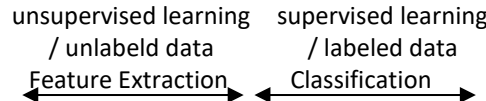
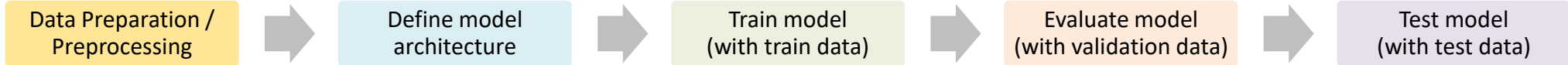


# DEEP LEARNING - OVERVIEW



Model-Selection / Cross-Validation



we need lots of data!

structured data -> tables/databases  
unstructured data -> audio/image/text

**Shuffle** data -> to get same characteristics

**Flatten** (from 28x28 to 784) -> dimension more dimension -> more data (curse of dimensionality)

**Split** data into 3 subsets -> avoid overfitting  
- training (large 98-60% small) -> train model  
- validation (1-20%) -> tune hyper-params  
- test set (1-20%) -> evaluate model

**Normalise** all data -> faster convergence  
- z-normalisation  
- shift -> improve robustness of learning  
- scale -> faster convergence and accuracy

$$x_k^{(i)} = \frac{x_k^{(i)} - \mu_k}{\sigma_k} \quad \mu_k = \frac{1}{N} \sum_{k=1}^N x_k^{(i)} \quad \sigma_k^2 = \frac{1}{N} \sum_{k=1}^N (x_k^{(i)} - \mu_k)^2$$

- min/max rescaling → [0,1]  
$$x_k^{(i)} = \frac{x_k^{(i)} - \min_j(x_k^{(j)})}{\max_j(x_k^{(j)}) - \min_j(x_k^{(j)})}$$

- min/max normalisation → [-1,1]  
$$x_k^{(i)} = 2 * \frac{x_k^{(i)} - \min_j(x_k^{(j)})}{\max_j(x_k^{(j)}) - \min_j(x_k^{(j)})} - 1$$

- modify output to 1-hot target (4 -> 00010)

**Data Augmentation** -> avoid overfitting

$$X = \begin{bmatrix} | & \dots & | \\ x^{(1)} & \dots & x^{(m)} \\ | & \dots & | \end{bmatrix} \in \mathbb{R}^{n \times m}$$

$$Y = [y^{(1)} \dots y^{(m)}] \in \mathbb{R}^{1 \times m}$$

**Type**

- Sequential (1 to 1 layer, stack)
- Functional (multiple paths/in/out)

**Layers** ( $w * x + b$ )

- Dense, fully connected (# neurons)
- Conv2D (#filter, #channel, filter\_size, ...)
- Activation
- MaxPooling2D (pool size)
- Dropout (dropout rate p)
- Flatten
- Batch-Normalisation
- Fall Back layers -> for deep networks
- Bottleneck layer -> reduce params
- GRU (Gated Recurrent Units)
- LSTM (Long Short Term Memory)

**Activation functions / Prediction:**

- > robustness and performance
- Heavyside
- Sigmoid
- Tanh
- Recti-Linear Unit (ReLU)
- Leaky ReLU ( $\alpha$ )
- Exponential Linear Unit (ELU) ( $\alpha$ )
- Identity
- Softmax → final layer

**Classification / Output**

- Binary classification  $\begin{cases} 1 = \text{yes} \\ 0 = \text{no} \end{cases}$
- Multi classification ( $K$  classes)

**Sample architectures**

Standard NN -> simple application	
CNN -> image recognizing	
RNN -> sequenced data	
Autoencoders -> reproduce input	
Transfer learning -> reuse feature extr	

**Model**

**Computation**  
Prefer TPU to GPU to CPU

**0. Parameter initialisation**  $\theta = (w, b)$   
rand initialisation -> break symmetries  
normalize per layer (Xavier)

**1. Propagate and make Prediction**  
 $\hat{y}(x) = \text{round}(h_\theta(x)) \in \{0,1\}$   
 $\hat{y}(x) = \sigma(w^T x + b) \rightarrow$  logistic regression  
 $\hat{y}(x) = \underset{l}{\text{argmax}} \{h_{\theta_l}(x)\} \in \{0 \dots K - 1\}$

**2. Measure Diff**  
Cost/Loss(Error) Functions  
- # of correctly classified samples  
- Mean Square Error

$$J_{MSE}(\theta) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

- Cross-Entropy -> convex function

$$J_{CE}(\theta) = -\frac{1}{m} \sum_{i=1}^m \log(h_{\theta, y^{(i)}}(x^{(i)}))$$

-> for classification

High Bias: Model is too simple for data.  
High Variance: Model adapts bad to new data.

**3. Backprop & update weights**  
**Learning algorithm (Correct)**

- Perceptron learning rule
- $w \leftarrow w - \alpha * (\hat{y}^{(i)} - y^{(i)}) * x^{(i)}$   
 $b \leftarrow b - \alpha * (\hat{y}^{(i)} - y^{(i)})$
- **Gradient descent**  $\theta_{t+1} \leftarrow \theta_t - \alpha * \nabla_{\theta} J(\theta_t)$ 
  - Batch Gradient Descent
  - Stochastic Gradient Descent
  - Mini-Batch Gradient Descent
- learning rate  $\alpha$   
too large -> not converge  
too small -> slow convergence
- batch size
- # epochs / iterations  
too small -> not optimal value  
too big -> overfitting
- **Regularisation** -> prefer simple models
  - Weight Decay -> add penalty  
 $J = J_0 + \lambda \Omega(W)$
  - L1  $\Omega(W) = \|W\|_1 = \sum_{l,k,j} |W_{kj}^{(l)}|$
  - L2  $\Omega(W) = \|W\|_2^2 = \sum_{l,k,j} |W_{kj}^{(l)}|^2$
- Dropout (temp ignored) -> more robust
- Early Stopping -> avoid overfitting
- Data Augmentation -> more tolerant

**Advanced Optimizer**

- Momentum / Nesterov (surpass flat, moment  $\beta_1$ )  
 $m \leftarrow \beta_1 m + \alpha \nabla J(\theta + \beta_1 m)$
- RMSprop (adjust  $\alpha$  per direction, RMS decay  $\beta_2$ )  
 $s \leftarrow \beta_2 s + (1 - \beta_2) \nabla J(\theta) \odot \nabla J(\theta)$   
 $\theta \leftarrow \theta - \frac{\alpha}{\sqrt{s + \epsilon}} \odot \nabla J(\theta)$
- Adam (combination of others,  $\beta_1, \beta_2$ )

test loss  
test accuracy

**Confusion matrix**

		predict		
		a	b	c
Act	a	120	21	7
	b	8	131	63
	c	12	30	80

Overall =  $\frac{\sum \text{diag elem}}{\# \text{samples}}$   
Accuracy = Overall  
Error = 1 - Accuracy

**Confusion Table**

		Predicted	
		Pos	Neg
Act	Pos	TP	FN
	Neg	FP	TN

$\text{accuracy} = \frac{TP + TN}{\# \text{sample}}$   
 $\text{sensitivity} = \frac{TP}{TP + FN}$   
 $= \text{recall} = \frac{TP + FP}{TP}$   
 $\text{precision} = \frac{TP}{TP + FP}$   
 $F \text{ Score} = \frac{2 \text{ prec} * \text{reca}}{\text{prec} + \text{reca}}$   
 $\text{specificity} = \frac{TN}{TN + FP}$

**Parameter learning**  
Hyper parameter tuning

**Calculus per layer**

one sample	multiple sample
$W^{(l)}: n_l \times n_{l-1}$	$W^{(l)}: n_l \times n_{l-1}$
$b^{(l)}: n_l \times 1$	$b^{(l)}: n_l \times m$ (broadc.)
$a^{(l-1)}: n_{l-1} \times 1$	$A^{(l-1)}: n_{l-1} \times m$
$a^{(l)}: n_l \times 1$	$A^{(l)}: n_l \times m$
$z^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)}$	$Z^{(l)} = W^{(l)} A^{(l-1)} + b^{(l)}$
$a^{(l)} = g^{(l)}(z^{(l)})$	$A^{(l)} = g^{(l)}(Z^{(l)})$

Chain rule:  
 $L(g(h(x))) \rightarrow dL = \frac{\delta L}{\delta g} \frac{\delta g}{\delta h} \frac{\delta h}{\delta x} dx$

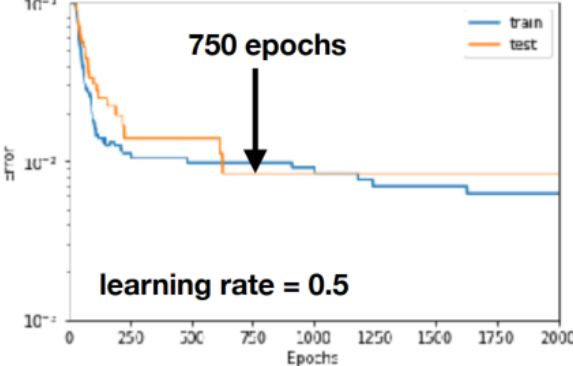
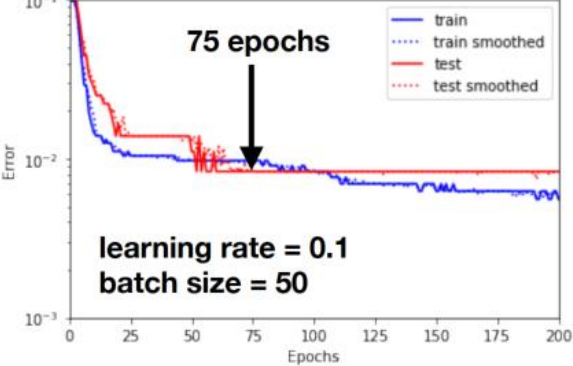
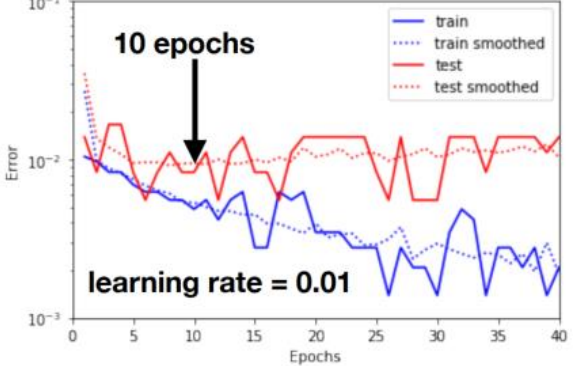
**Activation Functions**

Name	Output	function	derivate	grafic	good (+)	bad (-)	practical
<b>Heavyside</b>	[0,1]	$H(z) = \begin{cases} 1 & (z \geq 0) \\ 0 & (z < 0) \end{cases}$	$H'(z) = 0$		easy	not differentiable	no practical use
<b>Sigmoid</b>	[0..1]	$\sigma(z) = \frac{1}{1 + e^{-z}}$	$\sigma'(z) = \sigma(z)(1 - \sigma(z))$		smooth, gradient works	saturation regions leading to vanishing gradients	illustrative examples, binary classification (output)
<b>Tanh</b>	[-1..1]	$f(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$	$f'(z) = 1 - \tanh(z)^2$		smooth, gradient works output centred around 0	saturation regions leading to vanishing gradients	better than sigmoid except for output
<b>Recti-Linear unit (ReLU)</b>	[0..∞]	$f(z) = \max(0, z)$	$f'(z) = \begin{cases} 0 & (z < 0) \\ 1 & (z \geq 0) \end{cases}$		Alleviate the vanishing gradient problem	suffer from dying units problem $z < 0$	increasingly used as de facto standard
<b>Leaky Recti-Linear Unit (Leaky ReLU)</b>	[-∞..∞]	$f(z) = \max(\alpha z, z)$ $0 < \alpha < 1 \rightarrow \alpha = 0.01$	$f'(z) = \begin{cases} \alpha & (z < 0) \\ 1 & (z \geq 0) \end{cases}$		Alleviates both problems (vanishing gradient and dying units)		
<b>Exponential Linear Unit (ELU)</b>	[-∞..∞]	$f(z) = \begin{cases} \alpha(e^z - 1) & (z < 0) \\ z & (z \geq 0) \end{cases}$				more expensive to compute gradient vanishes at small negative values	
<b>Identity</b>	[-∞..∞]	$f(z) = z$	$f'(z) = 1$			network is linear -> useless to go deeper (lin+lin=lin)	only used in specific cases
<b>SoftMax</b>		$softmax(z) = \frac{e^z}{\sum_{j=1} e^z}$				used in last layer	multi classification (output)

**Layers**

Layer	Goal	Parameters	Output shape	# of Parameters, $n_0 = n_x, n_L = n_{classes}$
Dense, fully connected	connect two layer	# neurons	# of neurons	$(n_{l-1} + bias) * n_l$
Dropout	rand disable neurons	p dropout rate	no change	0
Conv2D	for image processing feature detection	image_size W, # channels C, # filters D, stride $s_w$ , padding $P_w$ , kernel size = filter size w	$\frac{(W - w + 2P_w)}{s_w} + 1$   $P_w = \frac{w - 1}{2}$	$D * \left( \underbrace{whC}_{weight} + \underbrace{1}_{bias} \right)$
Activation (e.g. ReLU)	insert non-linearity	only in Leaky-ReLU or ELU: $\alpha$	no change	0
MaxPooling2D	reduce output size	pool size	$\frac{input\ size}{pool\ size}$	0
Flatten	to calculate properly	-	$width * height * depth$	0
SimpleRNN	for sequential data	# units, # unique_chars (=features)	# of units	$(\# features + \# units + 1) * units$
Batch-Normalis. / Gradient clip	faster learning		no change	
Fall Back layers	for deep networks			
Bottleneck layer	reduce params			
GRU (Gated Recurrent Units)				
LSTM (Long Short Term Memory)	longer memory			

**Learning Schemes**

BGD: Batch Gradient Descent	MBGD: Mini-Batch Gradient Descent (Compromise)	SGD: Stochastic Gradient Descent
		
Summation over all the training set.	Summation only over subset of the training samples.	No summation, just a single sample.
Smooth	Wiggling curves, need smoothing.	Slightly wiggling curves.
Strictly decreasing curves	Cost curve not necessarily decreasing. Ends up wandering closely around the global minimum - not a problem in practise.	
Slowest, update for all training samples.	Faster for large m. Much less epochs than BGD, but more than SGD needed. Faster since parameters are updated for each mini-batch	Faster since parameters are updates for each training sample.
All data in RAM (No <b>out-of-core support</b> )	Can handle very large data sets, that do not fit in memory ( <b>out-of-core learning</b> )	
No incremental learning possible.	Allows incremental learning ( <b>online learning, on-the-fly</b> )	
Easily parallelised on GPU and HPC.		Not easily parallelised.
The learning principle is generalisable to many other "hypothesis families"		
for each epoch: $\Delta w_j := \alpha * \sum_i (target^{(i)} - output^{(i)}) x_j^{(i)}$ for each weight j: $w_j \leftarrow w_j - \Delta w_j$	for each epoch: $\Delta w_j = \alpha * \sum_i^b (target^{(i)} - output^{(i)}) x_j^{(i)}$ for each minibatch: for each weight j: $w_j \leftarrow w_j - \Delta w_j$	for each epoch: i := rand training sample: $\Delta w_j = \alpha * (target^{(i)} - output^{(i)}) x_j^{(i)}$ for each weight j: $w_j \leftarrow w_j - \Delta w_j$
1) Start with some initial value for the parameter vector $\theta_0$ 2) Iteratively update the parameter vector by <ol style="list-style-type: none"> <li>Computing the gradient of the cost function <math>J(\theta)</math> at the last position reached (<math>\theta_t</math>): <math>\nabla_{\theta} J(\theta_t)</math></li> <li>Stepping in the negative direction of the gradient according to <math>\theta_{t+1} = \theta_t - \alpha * \nabla_{\theta} J(\theta_t)</math></li> </ol> 3) Stop when change in parameter vector small	$W, b = \text{initialize\_weights}(nx, ny)$ $X = \text{shuffle}(X)$ foreach(epoch in epochs) minibatches = get_next_minibatch(X) foreach(minibatch in minibatches) $w \leftarrow w - \alpha (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}$ $b \leftarrow b - \alpha (h_{\theta}(x^{(i)}) - y^{(i)})$	$W, b = \text{initialize\_weights}(nx, ny)$ $X = \text{shuffle}(X)$ foreach(epoch in epochs) foreach(x in X) $w \leftarrow w - \alpha (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}$ $b \leftarrow b - \alpha (h_{\theta}(x^{(i)}) - y^{(i)})$