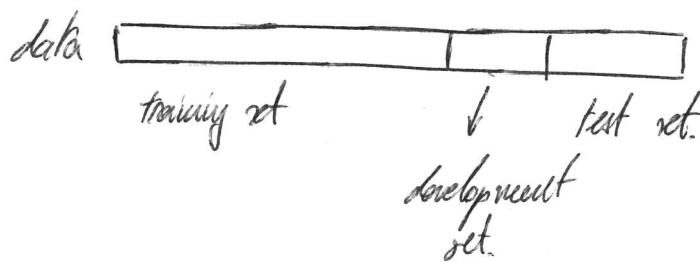


Hyperparameter training, regularization and optimisation

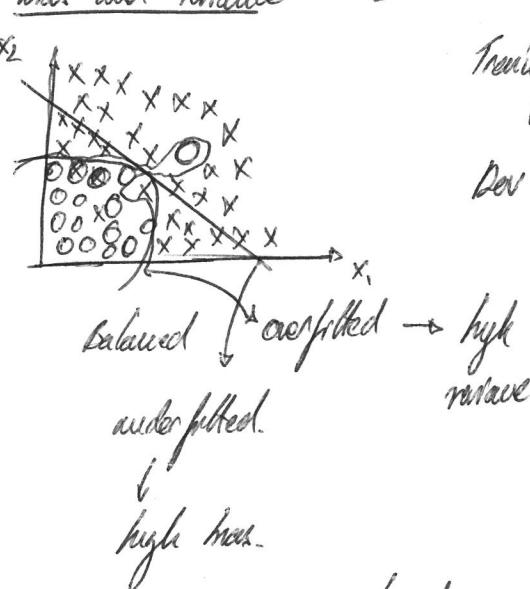
* Train, dev, test →



development set = shows which model performs best.

98% / 11% / 11%.

* Bias and variance →



Training set 1% | 15% | 15% | 0.5%

Dev set 11% | 16% | 30% | 1%

Overfitting high bias
high variance underfitting both low bias and variance.

* This is based on optimal error 10%, but if optimal error is 15%
(2) is a very good model.

A. High variance →

1. % error from 10% is large enough to optimal.

2. Improve →

A. Better Network.

B. Train longer

C. NN arch set.

B. High variance \rightarrow large % error dev compared to training.

1. Improvements \rightarrow

A. More data.

B. Regularization.

C. NN arch ret.

* Regularization \rightarrow

- prevent overfitting

$$J(w, b) = \frac{1}{M} \sum_{i=1}^M h(\hat{y}_{(i)}, y_{(i)}) + \frac{\lambda}{2M} \|w\|_2^2$$

$$\text{l}_2 \text{ regularization} \rightarrow \|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w$$

$$\text{l}_1 \text{ regularization} \rightarrow \frac{1}{M} \sum_{i=1}^M |w_i| = \frac{\lambda}{M} \|w\|_1$$

λ = Regularization parameter.

- Neural Network \rightarrow

$$J(w_{0:1}, b_{0:1}, \dots, w_{l:1}, b_{l:1}) = \frac{1}{M} \sum_{i=1}^M h(\hat{y}_{(i)}, y_{(i)}) + \frac{\lambda}{2M} \sum_{l=1}^L \|w_{l:1}\|_2^2$$

$$\|w_{l:1}\|_F^2 = \|w_{l:1}\|_2^2 = \sum_{i=1}^{n_{cl-1,l}} \sum_{j=1}^{n_{cl,l}} (w_{ij,l})^2 ; \quad w : (n_{cl,l}, n_{cl-1,l})$$

* gradient descent \rightarrow

$$-\underline{\lambda} \underline{dw}$$

$$w_{0:1} = w_{0:1} - \frac{\lambda}{M} w_{0:1} - \alpha [\text{back propagation}]$$

\swarrow weight decade.

* Vanishing / Exploding gradients →

- In very deep networks → slopes can get really high or low.
- Activation values → increase / decrease exponentially.

* Weight initialization for Deep Networks →

- Need to balance the $z = w_1x_1 + \dots + w_nx_n$ so the values are not very high, low. → control weight with variance.

$$\text{Var}(w_i) = \frac{1}{n} \rightarrow \text{w}_{i, \text{init}} = \text{Random number} * \text{np.sqrt}\left(\frac{1}{n \cdot \text{ncol} - 1}\right)$$

- ReLU → $\sqrt{\frac{2}{\text{ncol} - 1}}$

- Other probability →

- Tanh → $\sqrt{\frac{1}{\text{ncol} - 1}}$

$$\sqrt{\frac{2}{\text{ncol} * \text{ncol} - 1}}$$

* Gradient checking →

$$w_{0,1}, b_{0,1}, \dots, w_{L,1}, b_{L,1} \rightarrow \theta \Rightarrow f(w_{0,1}, b_{0,1}, \dots, w_{L,1}, b_{L,1}) = f(\theta)$$

$$f(\theta_1, \theta_2, \dots, \theta_L) \rightarrow$$

- for each layer →

$$\Delta \theta_{\text{approx}, i, l} = \frac{f(\theta_1, \theta_2, \dots, \theta_l + \epsilon, \dots) - f(\theta_1, \theta_2, \dots, \theta_l - \epsilon, \dots)}{2\epsilon}$$

$$\Delta \theta_{\text{exact}, i, l} = \frac{\partial f}{\partial \theta_l}$$

- Check if $\Delta \theta_{\text{approx}, i, l} \approx \Delta \theta_{\text{exact}, i, l}$? → Euclidean distance

$$\|\Delta \theta_{\text{approx}, i, l} - \Delta \theta_{\text{exact}, i, l}\|_2 / (\|\Delta \theta_{\text{approx}, i, l}\|_2 + \|\Delta \theta_{\text{exact}, i, l}\|_2) \leq \epsilon = 10^{-7}$$

1. Do not use in training \rightarrow overhead mainly.
2. If it fails \rightarrow check components on $i \rightarrow$ doc, do approx 0.7
3. Keep in mind regularization term \rightarrow

$$J(\theta) = \frac{1}{M} \sum_{i=1}^M h(\tilde{y}(i), y(i)) + \frac{\lambda}{2M} \sum_{j=1}^J \|w_j\|_F^2$$

4. Don't work with dropout.
5. Run at random initialization \rightarrow to check that it not only works for $w, b \geq 0$.

* Optimization Algorithms \rightarrow

- mini-batch gradient descent \rightarrow

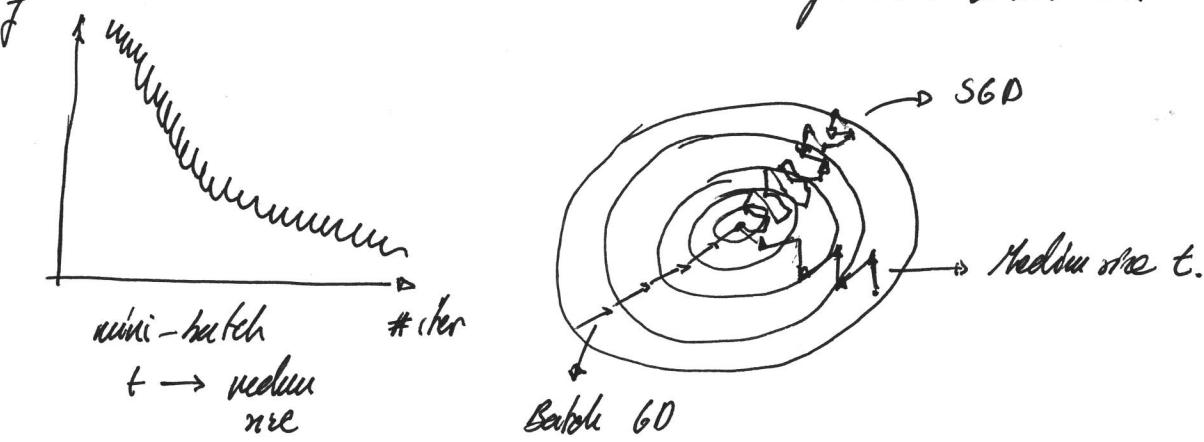
$$\mathbf{x} = \begin{bmatrix} \underbrace{x(1)}_{x^{(1)}} & \underbrace{x(2)}_{x^{(2)}} & \dots & \underbrace{x(M)}_{x^{(M)}} \end{bmatrix}$$

t -thre

* 1 epoch \equiv Pass through training set m .

$t=1 \rightarrow$ stochastic GD \equiv loses the advantage of regularization.

$t=m \rightarrow$ Batch GD \equiv gradient descent not updated until m .



* Mini-batches size $\rightarrow 2^{nk}$

- * The λ will control weight for each layer, making them fall within small values of z ; $a = g(z) = \tanh(z) \rightarrow$ linear regions.

- * λ = hyperparameter to be set; takes all medium quadratic values. A large value will control the gradient for weight update $\rightarrow w_t$

$$J(w_{CLS}, b_{CLS}) = \frac{1}{M} \sum_{i=1}^M \ell(\hat{y}_{(i)}, y_{(i)}) + \frac{\lambda}{2M} \sum_{l=1}^L \|w_l\|_F^2$$

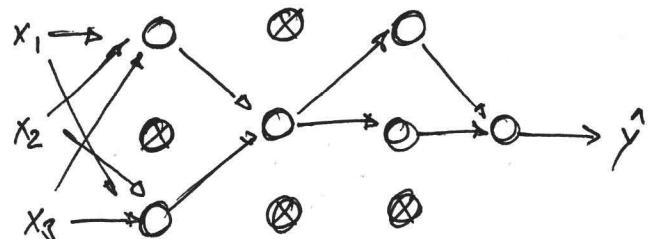
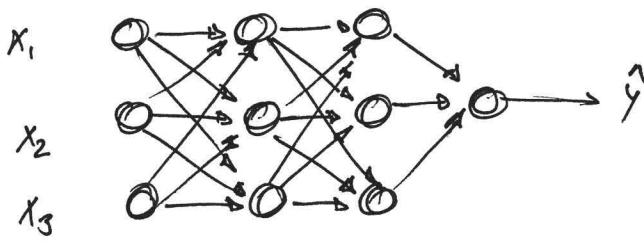
$$\|w_F\|^2 = \sum_{i=1}^{n_{CL-1}} \sum_{j=1}^{n_{CLS}} (w_{ij})^2$$

- * Large weights give too much importance to certain feature, we are looking for general features, not to stand out outliers.

$$w_{CLS} = w_{CLS} \underbrace{\left(1 - \alpha \frac{\lambda}{M}\right)}_{0.0} - \alpha [\text{backpropagation}]$$

- * Dropout \rightarrow

Each iteration, network changes.



- During training each neuron has a dropout probability on the network
- Only done in training; not in test.

- * Direct training cost function will change with each iteration.

- * Implementation →

- Forward propagation →

each layer [C]

$dl = np.random.rand(al.shape[0], al.shape[1]) < \text{keep-prob.}$

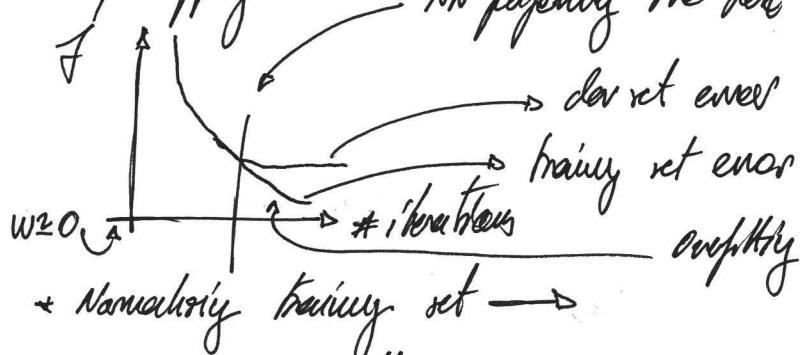
$al = up.multiply(al, dl)$

$al = al / \text{keep-prob}$ → since some neurons have been eliminated from network, give more importance to other values → balance outputs.

- Backward propagation →

$dA[C] = dA[C]/\text{keep-prob.}$

- * Early stopping → NN停止到的点



- * Normalizing training set →

$$\mu = \frac{1}{M} \sum_{i=1}^M x(i) \rightarrow x = x - \mu$$

$$\sigma^2 = \frac{1}{M} \sum_{i=1}^M x(i)^2 \rightarrow x = x / \sigma^2$$

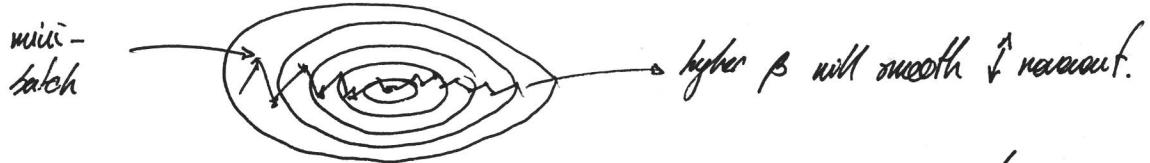
- Takes out the fact that the next are orthogonal.

1. Optimizing cost function J
2. Not oneflection.

- Centers and distributes points based on variance →

makes it faster to reduce cost function

- Gradient descent with momentum →



- * Idea → average out vertical axis, and move faster on axis.
 $\beta \geq 0.9$ default.

$$Vdw = \beta Vdw + (1-\beta) dw$$

$$Vdb = \beta Vdb + (1-\beta) db$$

$$Vdw^{\text{corrected}} = Vdw / (1 - \beta^t)$$

$$Vdb^{\text{corrected}} = Vdb / (1 - \beta^t)$$

$$\begin{aligned} w &= w - \alpha Vdw^{\text{corrected}} \\ b &= b - \alpha Vdb^{\text{corrected}} \end{aligned}$$

- RMS prop →

- * base idea → stochastic, speed-up w small

$$Sdw = \rho Sdw + (1-\rho) dw^2 ; dw^2 = \text{element wise.}$$

$$Sdb = \rho Sdb + (1-\rho) db^2$$

$$w = w - \alpha \frac{dw}{\sqrt{Sdw}}$$

$$b = b - \alpha \frac{db}{\sqrt{Sdb}}$$

- ADAM optimization algorithm \rightarrow

* Adam = Adaptive Moment Estimation.

* Momentum + RMS prop

$$V_{dw} = 0; \quad S_{dw} = 0; \quad V_{db} = 0; \quad S_{db} = 0$$

On every iteration t (mini-batch) \rightarrow

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dw$$

$$V_{db} = \beta_1 V_{db} + (1 - \beta_1) db$$

$$V_{dw}^{\text{conv}} = V_{dw} / (1 - \beta_1^t)$$

$$V_{db}^{\text{conv}} = V_{db} / (1 - \beta_1^t)$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dw^2$$

$$S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2$$

$$S_{dw}^{\text{conv}} = S_{dw} / (1 - \beta_2^t)$$

$$S_{db}^{\text{conv}} = S_{db} / (1 - \beta_2^t)$$

$$w = w - \alpha \left(\frac{V_{dw}^{\text{conv}}}{\sqrt{S_{dw}^{\text{conv}}} + \epsilon} \right)$$

$$b = b - \alpha \left(\frac{V_{db}^{\text{conv}}}{\sqrt{S_{db}^{\text{conv}}} + \epsilon} \right)$$

$$\text{default parameters} \rightarrow \beta_1 = 0.9; \quad \beta_2 = 0.999; \quad \epsilon = 10^{-8}$$

- learning rate decay \rightarrow

- * Homogeneous learning decay with time \rightarrow variation when you reach case of minimum cost is small.

$$\frac{|x_1(t)| \times |x_2(t)| \dots |x_m(t)|}{\text{epoch}} \quad t = \text{batch size.}$$

* Implementation \rightarrow

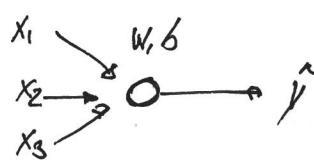
$$\alpha = \frac{1}{1 + \text{decay-rate} * \text{epoch \#}} \quad \alpha_0$$

$$\alpha = 0.95^{\text{epoch \#}} \quad \alpha_0$$

$$\alpha = \frac{K}{\text{epoch \#}} \quad \alpha_0 \quad \text{or} \quad K = \frac{K}{\text{epoch \#}} \quad \alpha_0$$

* Hyperparameter tuning, \rightarrow

- batch normalization \rightarrow



$$\mu = \frac{1}{M} \sum_i x(i) \quad ; \quad \sigma^2 = \frac{1}{M} \sum_i x(i)^2 - \mu^2 \quad \text{over all samples}$$
$$x = x - \mu \quad ; \quad x = x / \sigma^2 \quad \text{den. var.}$$

* layer l \rightarrow

$$\mu = \frac{1}{M} \sum_i z(i) \quad ; \quad \sigma^2 = \frac{1}{M} \sum_i (z(i) - \mu)^2 \quad ; \quad \hat{z}_{\text{norm}}^{(l)} = \frac{z(i) - \mu}{\sqrt{\sigma^2 + \epsilon}} \quad \text{Avoid } 0/0$$

$$\hat{z}_{\text{cont}}^{(l)} = \hat{z}_{\text{norm}}^{(l)} \cdot r + \beta \quad \text{learnable param are needed.}$$

* r, β allows you to move the distribution

* After iteration $\rightarrow d\beta^{(l)} \rightarrow \beta^{(l+1)} = \beta^{(l)} - \alpha d\beta^{(l)}$

* Normalizing with mini-batches \rightarrow

$$x_{itj} \rightarrow z_{itj} \xrightarrow[\text{wces, bces}]{BN} \hat{z}_{itj} \rightarrow g_{itj}(\hat{z}_{itj})$$

\hat{z}_{itj} \rightarrow Normalized, scaled out.

$$z_{itj} = w_{ces} \cdot a_{itj-1} + b_{itj}$$

$$z_{itj} = w_{ces} \cdot a_{itj-1}$$

$$z_{itj, \text{norm}} \rightarrow \hat{z}_{itj} = g_{itj}(\text{norm}(z)) + \beta_{itj}; \quad \beta \in (\text{nces}, 1)$$

wces, dpces, drtes

update parcell.

↑
Important

* Parcells \rightarrow

1. Can be represented by momentum, RLSprop, Adams.
2. Batch norm reduces the variance of activation function in deeper networks \rightarrow as you center in distribute data. \rightarrow More stable.

* Regularization effect \rightarrow

1. Many variables adds some noise, depending on me.
2. Relation between μ, τ^2 .
3. Larger size, less noise.

* Test time \rightarrow

- the weighted average across minibatches

$$x_{1tj}, x_{2tj}, \dots, x_{ntj} \rightarrow \mu_{1tj}, \sigma_{1tj}, \dots, \mu_{ntj}, \sigma_{ntj}$$

$$\hat{z}_{itj, \text{norm}} = \frac{z_{itj} - \mu}{\sqrt{\tau^2 + C}} \rightarrow \hat{z} = \gamma z_{itj, \text{norm}} + \beta$$

- softmax \rightarrow

$$z = W \cdot \text{prev} + B$$

$$t = e^z$$

$$\rightarrow a_i[\text{CL}] = \frac{t_i[\text{CL}]}{\sum_{j=1}^C t_j}; \quad \text{CLS} = \frac{t[\text{CLS}]}{\sum_{j=1}^C t_j[\text{CLS}]}$$

* Distributes probabilities between 0,1 for all possible classes in vector.

* Numerical stability \rightarrow

1. Exponentials can explode values.

$$\log C = -\max_i f(y_i)$$

2. Division causes stability.

Shift values to highest digit 0.

$$A = \frac{e^{f(y_i) + \log C}}{\sum_{j=1}^C e^{f(y_j) + \log C}}$$

* Use cross-entropy \rightarrow

$$h(y, \hat{y}) = - \sum_{i=1}^M y_i \log \hat{y}_i; \quad J(w, b) = \frac{1}{M} \sum_{i=1}^M h(\hat{y}(w), y(i))$$

$$\frac{\partial J(y, \hat{y})}{\partial \text{CLS}} = \hat{y} - y$$

- batch normalization points →

1. Improves gradient flow through network.
2. Allows higher learning rates.
3. Reduces dependency on initializations.
4. Acts as a form of weight normalization.

Appendix

- Hinge loss function →

- * the for multiclass support vector machine SVM.

$s \equiv$ score function $f(x_i; w)$

$$l = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + \Delta)$$

- * Expression that sums up all incorrect classes $j \neq y_i$, keeping only larger than one.

- * SVM wants the score of the correct class y_i to be larger than incorrect classes by at least Δ .

- * Squared hinge loss → $\text{max}^2(0, \dots)$ L2-SVM

- Cross entropy function →

$$H(p, q) = - \sum_x p(x) \cdot \log q(x)$$

$p \equiv$ true distribution.

$q \equiv$ estimated distribution.

def cross-entropy (x, y) :

$$x \rightarrow m, n_x$$

$$y \rightarrow m, 1$$

$$m = y \cdot \text{shape}(x)$$

$$p = \text{softmax}(x)$$

$$\text{log-likelihood} = -\text{np. log}(p[\text{range}(m), y])$$

$$\text{loss} = \text{np. sum(log-likelihood)} / m$$

* softmax derivative \rightarrow

$$A_{a,b} = \frac{e^{z_{ab}}}{\sum_{k=1}^{n_k} e^{z_{kb}}} ; \quad h(x) = f(x)/g(x) \rightarrow h'(x) = \frac{f'(x)g(x) - f(x)g'(x)}{g^2(x)}$$

$$\frac{\partial A_{ab}}{\partial z_{ij}} = \frac{(e^{z_{ij}})' \sum_{k=1}^{n_k} e^{z_{kj}} - e^{z_{ij}} \left(\sum_{k=1}^{n_k} e^{z_{kj}} \right)'}{\left(\sum_{k=1}^{n_k} e^{z_{kj}} \right)^2} = [\text{if } j \neq b \text{ then } 0.] =$$

$$\frac{\partial A_{ab}}{\partial z_{ib}} = \frac{(e^{z_{ab}})' \sum_{k=1}^{n_k} e^{z_{kb}} - e^{z_{ab}} \left(\sum_{k=1}^{n_k} e^{z_{kb}} \right)'}{\left(\sum_{k=1}^{n_k} e^{z_{kb}} \right)^2} =$$

$$\Rightarrow [i=a] = \frac{e^{z_{ab}}}{\left(\sum_{k=1}^{n_k} e^{z_{kb}} \right)^2} - \frac{e^{z_{ab}} \cdot e^{z_{ab}}}{\left(\sum_{k=1}^{n_k} e^{z_{kb}} \right)^2} = A_{a,b} - A_{ab}^2$$

$$\Rightarrow [i \neq a] = 0 - \frac{e^{z_{ib}} e^{z_{ib}}}{\left(\sum_{k=1}^{n_k} e^{z_{kb}} \right)^2} = - A_{ib} \cdot A_{ab}$$

$$\frac{\partial A_{a,b}}{\partial z_{ij}} = [\delta_{i,a} A_{ab} (1 - A_{ab}) - A_{ib} \cdot A_{ab} (1 - \delta_{ia})] \delta_{bi,j}$$

Cross entropy gradient \rightarrow

$$\hat{J}_{pb} = \frac{1}{m} \sum_{b=1}^m \left(- \sum_{p=1}^{n_y} y_{pb} \cdot \log \hat{y}_{pb} \right)$$

$$g = \log(\hat{y}) \rightarrow \frac{\partial g_{ke}}{\partial \hat{y}_{ij}} = \frac{1}{\hat{y}_{ke}} \delta_{k,i} \delta_{e,j}$$

$$\frac{\partial h}{\partial \hat{y}_{ij}} = \frac{1}{m} (-y_{ij} \cdot 1) = -\frac{1}{m} y_{ij}$$

$$\frac{\partial h}{\partial \hat{y}_{ij}} = \sum_{k=1}^{n_y} \sum_{e=1}^m \frac{\partial h}{\partial y_{ke}} \cdot \frac{\partial y_{ke}}{\partial \hat{y}_{ij}} = \sum_{k=1}^{n_y} \sum_{e=1}^m \left(-\frac{1}{m} y_{ke} \right) \cdot \frac{1}{\hat{y}_{ke}} \delta_{k,i} \delta_{e,j} =$$

$$= -\frac{1}{m} \frac{y_{ij}}{\hat{y}_{ij}}$$

- softmax \rightarrow

$$\frac{\partial h}{\partial z_{rj}} = \sum_{k=1}^{n_y} \sum_{e=1}^m \frac{\partial h}{\partial \hat{y}_{ke}} \cdot \frac{\partial \hat{y}_{ke}}{\partial z_{rj}} = \sum_{k=1}^{n_y} \sum_{e=1}^m \left(-\frac{1}{m} \frac{y_{ke}}{\hat{y}_{ke}} \right) \delta_{e,j} \left[\delta_{k,i} \hat{y}_{ke} (1 - \hat{y}_{ke}) - \hat{y}_{ij} \hat{y}_{ke} \right]$$

$$(1 - \delta_{k,i})] \delta_{e,j} = \sum_{k=1}^{n_y} \left(-\frac{1}{m} y_{kj} \right) \left[\delta_{k,i} (1 - \hat{y}_{kj}) - \hat{y}_{ij} (1 - \delta_{k,i}) \right]$$

$$= \left(-\frac{1}{m} y_{ij} \right) (1 - \hat{y}_{ij}) + \sum_{k=1}^{n_y} \left(-\frac{1}{m} y_{kj} \right) \hat{y}_{ij} - \left(\cancel{\frac{1}{m} y_{ij} \cdot \hat{y}_{ij}} \right) =$$

$$= -\frac{1}{m} y_{ij} + \frac{1}{m} \hat{y}_{ij} \underbrace{\sum_{k=1}^{n_y} y_{kj}}_{=1} = \frac{1}{m} (\hat{y}_{ij} - y_{ij})$$

softmax gradient →

$$A_{a,b} = \frac{e^{z_{ab}}}{\sum_{k=1}^{n_k} e^{z_{kb}}} ; \quad h(x) = g(x)/f(x) \rightarrow h'(x) = \frac{g''(x)f(x) - g(x)f'(x)}{f^2(x)}$$

$$\begin{aligned} \frac{\partial A_{a,b}}{\partial z_{ij}} &= \frac{(e^{z_{ab}})' \sum_{k=1}^{n_k} e^{z_{kb}} - e^{z_{ab}} \left(\sum_{k=1}^{n_k} e^{z_{kb}} \right)'}{\left(\sum_{k=1}^{n_k} e^{z_{kb}} \right)^2} = \begin{cases} 0 & \text{if } j \neq b \\ 1 & \text{if } j = b \end{cases} = \\ &= \frac{(e^{z_{ab}})' \sum_{k=1}^{n_k} e^{z_{kb}} - e^{z_{ab}} \left(\sum_{k=1}^{n_k} e^{z_{kb}} \right)'}{\left(\sum_{k=1}^{n_k} e^{z_{kb}} \right)^2} \delta_{j,b} = \end{aligned}$$

$$\underline{i=a} \Rightarrow \frac{e^{z_{ab}}}{\sum_{k=1}^{n_k} e^{z_{kb}}} - \frac{e^{z_{ab}} \cdot e^{z_{ab}}}{\left(\sum_{k=1}^{n_k} e^{z_{kb}} \right)^2} = A_{ab} - A_{ab}^2$$

$$\underline{i \neq a} \Rightarrow 0 - \frac{e^{z_{ab}} e^{z_{ib}}}{\left(\sum_{k=1}^{n_k} e^{z_{kb}} \right)^2} = - A_{ab} A_{ib}$$

$$\frac{\partial A_{ab}}{\partial z_{ij}} = [(A_{ab} - A_{ab}^2) \delta_{i,a} - A_{ab} A_{ib} (1 - \delta_{i,a})] \delta_{j,b}$$

Batch Normalization →

samples
↓
 $x \rightarrow (m, nx)$
 $y \rightarrow (m, ny)$

$$\mu_e = \frac{1}{N} \sum_p x_{pe} \cdot \frac{1}{N} ; \quad \sigma^2 = \frac{1}{N} \sum_p (x_{pe} - \mu_e)^2$$

$$\hat{x} = \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}} \rightarrow \hat{x}_{ke} = \frac{x_{ke} - \mu_e}{\sqrt{\sigma_e^2 + \epsilon}}$$

$$y = \gamma \hat{x} + \beta \rightarrow \begin{pmatrix} y_{11} & y_{12} & \dots \end{pmatrix} = \begin{pmatrix} \gamma_1 \hat{x}_{11} + \beta_1 & \gamma_2 \hat{x}_{12} + \beta_2 & \dots \end{pmatrix}$$

$$y_{ke} = \gamma_e \hat{x}_{ke} + \beta_e \rightarrow \frac{\partial y_{ke}}{\partial x_{ij}} = \gamma_e (\delta_{ik} \delta_{jk}) \rightarrow \text{Element wise.}$$

$$\frac{\partial h}{\partial x_{ij}} = \sum_{k=1}^M \sum_{l=1}^{n_y} \frac{\partial h}{\partial y_{ke}} \cdot \frac{\partial y_{ke}}{\partial x_{ij}} = \sum_{k=1}^M \sum_{l=1}^{n_y} \frac{\partial h}{\partial y_{ke}} \cdot \frac{\partial y_{ke}}{\partial \hat{x}_{ke}} \cdot \frac{\partial \hat{x}_{ke}}{\partial x_{ij}}$$

* Just translation →

$$\hat{x}_{ke} = x_{ke} - \mu_e = x_{ke} - \frac{1}{N} \sum_p x_{pe}$$

$$\frac{\partial \hat{x}_{ke}}{\partial x_{ij}} = \begin{cases} 1 - \frac{1}{N} & \text{if } k=i, l=j \\ -1/N & \text{if } k \neq i, l \neq j \\ 0 & \text{otherwise} \end{cases} = t_{ij} (\delta_{ik} - \frac{1}{N})$$

* Variance \rightarrow

$$\sigma_e^2 = \sum_p^M (x_{pl} - \mu_e)^2 \cdot \frac{1}{N}$$

$$\begin{aligned} \frac{\partial \sigma_e^2}{\partial x_{ij}} &= \frac{1}{N} \cdot 2 \sum_p^M \frac{\partial(x_{pl} - \mu_e)}{\partial x_{ij}} \cdot (x_{pl} - \mu_e) = \frac{2 s_{e,j}}{N} \sum_p^M (x_{pl} - \mu_e) (s_{p,i} - \frac{1}{N}) \\ &= \frac{2 s_{e,j}}{N} \left[(x_{ie} - \mu_e) + \underbrace{\frac{1}{N} \sum_p^M \mu_e - \frac{1}{N} \sum_p^M x_{pe}}_{\frac{N}{N} \mu_e - \mu_e} \right] = \\ &= \frac{2}{N} s_{e,j} (x_{ie} - \mu_e) \end{aligned}$$

* Norm \rightarrow

$$\begin{aligned} \frac{\partial \hat{x}_{ke}}{\partial x_{ij}} &= \frac{\partial(x_{ke} - \mu_e)}{\partial x_{ij}} (\tau_e^2 + \epsilon)^{-1/2} - \frac{1}{2} \frac{\partial \sigma_e^2}{\partial x_{ij}} (x_{ke} - \mu_e) (\tau_e^2 + \epsilon)^{-3/2} \\ &= s_{k,j} \left(s_{k,i} - \frac{1}{N} \right) (\tau_e^2 + \epsilon)^{1/2} - \frac{1}{N} s_{e,j} (x_{ie} - \mu_e) (x_{ke} - \mu_e) (\tau_e^2 + \epsilon)^{-3/2} \end{aligned}$$

* Cost function \rightarrow

$$\frac{\partial h}{\partial x_{ij}} = \sum_{k=1}^M \sum_{l=1}^{n_y} \frac{\partial h}{\partial y_{ke}} \cdot \frac{\partial y_{ke}}{\partial x_{ke}} \cdot \frac{\partial x_{ke}}{\partial x_{ij}} =$$

$$= \sum_{k=1}^M \sum_{l=1}^{n_y} \frac{\partial h}{\partial y_{ke}} \cdot \gamma_l \cdot s_{k,j} (\tau_e^2 + \epsilon)^{-1/2} \left[s_{k,i} - \frac{1}{N} - \frac{1}{N} (x_{ie} - \mu_e) (x_{ke} - \mu_e) (\tau_e^2 + \epsilon)^{-1/2} \right]$$

$$= \frac{\pi}{N} (\tau_j^2 + \varepsilon)^{-1/2} \left[\frac{\partial h}{\partial y_{ij}} - \sum_{k=1}^M \frac{\partial h}{\partial y_{kj}} - \sum_{k=1}^M \frac{\partial h}{\partial y_{kj}} (x_{ij} - \mu_j) (x_{kj} - \mu_j) (\tau_j^2 + \varepsilon)^{1/2} \right]$$

$$= \frac{\pi_j}{N} (\tau_j^2 + \varepsilon)^{-1/2} \left[\frac{\partial h}{\partial y_{ij}} - \sum_{k=1}^M \frac{\partial h}{\partial y_{kj}} - (x_{ij} - \mu_j) (\tau_j^2 + \varepsilon)^{-1/2} \sum_{k=1}^M \frac{\partial h}{\partial y_{kj}} (x_{kj} - \mu_j) \right]$$

$$\begin{aligned} \frac{\partial h}{\partial y_{ij}} &= \sum_{k=1}^M \sum_{\ell=1}^{n_y} \frac{\partial h}{\partial y_{k\ell}} \cdot \frac{\partial y_{k\ell}}{\partial y_{ij}} = \sum_{k=1}^M \sum_{\ell=1}^{n_y} \frac{\partial h}{\partial y_{k\ell}} \cdot \delta_{ij} \cdot \hat{x}_{k\ell} = \\ &= \sum_{k=1}^M \frac{\partial h}{\partial y_{kj}} \cdot \hat{x}_{kj} \end{aligned}$$

$$\frac{\partial h}{\partial \beta_{ij}} = \sum_{k=1}^M \sum_{\ell=1}^{n_y} \frac{\partial h}{\partial y_{k\ell}} \frac{\partial y_{k\ell}}{\partial \beta_{ij}} = \sum_{k=1}^M \sum_{\ell=1}^{n_y} \frac{\partial h}{\partial y_{k\ell}} \cdot \delta_{ij} = \sum_{k=1}^M \frac{\partial h}{\partial y_{kj}}$$

