## Column 1

$\sigma(x) = \frac{1}{1+\exp(x)}$ — saturate grad, kill neuron

$\frac{d}{dx}\sigma(x) = \sigma(x)(1-\sigma(x))$

tanh

Relu: Leaky Relu: $\max(x, 0.1x)$

Cuz Relu: dead neurons with 0 grad

**Backprop**

Q: $64\times64$ RGB to binary categories need $256^{3\times64\times64}$ bits.

He Initialization $(2/n^{[L-1]})$: Relu
Xavier Initialization: tanh/sigmoid

$-W_{i,j}^{[L]} = N(0, \frac{1}{n^{[L-1]}})$

- assumes tanh activ.
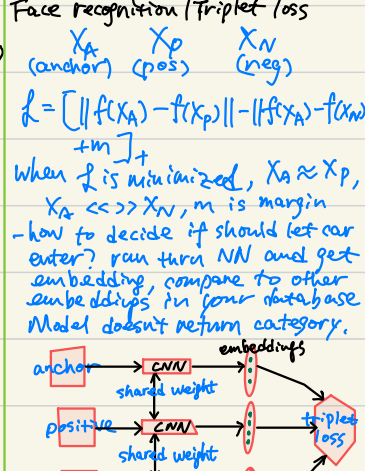- reduces vanishing/exploding gradients
- Xavier: $1/n_{prev}$
- He: $2/n_{prev}$
- Glorot: $2/(n_{prev}+n_{cur})$
- Derivation of Xavier:

$Var(a_i^{[L-1]}) = Var(a_i^{[L]})$
$\approx Var[Z_i^{[L]}]$

$= Var(\sum_{j=1}^{n^{[L-1]}} W_{ij}^{[L]} a_j^{[L-1]})$

$= \sum_{j=1}^{n^{[L-1]}} Var(W_{ij}^{[L]} a_j^{[L-1]})$

$= \sum_{j=1}^{n^{[L-1]}} E(W_{ij}^{[L]})^2 Var(a_j^{[L-1]})$

$+ E[a_j^{[L-1]}]^2 Var(W_{ij}^{[L]})^2$

$+ Var(W_{ij}^{[L]}) Var(a_j^{[L-1]})$

$= n^{[L-1]} Var(W_{ij}^{[L]}) Var(a_j^{[L-1]})$

$\Rightarrow Var(W) = \frac{1}{n^{[L-1]}}$

Uniform $(-d, d) \sim U/0, \frac{a^2}{3})$
$(a,b) \sim U(\frac{a+b}{2}, \frac{1}{12}(b-a)^2)$

Optimization: ① GD (+d-GD)
② Momentum $V_{dw} = \beta V_{dw} + (1-\beta)dW$
  $W = W - \alpha V_{dw}$
③ RMSProp $S_{dw} = \beta S_{dw} + (1-\beta)dW^2$
  $W = W - \alpha \frac{dW}{\sqrt{S_{dw}}+\epsilon}$

④ Adam $V_{dw} = \beta_1 V_{dw} + (1-\beta_1)dW$
  $S_{dw} = \beta_2 S_{dw} + (1-\beta_2)dW^2$
  $V_{corr\,dw} = V_{dw}/(1-\beta_1)^t$
  $S_{corr\,dw} = S_{dw}/(1-\beta_2)^t$
  $W = W - \alpha \frac{V_{corr\,dw}}{\sqrt{S_{corr\,dw}}+\epsilon}$

**Same for update b!**

Hyperparams:
① Batch size
② learning rate
③ learning rate decay

## Column 2

Mini batch / SGD better at avoiding saddle points

$L_1$ & $L_2$ Reg
$J_{L_1}(\vec{W}) = J(\vec{w}) + \lambda \sum_{i=1}^{\ell} |W_i|$
$J_{L_2}(\vec{W}) = J(\vec{w}) + \lambda \sum_{i=1}^{\ell} ||W_i||^2$
Update:
$W_i^{k+1} = W_i^k - d\lambda \, sign(w_i) - d\frac{\partial J}{\partial w_i}$
$W_i^{k+1} = W_i^k - 2d\lambda W_i - d\frac{\partial J}{\partial w_i}$ ← decay
Goal: limit weight growth (weight decay)
L1 encourages sparsity
Dropout: "Inverted Dropout" ÷ keep-prob
— drop $a^{[L]}$ after activation
Test: switch off Dropout
Early stopping: non-orthogonal approach to mitigate overfitting
Data augmentation:
— crop, flip, rotate, blur, change color
Python Numpy/Tensorflow
— np.var, np.mean, np.exp, np.sqrt
— matrix operations:
  a.dot(b) = np.dot(a,b)
— np.empty_like(x): creates empty matrix w/ x's shape
— broadcasting

Algo and Case studies:
① Neural Style Transfer
Style: $X_S$
$\mathcal{L}_S(f(X_S), f(\hat{X}))$
output: $\hat{X}$
context: $X_C$
$\mathcal{L}_C(f(X_C), f(\hat{X}))$
— NN: pre-trained network (VGG) for feature extraction
— $\mathcal{L}_C = \sum_{\ell \in \{l_c\}} ||F^\ell(I_c) - F^\ell(I)||^2$
— $\mathcal{L}_S = \sum_{\ell \in \{l_s\}} ||G(F^\ell(I_s)) - G(F^\ell(I)')||^2$ ← extract content
— $\mathcal{L}_C$: content loss, tensor $L^2$ norm loss, F: intermediate features
$\mathcal{L}_S$: style loss, G: graham matrix for style extraction

Face recognition/Triplet loss
$X_A$    $X_P$    $X_N$
(anchor) (pos)  (neg)

$\mathcal{L} = [||f(X_A) - f(X_P)|| - ||f(X_A) - f(X_N)|| + m]_+$

When $\mathcal{L}$ is minimized, $X_A \approx X_P$, $X_A \ll\gg X_N$, m is margin
— how to decide if should let car enter? run thru NN and get embedding, compare to other embeddings in your database. Model doesn't return category.

anchor → CNN → embeddings
  shared weight
positive → CNN →   → triplet loss
  shared weight
negative → CNN →

## Column 3

# CNN
- 3D: learnable params: $(f \times f \times n_c + bias)$
  • output shape: $(n-f+1)^2 \times n_f \times n_f^{+1}$
- "Valid": no padding
  "Same": pad s.t. output size = input size
- General formula for shap:
  $$\left\lfloor \frac{n-f+2p}{s} +1 \right\rfloor^2_{Floor} \times n_f$$

- Pooling layers:          usually:
  • no padding (normally)   • $f=2$, $S=2$
  • increases output sensitivity to position of image.

- Batchnorm layers
  • mini batch norm
  • apply before activation (on $Z_i^{(L)}$)
  • Test time: use running avg of mean/var
  • $\mu = \frac{1}{m}\sum_i^m Z^{(i)}$, $\sigma^2 = \frac{1}{m}\sum_i^m (Z^{(i)}-\mu)^2$
    $Z_{norm} = \frac{Z^{(i)}-\mu}{\sqrt{\sigma^2+\epsilon}}$, $\tilde{Z} = \gamma Z_{norm} + \beta$
  • $\gamma$ and $\beta$ are learnable parameters
  • w/o $\gamma$ and $\beta$: restrict network from learning larger values
  • why does batch norm work?
    — faster learning: each dim take similar values
    — mitigate covariance shift
    (make weights in later layers more robust to weight changes in earlier layers)
    — slight regularization effect: mean and var of mini-batch adds noise

**GANs:** $\min_{\theta_g} \max_{\theta_d} [E_x \log D_{\theta_d}(x) + E_z \log(1 - D_{\theta_d}(G_{\theta_g}(z)))]$

$G(z)$    Discriminator output for real data X    Discriminator for fake data G(z)
G →  → O ← if x=G(z) 0 otw
  grad → O
real image X
— $\theta_d$ wants to maximize objective s.t. $D(x)$ is close to 1 (real) and $D(G(z)) \approx 0$ (fake)
— $\theta_g$ wants to minimize objective s.t. $D(G(x))$ close to 1 (discriminator fooled into thinking G(z) is real)

Defense: ① safety net: network discerning fakes
② generate adversarial examples and label correctly
③ adversarial training: $L_{new} = L(w,b,x,y) + \lambda L(w,b,x_{adv},y)$
④ logit pairing: $L_{new} = L(w,b,x,y) + \lambda ||f(x|w,b) - f(x_{adv}|w,b)||_2^2$

Alternative cost function:
$J^{(D)} = -\frac{1}{m_{real}}\sum_{i=1}^{m_{real}} y_{real}^{(i)} \log(D(x^{(i)})) - \frac{1}{m_{gen}}\sum_{i=1}^{m_{gen}}(1-y_{gen}^{(i)})\log(1-D(G(z^{(i)})))$

$J^{(G)} = -J^{(D)} = \frac{1}{m_{gen}}\sum_{i=1}^{m_{gen}} \log(1-D(G(z^{(i)})))$

Label: $y_{real} = 1$, $y_{gen} = 0$    saturating cost
Non-saturating cost of $J^{(G)} = -\frac{1}{m}\sum_{i=1}^m \log(D(G(z^{(i)})))$
Methods of generating adversarial examples:
— Fast gradient sign: $X^* = X + \epsilon \, sign(w)$
— Input optimization

$D(G(z))$

# Loss functions:
- logistic:
$$J(\theta) = \frac{1}{m}\sum_{i=1}^{m}\left[-y^{(i)}\log(h_\theta(x^{(i)})) + (1-y^{(i)})\log(1-h_\theta(x^{(i)}))\right]$$
$$+ \frac{\lambda}{m}\sum_{j=1}^{n}|\theta_j| + \frac{\lambda}{2m}\sum_{j=1}^{n}\theta_j^2$$

- if output is probability, use logistic/softmax + cross entropy for loss.

$\log(\cdot)/\log(1-\cdot)$ input $\in(0,1)$

- tanh has vanishing grad
- ReLu has dead neurons.

Softmax: $S(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}}$

loss: $\sum_{i=1}^{n_g} -(y_i \log \hat{y}_i)$

---

Trigger Word Detection + Data acquisition

- need overfit / bias
- transition invariant
- covariance shift
- $\int\int$ over adam again
- backprop derivation

---

Recall: $\frac{TP}{TP+FN}$ ... Predicted

Precision: $\frac{TP}{TP+FP}$ ... Actual

Accuracy: $\frac{TP+TN}{Total}$

$F_1$-Score: $2\frac{precision \times recall}{precision + recall}$

Saliency map: $\frac{\partial \hat{y}}{\partial x}$

Occlusion sensitivity: confidence of pred w/o area

Class Activation map:



---

# Forward Prop:
$$Z = W_1 x + b, \quad A = \sigma(Z)$$
$$\hat{y} = W_2 A + b_2$$
$$L = \frac{1}{m}\|\hat{y} - y\|^2$$

# Back prop:
$$\frac{\partial L}{\partial W_2} = \frac{2}{m}(\hat{y}-y)A^T$$
$$\frac{\partial L}{\partial b_2} = \frac{2}{m}(\hat{y}-y)\vec{1}$$
$$\frac{\partial L}{\partial W_1} = \left(W_2^T \frac{2}{m}(\hat{y}-y)\odot A \odot(1-A)\right)x^T$$
$$\frac{\partial L}{\partial b_1} = \left(W_2^T \frac{2}{m}(\hat{y}-y)\odot A \odot(1-A)\right)\vec{1}$$