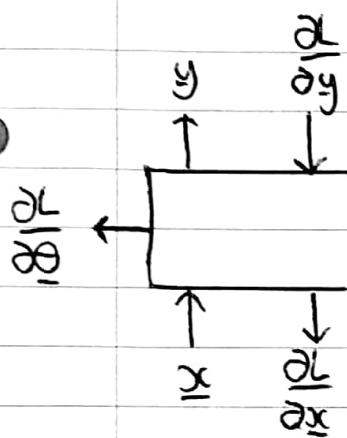
cross entropy loss →  $\ell(y, \hat{p}) = y \log \hat{p} + (1-y) \log (1-\hat{p})$

$$\begin{aligned} L(w, S) &= -\sum_i \ell(y_i, \hat{p}) \\ &= \sum_i \log(1 + e^{w^T x_i}) - y_i w^T x_i \end{aligned}$$

$$\begin{aligned} \nabla_w L(w, S) &= \sum_i \frac{x_i e^{w^T x_i}}{1 + e^{w^T x_i}} - y_i x_i \\ &= \sum_i \left( \frac{1}{1 + e^{-w^T x_i}} - y_i \right) x_i \end{aligned}$$

# Neural Networks: Foundations

- Composition of linear transforms + non-linear functions
  - learn by chain rule derivatives & SGD on loss function
  - each node in compute graph has:
    - forward-pass:  $\underline{x} \rightarrow y$  (computing output & loss)
    - backward-pass:  $\frac{\partial L}{\partial y} \rightarrow \frac{\partial L}{\partial \underline{x}}$  (backpropagation)
    - compute gradients:  $\underline{x}, \frac{\partial L}{\partial y} \rightarrow \frac{\partial L}{\partial \theta}$  (for param updates)
- Module Zoo:



- Linear layer:

$$y = \underline{W} \underline{x} + b$$

$$\frac{\partial L}{\partial \underline{x}} = \frac{\partial L}{\partial y} \underline{W}$$

$$\frac{\partial L}{\partial \underline{W}} = \left( \frac{\partial L}{\partial y} \right)^T \underline{x}^T$$

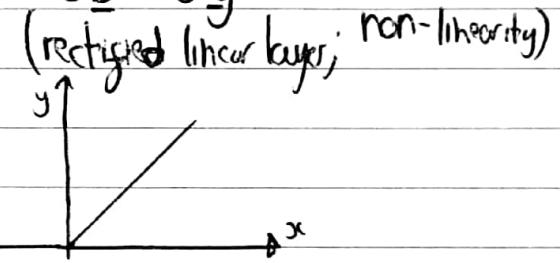
$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial y}$$

(applying linear transform)

- ReLU:

$$y_i = \max(x_i, 0)$$

$$\frac{\partial L}{\partial x_i} = (y_i > 0)$$



- Softmax:

$$y_i = \frac{e^{x_i}}{\sum e^{x_i}}$$

$$\frac{\partial L}{\partial x_i} = S - \sum_i S_i, \quad S_i = \frac{\partial L}{\partial y_i}, \quad y_i$$

- Cross-entropy loss:

$p$  = label probabilities

$$y = -\sum_i p_i \log x_i$$

$$\frac{\partial L}{\partial x_i} = -\frac{p_i}{x_i}$$

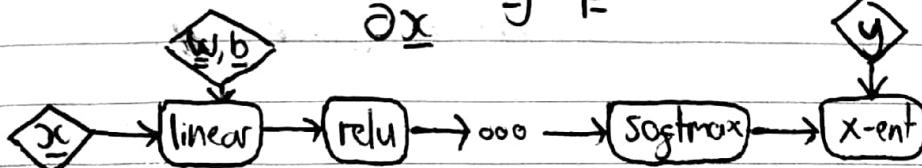
good to prevent overflow

- Cross-entropy loss w/ logits:

$$y = -\sum_i p_i \log \left( \frac{e^{x_i}}{\sum_j e^{x_j}} \right)$$

(softmax + Xent loss in 1)

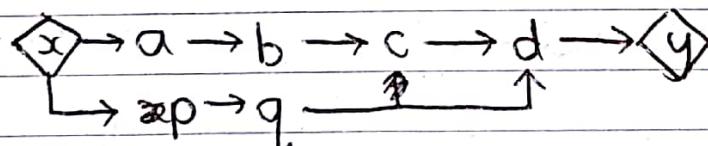
$$\frac{\partial L}{\partial \underline{x}} = y - p$$



element wise	• Add:	$y = \underline{a} + \underline{b}$	$\frac{\partial L}{\partial \underline{a}} = \frac{\partial L}{\partial \underline{y}}$
	• Mul:	$y = \underline{a} \odot \underline{b}$	$\frac{\partial L}{\partial \underline{a}} = \underline{b} \odot \frac{\partial L}{\partial \underline{y}}$
group wise	• Sum:	$y = \sum x_i$	$\frac{\partial L}{\partial \underline{x}} = \frac{\partial L}{\partial \underline{y}} \mathbf{1}^T$
	• Max:	$y = \max(x_i)$	$\frac{\partial L}{\partial x_i} = \frac{\partial L}{\partial y_i}$ if $x_i$ max, else 0
	• Leaky - ReLU:		
	• Squared loss:	$y = \ \underline{t} - \underline{x}\ ^2$	$\frac{\partial L}{\partial \underline{x}} = -2(\underline{t} - \underline{x})^T$

- Optimisation:
    - Forward pass to find loss
    - Backward pass to find gradients w.r.t. inputs
    - Compute gradients w.r.t. parameters & update by gradient descent

( backprop thru  
forward execution pass )



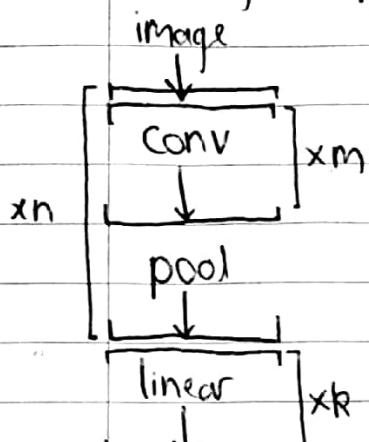
$$y = d(g(c(b(a(x))), q(p(x))))$$

- Issues:
    - Overfitting: regularisation (weight decay, dropout)
    - initialisation: Scaling, batch-norm
    - hyperparameters: random vs grid search

# Image Recognition with Convolutional Networks

- Aim for translational invariance
  - make fully connected layer locally connected by sharing weights & removing connections
  - only 4, not 16 connections
  - and share weights
  - not flat in practice; have multiple channels (e.g. RGB=3 channels)
- Convolutional layer:
$$y = \sum_{i \in \text{Receptive field}} w_i x_i + b = w * x + b$$
(slide receptive field over image & compute dot product)
- for  $N \times N$  input &  $k \times k$  kernel, output  $\rightarrow M \times M$ , where  $M = N - k + 1$
- Padding:
  - VALID: no padding
  - SAME: pad with  $p = (k-1)/2$  zeros on all sides, so output same size as input
- Stride:
  - filter takes steps larger than 1
  - reduce computational cost / spatial resolution, & invariance to local translation
- Dilation:
  - ratio  $r =$  enlarge receptive field but keep same num parameters
  - $R' = k + (k-1)(r-1)$
- $M = \left\lceil \frac{N+2p-k'}{s} \right\rceil + 1$ , for input  $N \times N$ ,  $p$  padding,  $k'$  dilated kernel size &  $s$  stride length
- Deconvolution:
  - swap direction of convolution; increases resolution
  - used in backprop
- Pooling layer:
  - reduces spatial resolution
  - average, or max pooling over receptive field
  - can do global pooling too; full invariance to location

## • Building deep networks:



- usually shallow layers with a ReLU
- deeper usually better; More invariance & non-linearities
- typically build up no. channels over depth, before linear, FC layers
- main workhorse:  $3 \times 3$  convolution layers
- stacking these up gives effective larger receptive fields

## • Softmax

## • Training:

- hard to train very deep nets; vanishing or exploding gradients

## • Fixes:

### • Weight initialisations:

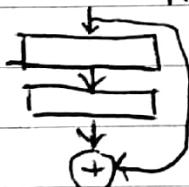
- 0 mean Gaussian low variance, or, eg, variance to preserve gradient magnitude

### • Batch normalisation:

- layer in network

- normalising input to zero mean, unit variance

### • Residual connections:

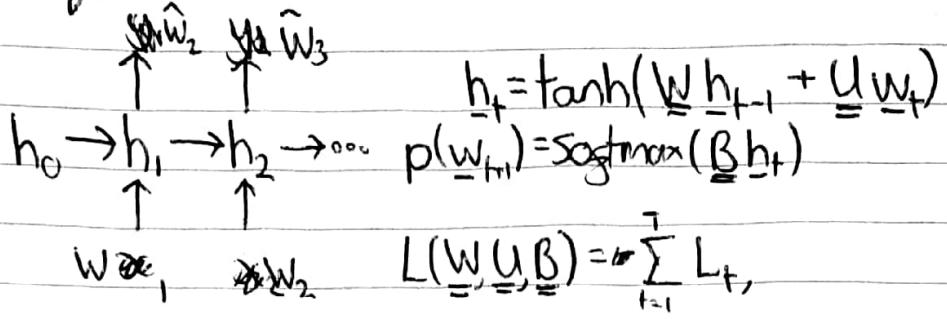


- very deep networks get higher training error, as hard to learn identity function

- add connections which skip layers

- also improves gradient flow, as can skip layers

# Sequences & Recurrent Neural Nets



$$\frac{\partial L_t}{\partial B} = \frac{\partial L_t}{\partial p_t} \frac{\partial p_t}{\partial h_t} \frac{\partial h_t}{\partial B}$$

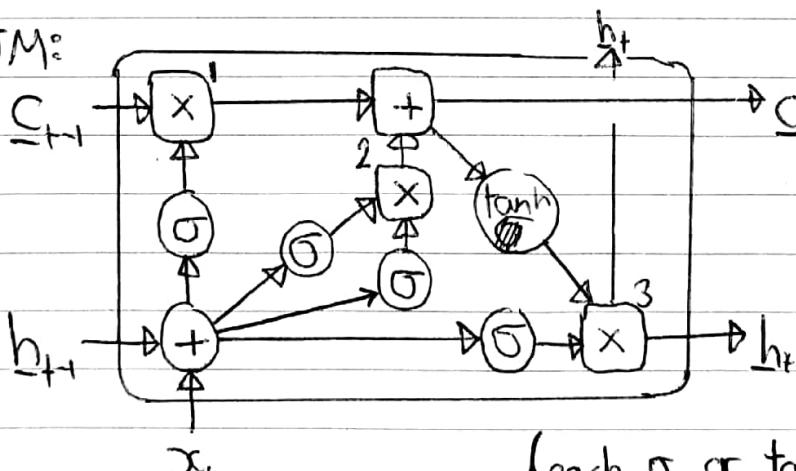
$$\begin{aligned} &= (p_t - \text{label}_t) h_t^\top \Rightarrow \frac{\partial L}{\partial B} = \sum_{t=1}^T \frac{\partial L_t}{\partial B} \\ \frac{\partial L_t}{\partial W} &= \frac{\partial L_t}{\partial p_t} \frac{\partial p_t}{\partial h_t} \frac{\partial h_t}{\partial W} = \sum_{k=1}^T \frac{\partial L_t}{\partial p_t} \frac{\partial p_t}{\partial h_t} \frac{\partial h_t}{\partial h_k} \frac{\partial h_k}{\partial W} \end{aligned}$$

- problem:  $h_t = W^\top h_{t-1}$ ,  $h_t \rightarrow \infty$  if  $\|W\| > 1$   
 $h_t \rightarrow 0$  if  $\|W\| < 1$

(vanishing or exploding gradients)

- makes learning long term dependencies very hard

- LSTM:



- cell state  $s_t$  allows for long term dependencies
- 1: forget gate
- 2: input gate
- 3: output gate

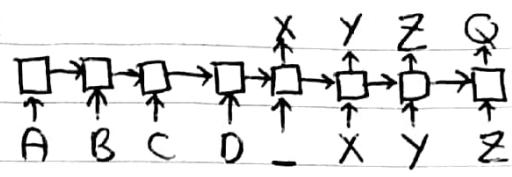
(each  $\sigma$  or  $\tanh$  is linear layer ( $w$  bias) followed by activation function)

- Can also help solve exploding gradients by thresholding:

$$\text{if } \|g\| > \text{threshold} \\ g \leftarrow \frac{\text{threshold}}{\|g\|} \times g$$

- Applications:

- Modelling languages; sequence of words
- Pixel prediction; image as sequence
- Sequence to sequence:
  - image captioning, speech
  - Machine translation



# Beyond Image Recognition: End-to-End Learning, Embeddings

- **Pretraining:**

- use a ready trained network, and use as base to build on
- don't update these params; just feed output into new network (or update slowly; low learning rate)
- e.g. preusing network pre-trained on ImageNet for other image datasets

- **Object detection:**

- As classification; classify all boxes (may be repeated objects; slow)
- As prediction; regres 4 corners of box; unknown no. boxes
- Bounding box proposals:
  - module predicts which boxes may be objects
  - then passes each box to a classifier

- **Semantic Segmentation:**

- label each pixel into an object class
- output full resolution map with 1 channel per class & probabilities for each pixel
  - 'fully convolutional networks'; no linear layers in network

- **Spatial transformers:**

- learn to transform input, e.g. rotate, to some canonical representation

- **Similarity / dissimilarity**

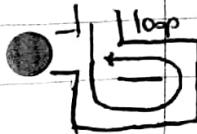
- inputs associated with each other; e.g. all images of a given persons face

- use different loss functions based on predicted similarity

- **Maze navigation:**

- problem: sparse rewards; we want some form of spatial knowledge;

- add auxiliary rewards for auxiliary tasks; depth prediction or loop closure prediction



# Optimisation for Machine Learning

objective

parameters

- Gradient Descent:  $\underline{\Theta}_{k+1} = \underline{\Theta}_k - \alpha_k \nabla h(\underline{\Theta}_k)$

Motivated by 1<sup>st</sup> order Taylor expansion  $h(\underline{\Theta} + \underline{\delta}) \approx h(\underline{\Theta}) + \nabla h(\underline{\Theta})^T \underline{\delta}$  greatest reduction in  $h$

reasonable for small  $\underline{\delta}$ , ie low learning rate

problem: oscillates along a valley

Cannot solve by fine tuning learning rate; either big zig-zag or slow progress

reason: we minimize local quadratic approximation to  $h(\underline{\Theta})$

$$h(\underline{\Theta} + \underline{\delta}) \approx h(\underline{\Theta}) + \nabla h(\underline{\Theta})^T \underline{\delta} + \frac{1}{2} \underline{\delta}^T H(\underline{\Theta}) \underline{\delta}$$

approximate  $H(\underline{\Theta})$  with  $L\mathbb{I}$ ; treat all directions

as having same curvature,  $\frac{1}{2} \underline{\delta}^T (L\mathbb{I}) \underline{\delta} = \|\underline{\delta}\|^2 \cdot \frac{L}{2}$

$$h(\underline{\Theta} + \underline{\delta}) \approx h(\underline{\Theta}) + \nabla h(\underline{\Theta})^T \underline{\delta} + \frac{L}{2} \|\underline{\delta}\|^2$$

$$\Rightarrow -\frac{1}{L} \nabla h(\underline{\Theta}) = \operatorname{argmin}_{\underline{\delta}} \left( h(\underline{\Theta}) + \nabla h(\underline{\Theta})^T \underline{\delta} + \frac{L}{2} \|\underline{\delta}\|^2 \right)$$

$\Rightarrow$  struggle when curvature varies a lot, ie  $H(\underline{\Theta}) \neq L\mathbb{I}$

Asymptotic performance;  
often stop before this  
point anyway

Worst case

Very pessimistic

- Convergence theory:

$$h(\underline{\Theta}_k) - h(\underline{\Theta}^*) \leq \frac{L}{2} \left( \frac{k-1}{k+1} \right)^{2k} \|\underline{\Theta}_0 - \underline{\Theta}^*\|^2 \quad (\text{for } \alpha_k = \frac{2}{L+1})$$

worst case  
upper  
bound

$\Rightarrow$  iterations for  $h(\underline{\Theta}_k) - h(\underline{\Theta}^*) \leq \epsilon$ :  $k \in \mathcal{O}(k \log \frac{1}{\epsilon})$   $k = \frac{L}{\epsilon}$

- Momentum:

$$\underline{v}_{k+1} = \eta_k \underline{v}_k - \alpha_k \nabla h(\underline{\Theta}_k)$$

$$\underline{\Theta}_{k+1} = \underline{\Theta}_k + \underline{v}_{k+1}$$

$\eta_k$ : friction constant  
 $\alpha_k$ : learning rate

• helps to fix zig-zagging problem by updating w/ moving average of gradient

- Convergence theory:

• another first order method (update is linear combination of previous gradient estimates)

•  $k \in \mathcal{O}(\sqrt{F} \log \frac{1}{\epsilon})$ ;  $\sqrt{F}$  factor speedup

worst case lower bound

upper bound for Nesterov, but

similar performance in practise



- Second order methods:

- Approximate  $h(\underline{\theta} + \underline{\delta})$  w/ 2<sup>nd</sup> order Taylor expansion

$$h(\underline{\theta} + \underline{\delta}) \approx h(\underline{\theta}) + \nabla h(\underline{\theta})^T \underline{\delta} + \frac{1}{2} \underline{\delta}^T \underline{H} \underline{\delta}$$

$$\Rightarrow \underline{H}^{-1} \nabla h = \underset{\underline{\delta}}{\operatorname{argmin}} (h(\underline{\theta}) + \nabla h(\underline{\theta})^T \underline{\delta} + \frac{1}{2} \underline{\delta}^T \underline{H} \underline{\delta})$$

$$\underline{\theta}_{k+1} = \underline{\theta}_k - \underline{H}^{-1} \nabla h(\underline{\theta}_k)$$

- Problems:

- approximation only valid near  $\underline{\theta}_k$

- 1<sup>st</sup> order methods approximate  $\underline{H}$  w/  $\underline{L} \underline{I}$ , here

we approximate the real  $H(\underline{\theta})$ , but this may

underestimate curvature as we move away from  $\underline{\theta}$

- Solution: constrain  $\underline{\delta}$  to be in some region  $R$ , near  $\underline{\theta}$

where our approximation remains good

- if  $R = \{\underline{\delta} : \|\underline{\delta}\|^2 \leq r\}$ , then doing our  $\underset{\underline{\delta} \in R}{\operatorname{argmin}}(\dots)$

is approximately equivalent to doing  $\underset{\underline{\delta}}{\operatorname{argmin}}(\dots)$ ,

but replacing  $H(\underline{\theta})$  with  $(H(\underline{\theta}) + \lambda^2 \underline{I})$ , where

$\lambda$  is some function of  $r$

- can also instead replace  $\underline{H}$  with another matrix, upper bounding curvature in over larger regions

- e.g. Generalized Gauss Newton Matrix

Fisher information matrix

Empirical Fisher info Matrix

- Generalized Gauss Newton:

- assume  $h(\underline{\theta}) = \sum h_i(\underline{\theta}) = \sum l(y_i, g(x_i, \underline{\theta}))$

- $\underline{G} = \sum \underline{J}_i^T \underline{H}_i \underline{J}_i$ ,  $\underline{J}_i$  is Jacobian of  $g(x_i, \underline{\theta})$  wrt  $\underline{\theta}$   
 $\underline{H}_i$  is Hessian of  $l(y_i, z_i)$  wrt  $z_i$

- $\underline{G} = \underline{H}$  if replace  $g(x_i, \underline{\theta})$  w/ 1<sup>st</sup> order approximations

- for square loss,  $\underline{H}_i = \underline{I} \Rightarrow \underline{G} = \sum \underline{J}_i^T \underline{J}_i$

- $\underline{G}$  always PSD, more conservative than Hessian & works much better in practise

## Optimisation 2

### • Curvature matrix approximations:

- or some other curvature matrix  $\underline{B}$   $\rightarrow$
- $\underline{\nabla}$  can have 10s of millions of dimensions
  - $\underline{H}$  is  $n \times n$  &  $O(n^3)$  to invert; impractical
  - Approximations:
    - diagonal:  $O(n)$  storage,  $O(n)$  to invert
      - only reasonable if eigenvalues of  $\underline{B}$  closely aligned w/ coordinate axes
    - popular is empirical Fisher (used by RMS-prop & Adam)
    - low-rank:  $O(rn)$  storage, computation & to invert
      - diagonal + rank  $r$  correlations
      - best if  $\underline{B}$  has few important eigenvectors w/ large eigenvalues
      - e.g. L-BFGS
    - block-diagonal:  $O(bn)$  to store,  $O(b^2n)$  to invert
      - can be awkward to compute (similar to diagonal)
      - e.g. TONGA
    - Kronecker-product:  $O(n)$  to store/compute,  $O(\underline{B}n)$  to invert
      - block diagonal of GGN / Fisher
      - blocks correspond to layers, each approximated by Kronecker-product of 2 smaller matrices

### • Stochastic optimisation:

- randomly sample mini-batch to compute gradient
- useful for momentum, but must be careful w/  $\alpha$  &  $\eta$  choice  $\star$  (and momentum)
- for second order methods, gives no asymptotic improvement
  - but pre-asymptotically can help, which is what we care most about

# Attention & Memory in Deep Learning

- Deep networks have implicit attention; respond more to certain parts of input than to others (shown by Jacobian)

$$J_{ij} = \frac{\partial y_i}{\partial x_j} \quad (\text{compute w/ backprop; set loss to output activation})$$

- In RNNs, attention <sup>shows</sup> a form of memory; attention to past inputs to produce present output
  - shown by sequential Jacobian;  $J_k^t = \left( \frac{\partial y_k^t}{\partial x_1}, \frac{\partial y_k^t}{\partial x_2}, \dots \right)$
  - e.g. delayed dot of  $i$  in handwriting, or out-of-order machine translation

- Explicit attention:

- limit data presented to network in some way
  - + ↑ computational efficiency, scalability, interpretability

Extra

data  
↓

Attention  
model

Output

- attention model chooses what  $bt$  should be visible to network in next time step; the 'glimpse'

• whole system recurrent, even if network isn't

Network

Glimpse

Input

• Glimpse:  $P(g|a)$

- generally define probability distribution over glimpses of data, given some set of attention outputs from the network

not  
differentiable

↓  
hard  
attention;  
specific  
window  
passed to  
network

- can train this w/ RL techniques; reward = loss induced by the glimpse (g can use RL whenever a module not differentiable)
  - glimpse could e.g. be Gaussian w/ given mean & variance, over coordinates of an image, defining a window to pass to network
  - Generally prefer soft attention; differentiable, can train end-to-end with backprop!
  - take expectation from  $P(g|a)$  rather than a sample

# Deep Learning for NLP

## • Word Vectors:

- sparse, orthogonal & semantically weak
- we want dense representation based on context/use of words
- Count-based:
  - $C$  = context words
  - count words in  $C$  in window around a given word to determine its vector representation
  - ~~also~~ need to normalise counts; distinguish informative high counts from uninformative (eg tf-idf, pmi...)
  - Measure Similarity by cosine similarity  $(\text{cosine}(u, v) = \frac{u \cdot v}{\|u\| \|v\|})$

## • Neural embedding models:

$$\text{arg}\underset{\Theta, E}{\min} L(\Theta, E), L(\Theta, E) = - \sum_t \sum_{t_i} \text{score}(t_i, c(t_i), \Theta, E)$$

$E$  is embedding

$t$  is a word in vocab  $t_i$  is an instance of  $t$  in data

Score function important;

(context)  $\cdot$  must embed  $t_i$  with  $E$ , be differentiable & take  $c(t_i)$  into account

## • Discrete Language Models:

• Goal: estimate  $P(S_{1:n}) = \prod_{i=1}^n P(S_i | S_{1:i-1})$

•  $S_{1:i-1}$  variable length; use  $n$ -gram (fixed window size)

• too small  $n$  = bad model, too large = counts too low; suffer from sparsity

## • Continuous Condition Language Models:

### • LBL/NLM:

• co-learn word embeddings, then use  $n$ -gram model on them

• log-bilinear model = linear gram embeddings to probabilities

• neural language model = sigmoid activation function

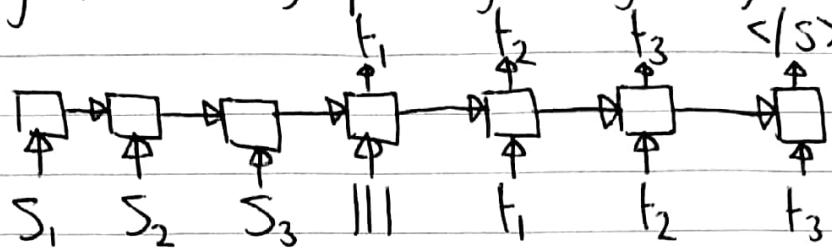
### • RNN-LMs:

• no longer need  $k$ -Markov assumption

(in practice still some soft limitations on  $k$  based on data, memory size etc)

- Sequence-to-sequence Mapping w/ RNNs:

- eg translation, parsing (string  $\rightarrow$  tree), or more general computation



- train to Maximise likelihood of  $t$  given  $s$

- problems:

- non-adaptive capacity, & modelling target sentence dominates training

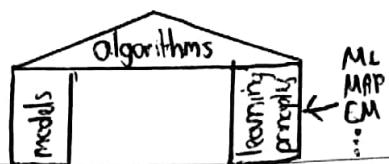
- encoder "gradient starved"  $\leftarrow$  (attention models can help here)

- Semantic Composition:

- language tree structure; get parse tree distribution & exploit this with an RNN

- ~~con, eg, annotat~~

# Unsupervised Learning & Generative Models



## • Probabilistic Machine Learning:

### • Model types:

- Unconditional:  $p(x)$ , no targets / labels

- what generative model usually means

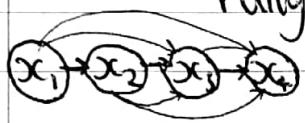
- Conditional  $p(z|x)$

- e.g. supervised learning; regression & classification

all are  
actually  
generative  
models

### • Types of generative Models:

#### • Fully observed models:



- model data directly; don't introduce any unobserved variables

#### • directed graphical models:

- Simple to learn parameters (log-likelihood computable)

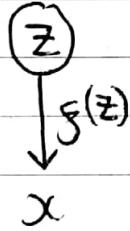
- scales well

- order sensitive

- undirected hard to learn parameters

- slow to sample generativ; need to go thru elements sequentially

#### • Implicit models:



- transform unobserved noise source using parameterised function

- easy to sample & compute expectations

- hard to maintain invertibility & to optimise

- lack of a noise model

#### • Latent variable models

- local random variables representing hidden causes

- easy to sample

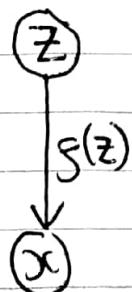
- encode hierarchy, depth, structure

- Latents compress & represent data

- can score, compare & select models by likelihood

- inversion difficult

- Marginal likelihood requires approximations



## • Inference in Prescribed Probabilistic Models:

### • Model evidence:

$$p(\underline{x}) = \int p(\underline{x}, \underline{z}) d\underline{z}$$

• aim to improve this for given data sets

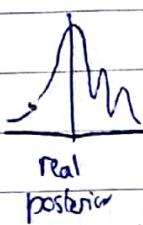
• intractable in general; transform into expectation over known distribution

### • Jensen's inequality:

• for concave functions  $g$ ,

$$g(E(\underline{x})) \geq E(g(\underline{x}))$$

### • Variational Inference



$$p(\underline{x}) = \int p(\underline{x}|\underline{z}) p(\underline{z}) d\underline{z}$$

$$p(\underline{x}) = \int p(\underline{x}|\underline{z}) \frac{p(\underline{z})}{q(\underline{z})} q(\underline{z}) d\underline{z}$$

Jensen's inequality

$$\log p(\underline{x}) \geq \int q(\underline{z}) \log(p(\underline{x}|\underline{z}) \frac{p(\underline{z})}{q(\underline{z})}) d\underline{z}$$

$$= \int q(\underline{z}) \log p(\underline{x}|\underline{z}) - \int q(\underline{z}) \log \frac{p(\underline{z})}{q(\underline{z})} d\underline{z}$$

$$\rightarrow = E_{q(\underline{z})} [\log p(\underline{x}|\underline{z})] - \text{KL}[q(\underline{z}) || p(\underline{z})]$$

•  $q(\underline{z})$  is our approximation to  $p(\underline{z})$ , using some known, simple distribution

•  $\text{KL}(q || p)$  is a 'penalty' term, as  $q \neq p$

• optimise for  $q$ ; eg if Gaussian, variational parameters are mean & variance to optimise for

• this can tighten the bound, getting closer to true marginal likelihood

max likelihood for other

• Only some parts of model given prob distributions, eg latent vars (in Variational Bayes, all unknown quantities are probability distributions)

+ applicable to all probabilistic models

+ convert integral problem to optimisation

+ numerically stable & parallelisable, fast to converge

→ need to choose form of  $q$ , and only approximate

- local minima

- can under estimate variance in posterior

# Unsupervised Learning & Generative Models

## • Family of Approximate Posterior distributions:

### • Free form:

$$\cdot \text{get } q(z) \propto p(z) e^{\log p(x|z, \theta)}$$

• problem: Solving is normalisation of original problem

### • Fixed form:

• specify explicit form of  $q$ ;  $q_{\phi}(z) = f(z; \phi)$

### • Mean field methods:

• assume distribution is factorised

• can introduce structure to factorisation (dependencies) to be more expressive

• e.g. autoregressive approximation: (condition on all previous variables)

$$q(z) = \prod_i q_i(z_i | z_{-i})$$

### • Variational Latent Gaussian models:

•  $z \sim N(z | 0, I)$ ,  $y \sim p(y | f_{\phi}(z))$ ,  $q(z) = \prod N(z_i | \mu_i, \sigma^2)$

$$F(y, q) = E_{q(z)} [\log p(y|z)] - KL[q(z) || p(z)]$$

$$= E_{q(z)} [\log p(y|z)] - \sum_i KL[N(z_i | \mu_i, \sigma^2) || N(z_i | 0, 1)]$$

$$= E_{q(z)} [\log p(y|z)] - \frac{1}{2} \sum_i [0_i + \mu_i^2 - 1 - \ln \sigma_i^2]$$

## • Variational Optimisation:

### • EM algorithm

### • E Step:

$$\phi \propto \nabla_{\phi} F(x, q)$$

(Variational params)

$$\cdot \phi_n \propto \nabla_{\phi} E_{q_n(z)} [\log p_{\phi}(x_n | z_n)] - \nabla_{\phi} KL[q(z_n) || p(z_n)]$$

### • M step:

$$\Theta \propto \nabla_{\Theta} F(x, q) \leftarrow (\text{model params})$$

$$\Theta \propto \frac{1}{n} \sum_n E_{q_n(z)} [\nabla_{\Theta} \log p_{\Theta}(x_n | z_n)]$$

• can do stochastic approach; do EM on minibatches

• more scalable

• if expected log likelihood & KL-divergence can't be calculated, use Monte-Carlo methods

• doubly stochastic estimation; stochastic sum minibatch & Monte-Carlo estimation

## • Learning in Implicit Probabilistic Models:

### • Estimation by comparison:

• compare estimate  $q(x)$  with true  $p^*(x)$

$$\frac{p^*(x)}{q(x)} = 1 \quad p^*(x) = q(x)$$

• test a hypothesis, then modify  $q$  slightly to better match  $p^*$  result

### • Adversarial Learning:

Comparison loss  $\rightarrow \Theta \propto \nabla_{\Theta} E_{p^*(x)} [\log D_{\Theta}(x)] + \nabla_{\Theta} E_{q_{\Phi}(x)} [\log (1 - D_{\Theta}(x))]$

Generative loss  $\rightarrow \Phi \propto -\nabla_{\Phi} E_{q_{\Phi}(z)} [\log (1 - D_{\Theta}(f_{\Phi}(z)))]$

$$\left( \frac{p^*(x)}{q(x)} = \frac{p(x|y=1)}{p(x|y=-1)} = \frac{p(y=1|x)}{p(y=-1|x)} = \frac{D_{\Theta}(x)}{1 - D_{\Theta}(x)} \right)$$