Department of Medical Physics and Biomedical Engineering

Centre for Medical Image Computing (CMIC)

Wellcome / EPSRC Centre for Interventional and Surgical Sciences (WEISS)



Deep Learning

MPHY0041 Machine Learning in Medical Imaging

Yipeng Hu yipeng.hu@ucl.ac.uk

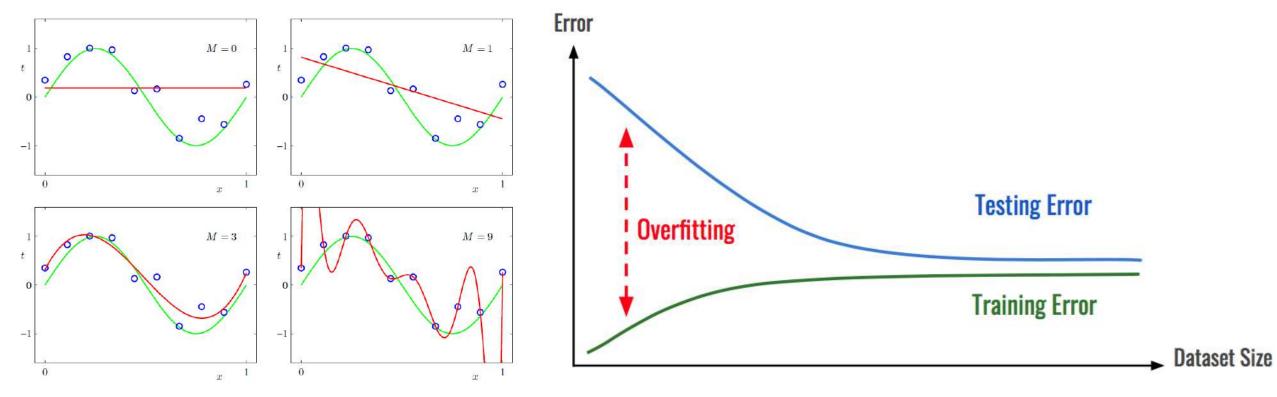


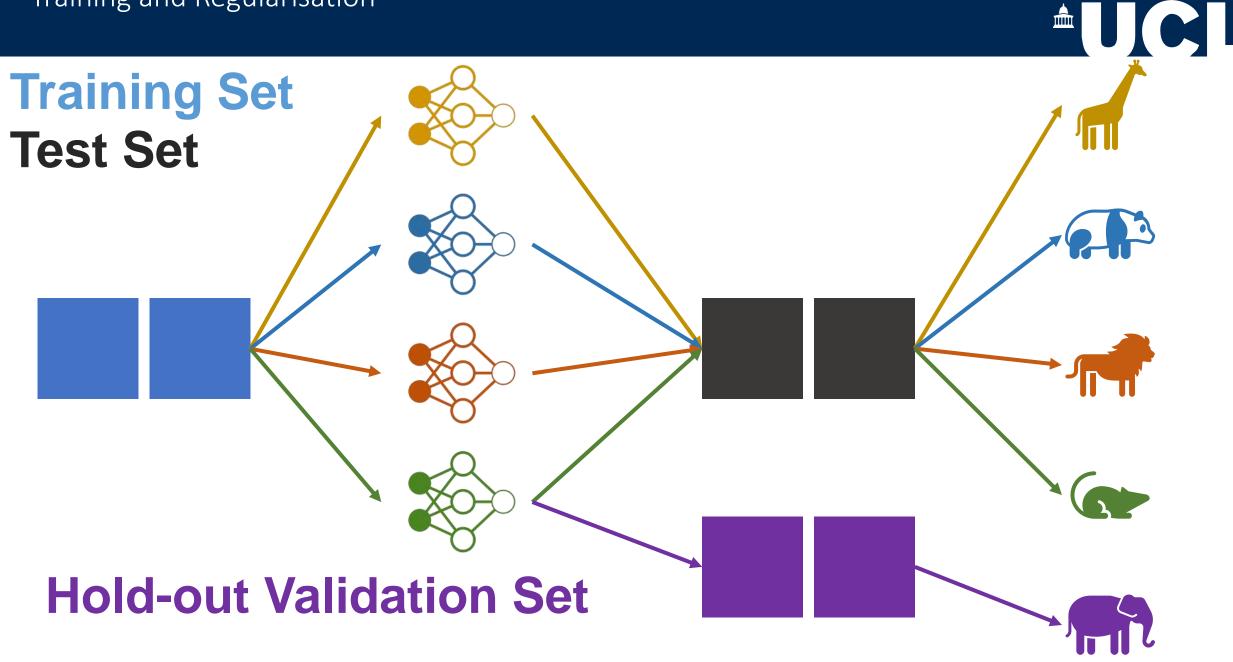


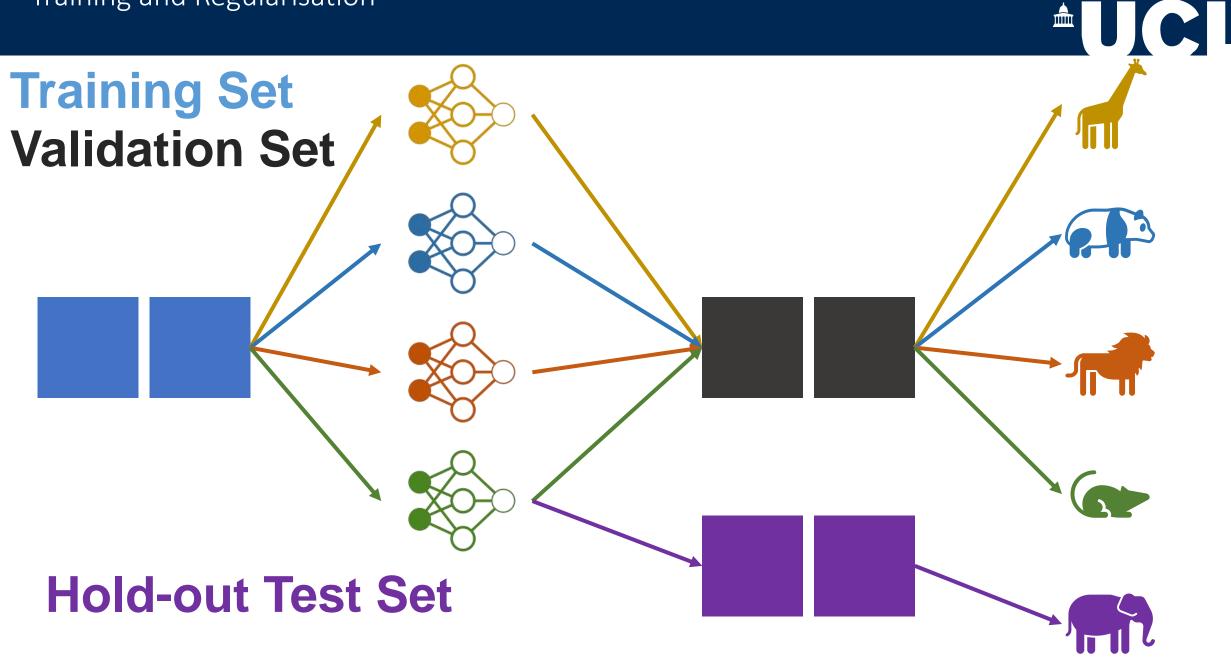
Training and Regularisation

Training and Regularisation











Training and Regularisation | Optimisation

Training and Regularisation | Optimisation



Data generating distribution

$$J^*(\boldsymbol{\theta}) = \mathbb{E}_{(\boldsymbol{x}, \mathbf{y}) \sim p_{\text{data}}} L(f(\boldsymbol{x}; \boldsymbol{\theta}), y).$$

Training data distribution

$$J(\boldsymbol{\theta}) = \mathbb{E}_{(\boldsymbol{x}, y) \sim \hat{p}_{\text{data}}} L(f(\boldsymbol{x}; \boldsymbol{\theta}), y),$$

Empirical risk

wrt. mini-batches

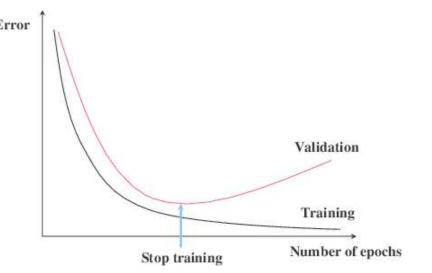
$$\mathbb{E}_{\boldsymbol{x}, \mathbf{y} \sim \hat{p}_{\text{data}}(\boldsymbol{x}, y)}[L(f(\boldsymbol{x}; \boldsymbol{\theta}), y)] = \frac{1}{m} \sum_{i=1}^{m} L(f(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), y^{(i)})$$

Large batches and small batches

"Regulairsing for generalisability"

Surrogate losses

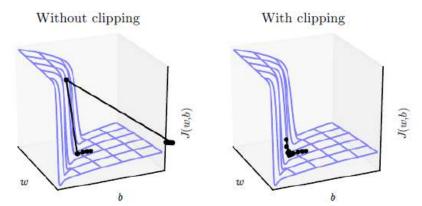
CE and Early stopping based on validation set





Local minima, non-identifiability and "a good enough solution"

Gradient explosion



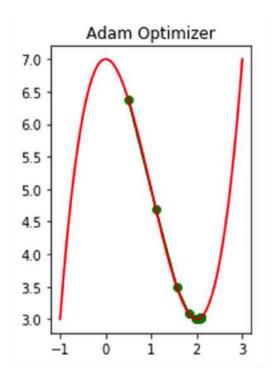
Algorithms

- SGD and mini-batch gradient descent

$$oldsymbol{v} \leftarrow lpha oldsymbol{v} - \epsilon
abla_{oldsymbol{ heta}} \left(rac{1}{m} \sum_{i=1}^{m} L(oldsymbol{f}(oldsymbol{x}^{(i)}; oldsymbol{ heta}), oldsymbol{y}^{(i)})
ight)$$

- Momentum

- $\theta \leftarrow \theta + v$.
- Parameter initialisation, wrt. Symmetry, depth etc
- Adaptive methods
 - high-order derivatives (Newton's Method, Conjugate gradient)
 - Past derivative values (AdaGrad, RMSProp)
 - Adaptive momentum (Adam)

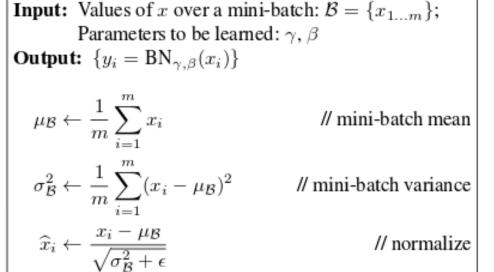




Reparameterisation for optimising deep networks

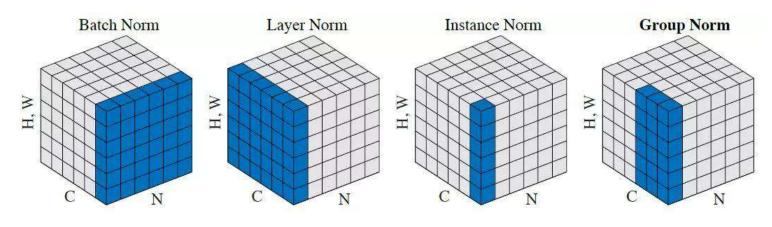
// scale and shift

Batch normalisation and variants



Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.

 $y_i \leftarrow \gamma \hat{x}_i + \beta \equiv BN_{\gamma,\beta}(x_i)$





Training and Regularisation | Constrain Parameters



Parameter norms

Weight decay by penalising L²- and L¹-norms

$$\ln p(\mathbf{w}|\mathbf{t}) = -\frac{\beta}{2} \sum_{n=1}^{N} \{t_n - \mathbf{w}^{\mathrm{T}} \phi(\mathbf{x}_n)\}^2 - \frac{\alpha}{2} \mathbf{w}^{\mathrm{T}} \mathbf{w} + \text{const.}$$

Difference between L2- and L1-norms

$$\Omega(\boldsymbol{\theta}) = ||\boldsymbol{w}||_1 = \sum_i |w_i|$$
 $\nabla_{\boldsymbol{w}} \tilde{J}(\boldsymbol{w}; \boldsymbol{X}, \boldsymbol{y}) = \alpha \operatorname{sign}(\boldsymbol{w}) + \nabla_{\boldsymbol{w}} J(\boldsymbol{X}, \boldsymbol{y}; \boldsymbol{w})$

 L^2 norm gradient scales with w, therefore less likely to become zero, i.e. sparsity.

Training and Regularisation | Constrain Parameters



Parameter sharing

- CNN
- RNN
- Tiring: e.g.

$$\Omega(\boldsymbol{w}^{(A)}, \boldsymbol{w}^{(B)}) = \|\boldsymbol{w}^{(A)} - \boldsymbol{w}^{(B)}\|_{2}^{2}$$

- Multi-task learning*
- Semi-supervised learning*



Training and Regularisation | Augment Data

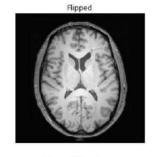
Training and Regularisation | Augment Data

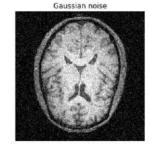


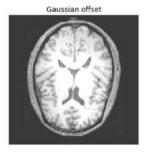
- Geometric transformation,
 - e.g. flipping, cropping, rotation
- Colour space,
 - e.g. RGB, intensity, contrast
- Imaging-specific,
 - e.g. photometric transformation for camera, speckle for ultrasound, bias field for MR, Projective transformation for CT
- Spatial transformation
 - inc. geometric transformation, affine, nonlinear deformation

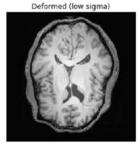
Generative models*

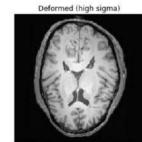














The "mixup"

Image

Label

$$\tilde{x} = \lambda x_i + (1 - \lambda) x_j,$$

where x_i, x_j are raw input vectors $\tilde{y} = \lambda y_i + (1 - \lambda)y_i$, where y_i, y_i are one-hot label encodings

[0.7, 0.3][1.0, 0.0] [0.0, 1.0] cat dog cat dog cat dog



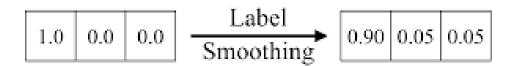
Training and Regularisation | Add Randomness



Augment input*

Add noise to weights

Label smoothing (assuming noisy output)



Dropout*



Training and Regularisation | Ensemble



Model Averaging

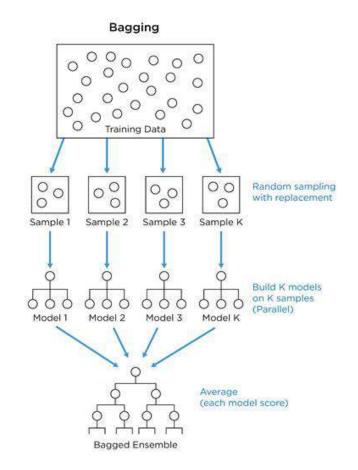
$$p(X) = \sum_{Y} p(X|Y)p(Y)$$

e.g. Bagging/bootstrap aggregating

- Weighted averaging
- Majority voting

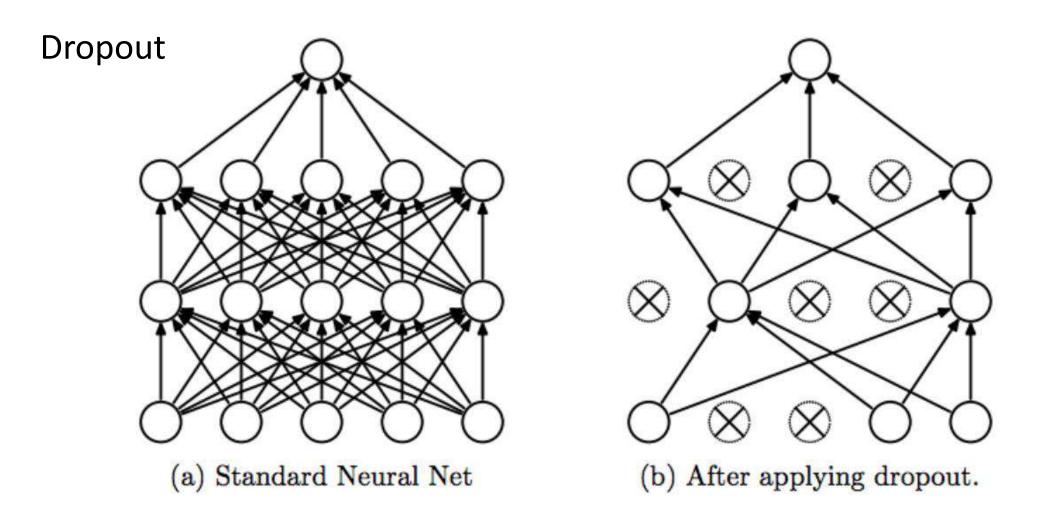
Boosting

Weak learners -> strong learners, i.e. a meta-algorithm for reducing bias



Competition winners!





- Test-time dropout for ensemble and uncertainty estimation



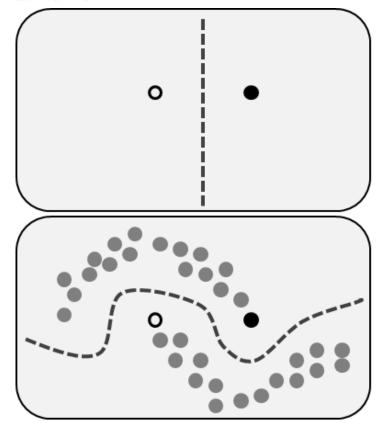
Training and Regularisation | Semi-Supervised Learning

Training and Regularisation | Semi-Supervised Learning



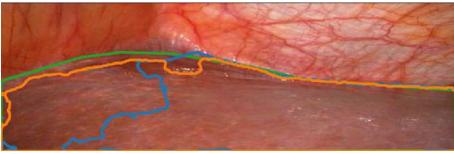
Supervised learning: labelled data = data + labels $P(\mathbf{x}, \mathbf{y})$ Semi-supervised learning: labelled data + unlabelled data $P(\mathbf{x})$

 $P(\mathbf{y} \mid \mathbf{x})$



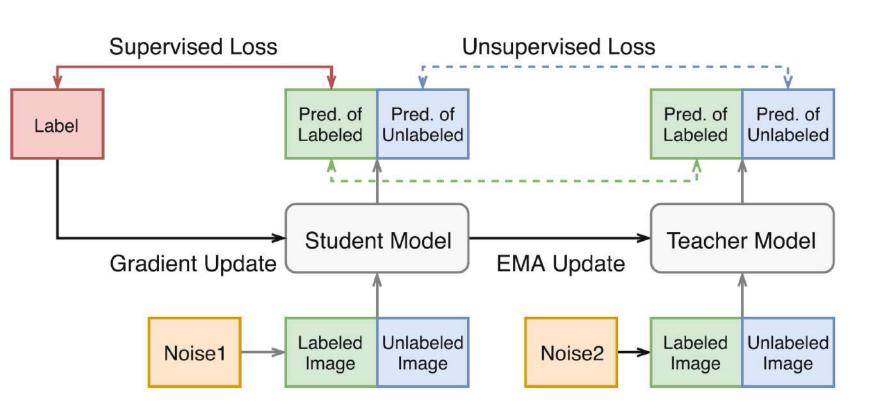
Laparoscopic video image segmentation





Training and Regularisation | Semi-Supervised Learning





Student model is guided by

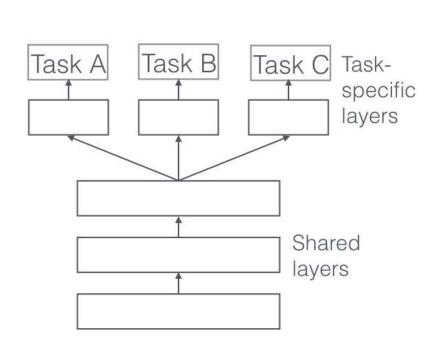
- Supervised label
- Teacher's prediction

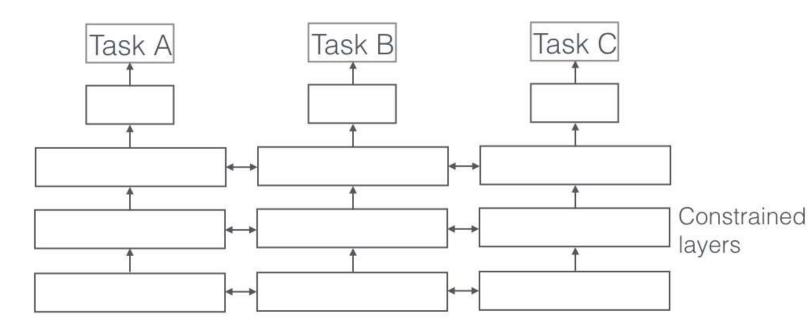
Teacher model's weight is the moving average of student model.

Random affine transformation is used as noise.



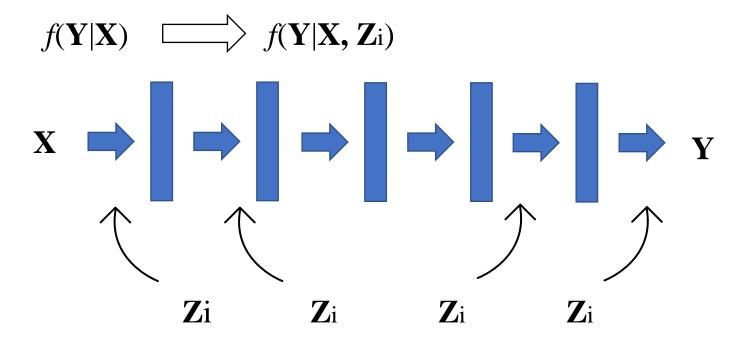






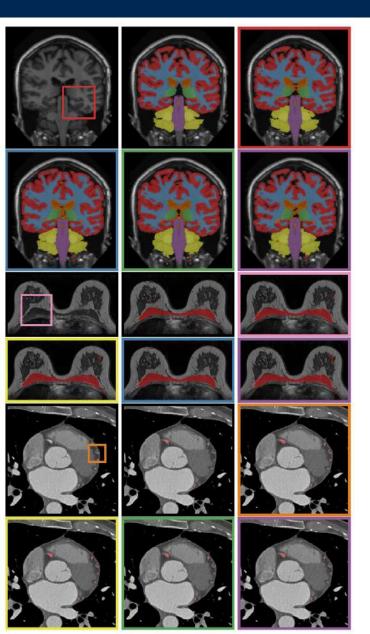


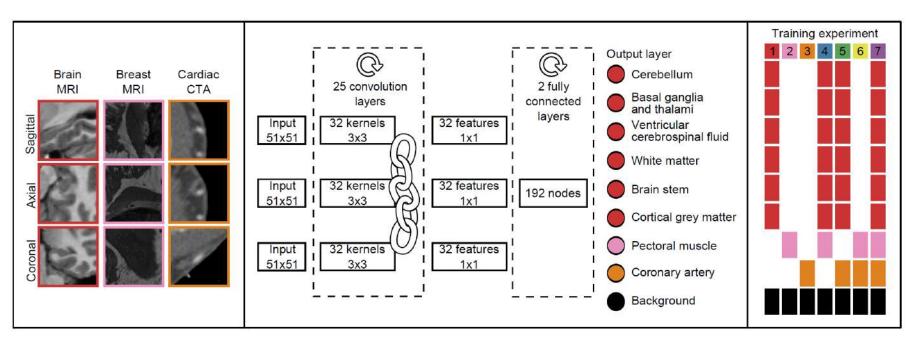
Zi: One-hot task index / indicator / descriptor



- Shared and task specific parameters (negative transfer)
- Concatenation/summation
- Multiplication conditioning (gating/multi-head)









Links between regularisation methods

- Regularisation, e.g. architecture, optimisation strategies
- Parameter sharing, e.g. semi-supervised, multi-task
- Randomness/noise, e.g. data augmentation, dropout
- Ensemble, e.g. dropout