

Negative log-likelihood loss:  $L = -\sum y_i \cdot \log \hat{y}_i$

Tomas B.

Sigmoid: 0-1 / like a neuron / get small gradients / not 0 centered

$$f(x) = \sigma(x) = 1/(1+e^{-x}) / f'(x) = f(x) \cdot (1-f(x))$$

Softmax: multiclass /  $\sum_i \sigma(x_i) = 1$  /  $f(x_i) = \exp(\eta_i)/\sum_j \exp(\eta_j)$

tanh: -1-1 / zero centered / steeper gradients

$$f(x) = \tanh(x) = (e^x - e^{-x})/(e^x + e^{-x}) / f'(x) = 1 - f(x)^2$$

ReLU: [0, inf] / can slow up / accelerate convergence / negative units get no update

$$f(x) = \max(0, x) / f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$$

Leaky ReLU: solves dying ReLU problem / or small positive constant

$$f(x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases} / f'(x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$$

Gradient Descent:  $\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta)$   $\rightarrow \theta = \theta - \eta \cdot \nabla_{\theta} J(\theta; x^{(i)}, y^{(i)})$

Stochastic Gradient Descent: Mini-batches / can jump / risk of overshooting

AdaGrad: Modifies learning rate for each parameter  $\theta_j$  based on past gradients

Max Likelihood:  $L(\theta) = \prod p(x_i; \theta) \Leftrightarrow \log L(\theta) = \log \prod p(x_i; \theta) \Leftrightarrow \log L(\theta) = \sum_i \log p(x_i; \theta)$

Feed forward Network Block:  $y = \sigma(w^T x + b)$  / Perceptron

$$\text{XOR}: \begin{matrix} x = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} & y = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} \end{matrix} \cdot \text{Linear: } J(\theta) = \frac{1}{N=4} \sum_{i=1}^4 (f(x^{(i)}) - f(x^{(i)}, \theta))^2 \quad f(x, w, b) = x^T w + b$$

Universal Approx Thm: Feed forward neural network with 1 hidden layer and continuous non-linear activation function can approximate any continuous function

With arbitrary precision / leading edge:  $-b/w$  je grüßer steiler /  $b/w$  je negativer desto rechter

Simple Cell: CNN / Small receptive field / at the beginning / tuned to specific stimuli

Complex Cell: combine output from simple cells / increase variance and receptive field

Convolution: linear:  $J(\alpha u + \beta v) = \alpha J(u) + \beta J(v)$  / Invariant:  $J(f(u)) = J(u)$  / equivariant  $J(f(u)) = f(J(u))$

Correlation:  $I''(i, j) = \sum_{m=-k}^k \sum_{n=-k}^k K(m, n) \cdot I(i+m, j+n)$   $\{K(i, j) = K(-i, -j) \rightarrow \text{Correlation} \equiv \text{convolution}$

Convolution:  $I''(i, j) = \sum_{m=-k}^k \sum_{n=-k}^k K(-m, -n) \cdot I(i+m, j+n)$

Semantic Segmentation: Classify each pixel / Downsampling (max-pool / strided) then Upsample again

Upsampling: Unpadding (reverse convolution) / Max Unpooling (remember which element max was)

Learnable Upsampling (learn transpose convolution filter)  $\xrightarrow{\text{Hidden state}}$

RNN:  $h^+ = \tanh(h \cdot W_{hn} \cdot h^{t-1} + W_{xh} \cdot x^+) / \hat{y}^+ = W_{ny} \cdot h^+$   $\xrightarrow{x} \xrightarrow{h} \xrightarrow{\oplus} \xrightarrow{h^+} \xrightarrow{\oplus} \xrightarrow{y}$

$\rightarrow$  1:1 = Vanilla / 1:m = Image Captioning / m:1 = Sentiment Classification

$\rightarrow$  m:n = Machine Translation / n:n = Video classification for each frame

Parameters shared between time steps / impossible to capture long-term dependencies

gradient multiplication for long sequences: vanish  $\Rightarrow$  LSTM / explodes  $\Rightarrow$  gradient clipping

LSTM: forget gate: what forgotten from previous cell state  $c_{t-1} \xrightarrow{\oplus} c_t$

input gate: how much input we should take

output gate: how much memory affects output

$$f_t = \sigma(W_f \cdot x_t + U_f \cdot h_{t-1} + b_f)$$

$$i_t = \sigma(W_i \cdot x_t + U_i \cdot h_{t-1} + b_i)$$

$$o_t = \sigma(W_o \cdot x_t + U_o \cdot h_{t-1} + b_o)$$

$$c_t = f_t \cdot c_{t-1} + i_t \cdot o_t \cdot \sigma(c_t) / h_t = o_t \cdot \sigma_h(c_t)$$

Cell state is memory / hidden state is for output

GRU: cell state stored in hidden state / lacks output gate

because memory already in hidden state

Regularization: Data standardization for better training convergence & local minima avoidance

$$\rightarrow x_s = \frac{x - \mu}{\sigma} / \mu = \frac{1}{n} \cdot \sum x_i / \sigma = \sqrt{1/n-1} \cdot \sqrt{\sum (x_i - \mu)^2}$$

Data augmentation: Real/synthetic discrepancies

Parameter norm penalties:  $L(\theta) = L(\theta) + \alpha \cdot \|L\|_1$  L1: Lasso / sparse

Dropout: Models share weights / trains part of model L2: ridge / shrinkage

Bagging: Models are independent / train all models until convergence

Pretraining with general big dataset or pretrained model

Autoencoder: Encoder encodes  $x$  into latent space  $\tilde{x} = f(x)$  / Decoder decodes  $\tilde{x}$  into  $\hat{x} = g(\tilde{x})$

$\cdot g \circ f$  approximates identity /  $\hat{f}, \hat{g} = \arg \min_{f, g} \sum_{i=1}^N \|x_i - \hat{x}_i\|^2$

$\cdot$  If  $f$  and  $g$  linear  $\Rightarrow$  optimal  $\left\{ \hat{f}, \hat{g} = \arg \min_{f, g} \sum_{i=1}^N \|x_i - g(f(x_i))\|^2 \right.$

solution given by PCA

Kullback-Leibler divergence:  $D_{KL}(P||Q) = \sum_x p(x) \cdot \log \frac{p(x)}{q(x)} \cdot dx$  (How wide apart are dist.  $p(x)$  and  $q(x)$ )

Batch normalization: distribution shifts at each training iteration  $\Rightarrow$  weights need to adjust constantly

$\rightarrow$  standardize each layer distribution / hard to compare models since  $\log \mathcal{E}$

$$\rightarrow E_z[\log p_\theta(x^{(i)}|z)] - D_{KL}(q_\theta(z|x^{(i)})||p_\theta(z))$$

Blurry images due to injecting noise and inform model

By introducing density model (Gaussian) over latent space  $z \sim f(x) \sim N(\mu, \sigma^2)$  where  $\mu$  and  $\sigma^2$  are estimated on training data  $\Rightarrow$  blurry images because  $z$  discrete  
Variational Autoencoder: To get  $x$ : sample from  $p_\theta(x|z^{(i)}) / p(z)$  is simple prior we parametrize  $p(x|z)$  by neural network since  $\text{complex } p(x) = \int p(x|z) \cdot p(z) dz$   
 Use  $p(z|x) = p(x|z) \cdot p(z) / p(x)$   $p(z|x)$  is approximated by encoder network

B-VAE: Learn disentangled representation without supervision / provide framework for automated discovery of interpretable factorised latent representation / modify VAE, add adjustable hyperparameter  $\beta$  balances latent channel capacity and constraints (GAN): Generator  $G: \mathbb{R}^D \rightarrow x$  and Discriminator  $D: x \rightarrow [0, 1]$

$$\min_{\theta} \max_{\phi} \mathbb{E}_{x \sim p(x)} [\log D_\phi(x)] + \mathbb{E}_{z \sim p_\theta(z)} [\log (1 - D_\phi(G_\theta(z)))]$$

1. Gradient ascend on  $D$ :  $\max_{\phi} \mathbb{E}_{x \sim p(x)} [\log D_\phi(x)] + \mathbb{E}_{z \sim p_\theta(z)} [\log (1 - D_\phi(G_\theta(z)))]$

2. Gradient ascend on  $G$ :  $\max_{\theta} \mathbb{E}_{z \sim p_\theta(z)} [\log D_\phi(G_\theta(z))]$  sample generated

K-steps do sample  $x_i, x_G$ , update discriminator  $D_\phi$   $\sum_{i=1}^m [\log D_\phi(x^{(i)}) + \log (1 - D_\phi(G_\theta(z^{(i)})))]$   
 do one update generator  $\log \prod_{i=1}^m \log(D_\phi(G_\theta(z^{(i)})))$  / Nash-Equilibrium - 2 player-game  
 Mode collapse: generator collapses and produces limited varieties of samples due focus on one mode

GAN < VAE: sharper images / GANs only care about samples / No variational bound needed

Autoregressive property: linear reg.  $\hat{y} = b_0 + \alpha_i x_i$  / autoregr. used for time series

Tadular approach: chain rule probabilities:  $p(x) = \prod p(x_i|x_1, \dots, x_{i-1}) = \prod p(x_i|x_0)$

Fully visible belief network:  $\hat{x}_i = p(x_i = 1 | x_{\neq i}) \Rightarrow$  logistic reg.:  $f(x_1, \dots, x_{i-1}) = \sigma(x_0 + \alpha_i x_i + \dots)$

NADE: Neural Autoregressive Density Estimator / Autoencoder-like neural network

Max. average log-likelihood:  $\frac{1}{T} \sum_i \log(p(x^i)) = \frac{1}{T} \sum_{i=1}^T \sum_{j=1}^D \log(p(x_j^i | x_{\neq j}^i))$

Efficient O(TD) / could make use of second order optimizers

MADE: Masked Autoencoder Distributor or Estimator / constrain Autoencoder that output can be used as conditionals  $p(x_i|x_{\neq i})$  / No computation path between  $x_d$  and input  $x_1, \dots, x_D$  must exist / Training same complexity as Autoencoder very large hidden layers necessary / computing  $p(x)$  just forward pass

Pixel RNN:  $p(x) = \prod p(x_i|x_1, \dots, x_{i-1})$  (likelihood of  $i$ -th pixel given all previous pixels)  
 Sequential generation is slow due to pixel dependencies

Pixel CNN: Dependencies on previous pixels as CNN over context regions

Reinforced learning: States  $S$  / Actions  $A$  / reward  $r(S, A) \rightarrow R$  / transition  $p(S, A) \rightarrow S'$  / init  $s_0$

Value function of policy  $\pi$ :  $V_\pi: S \rightarrow R$ , expresses expected cumulative reward

Bellman equation:  $G_{t+1} = \sum_a \pi(a|s_t) \cdot \sum_s \sum_r p(s_t, r | s_t, a) \cdot (r + \gamma \cdot \mathbb{E}_{s_{t+1}}(G_{t+1} | s_{t+1} = s'))$

$\hookrightarrow$  Prob. of taking an action given state  $\hookrightarrow$  Joint prob. of next state and reward

\* expected reward given a new state given current state and action (transition matrix)

Dynamic Programming: Value iter: 1. Compute optimal state value function  $V_\pi$   
 2. Obtain policy via back tracing

$r(s, a)$ : reward or cost Policy iter: 1. Compute state value function  $V_\pi$  for arbitrary policy  $\pi$   
 $V_\pi(p(s, a))$ : set value in neighborhood 2. Update policy  $\pi$  given  $V_\pi \rightarrow \pi'$   $\Rightarrow$  Iterate till  $\pi = \pi'$   
 $\hookrightarrow$  Update step for dyn. prog. step:  $\pi'(s) = \arg\max_a (r(s, a) + \gamma \cdot V_\pi(p(s, a)))$

+ exact, guaranteed to converge / -: know transition prob. matrix, iterate over whole state space

Temporal difference learning: If state space is too big / Random policy: exploration

greedy policy: exploitation, can get stuck in local minima ( $\Delta V(s) = r(s, a) + \gamma V(s') - V(s)$ )

Update only visited state spaces, use  $\epsilon$ -greedy-strategy /  $V(s) \leftarrow V(s) + \alpha \cdot \Delta V(s) \quad \alpha: \text{Learning rate}$

Q-Learning:  $\Delta Q(s, A) = R_{t+1} + \gamma \cdot \max_a Q(s', a) - Q(s, A) \quad Q(s, A) \leftarrow Q(s, A) + \alpha \cdot \Delta Q(s, A)$

Off-policy: policy to update Q-function different from policy we use to collect data (=  $\epsilon$ -greedy)

Tadular Learning: Learn each state separately / no generalization  $\Rightarrow$  Use function approx. to learn value func.

Deep Reinforced Learning: Use neural network to learn mapping state-action pairs  $(s, a)$  and value

Loss( $\theta$ ) =  $(R + \gamma \max_a Q_\theta(s', a) - Q_\theta(s, A))^2$  / Q-Networks: samples need i.i.d.  $\rightarrow$  Replay Buffer =  $S, a, r, s'$

Actor critic: Builds on policy gradient methods. Instead of using rewards, a value function from NN is used as an estimate

$\frac{\partial L}{\partial W} = \sum_{t=1}^T \frac{\partial L^t}{\partial W}$  On policy: compute Q-value according to policy and agent follows policy

$\frac{\partial L^t}{\partial W} = \sum_{k=1}^K \frac{\partial L^t}{\partial y^t} \cdot \frac{\partial y^t}{\partial h^t} \cdot \frac{\partial h^t}{\partial W}$  Off-policy: Q-value greedy but agent follows different exploration policy

$$\frac{\partial \tanh(x)}{\partial x} = (1 - \tanh^2(x)) \cdot x' \quad \frac{\partial \sigma(x)}{\partial x} = (\sigma(x) \cdot (1 - \sigma(x))) \cdot x'$$

$$(\log x)' = 1/x$$