

Lecture 2

Unsupervised Learning

- No label or target class
- Find out properties of the structure of the data → Autoencoder
- Clustering

Supervised Learning

- Labels or target classes
- Learn from train data → Learning = how well the model do on the test data

Reinforcement Learning



Image Classification → Supervised learning

Nearest Neighbor → Supervised Learning

↳ Which of the other images in the set is the closest to this one (compute distance)

K - Nearest Neighbor → take k nearest neighbor and classify image based on majority

↳ higher k can give less variance

Split data {

- train data → 66% / 70% → train with train data
- validation data → 20% / 5% → improve performance with the validation data
- test data → 20% / 5% → check final performance with test data

Cross-validation

Split data into n-fold ↗ random k-fold for validation

- train data with another n-k fold
- use k validation test set to optimize hyperparameter
- Train data and test again for the next split/run
- Do it n runsplit, each time each validation fold

Then pick the hyperparameter with the least average loss

Linear Regression → Supervised Learning

$$\hat{y}_i = \theta_0 + \sum_{j=1}^n x_{ij} \theta_j$$

bias weight
weight
no. of feature inputs (column) feature

i-th sample that we predict

x_i = i-th row of x

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1d} \\ 1 & x_{21} & \dots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{nd} \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$$

bias

Loss Function (Optimization)

① Linear Least Square $\min_{\theta} J(\theta) = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$ $\theta = (x^T x)^{-1} x^T y$

② Maximum Likelihood Estimate $\rightarrow \theta = \max \log p(y|x, \theta)$
 ↳ based on probability distribution

if y is normal distribution and $y_i = x_i \theta \rightarrow \theta = (x^T x)^{-1} x^T y$

$$p(y_i) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y_i - \mu\theta)^2}$$

$$y_i \sim N(\mu\theta, \sigma^2)$$

Regression \rightarrow predict a continuous value (temp. of a room) \rightarrow We can transform regression output to
 Classification \rightarrow predict a discrete value binary output by interpreting output of
 ↳ Binary \rightarrow predict either 0 or 1 regression as probability distribution
 ↳ Multiclass \rightarrow predict set of N classes ($P > 50\% \rightarrow \text{Class A}$, $P \leq 50\% \rightarrow \text{Class B}$)

Logistic Regression

Estimate probability of an event, bounded between 0 and 1

Loss function with $p(y) = \text{Bernoulli trial}$ and activation function σ sigmoid

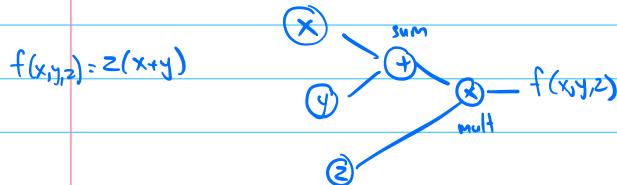
Binary Cross-Entropy Loss = $L(y_i, \hat{y}_i) = y_i \log \hat{y}_i + (1-y_i) \log(1-\hat{y}_i)$ \rightarrow Loss of all training
 (BCE)

Cost function = $-\frac{1}{n} \sum_{i=1}^n L(y_i, \hat{y}_i)$ \rightarrow Average loss over all training

Lecture 3

Computation Graphs

Forward Pass



Regression Loss

$$L_1 \text{ Loss} : L(y, \hat{y}; \theta) = \frac{1}{n} \sum_{i=1}^n \| \hat{y}_i - y_i \|_1,$$

$$\text{MSE Loss} : L(y, \hat{y}; \theta) = \frac{1}{n} \sum_{i=1}^n \| \hat{y}_i - y_i \|^2$$

Binary Cross Entropy → for binary classification

$$L(y, \hat{y}; \theta) = -\frac{1}{n} \sum_{i=1}^n y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)$$

For multiple class Cross Entropy Loss

$$L(y, \hat{y}; \theta) = \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^C -y_{ik} \log \hat{y}_{ik}$$

this generalise binary case → $y_{ik} = 1$ iff the true label (class) of the i -th sample

is the same as class k . 0 otherwise

(one-hot encoding)

Minimize the loss → find better prediction

Minimize Loss with respect to θ

Learning rate

$$\theta = \theta - \alpha \nabla L(y, \hat{y}) \rightarrow \text{find optimum } \theta \text{ that result in minimum loss} \rightarrow \text{Gradient Descent}$$

Lecture 4 - Optimization and Backpropagation

Neural Network

$$5 \text{ layers} \rightarrow f = W_5 \sigma \left(W_4 \tanh \left(W_3 \max (0, W_2 \max (0, W_1 x)) \right) \right)$$

of depth of NN = # of hidden + output layer

of node in input layer = # of column/feature in data

ReLU is the most used Activation Function

$$\begin{aligned} \text{Regression Loss} &\rightarrow L_1 \text{ loss } \rightarrow \sum_{i=1}^n \|y_i - \hat{y}_i\|_1 \\ &\rightarrow \text{MSE loss } \rightarrow \sum_{i=1}^n \|y_i - \hat{y}_i\|_2^2 \end{aligned}$$

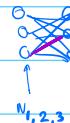
$$\text{Classification loss} \rightarrow \text{Cross Entropy Loss} = - \sum_{i=1}^n \sum_{k=1}^C y_{ik} \log \hat{y}_{ik}$$

Backpropagation $\nabla_{\theta} L(\theta)$

x_k - input variables

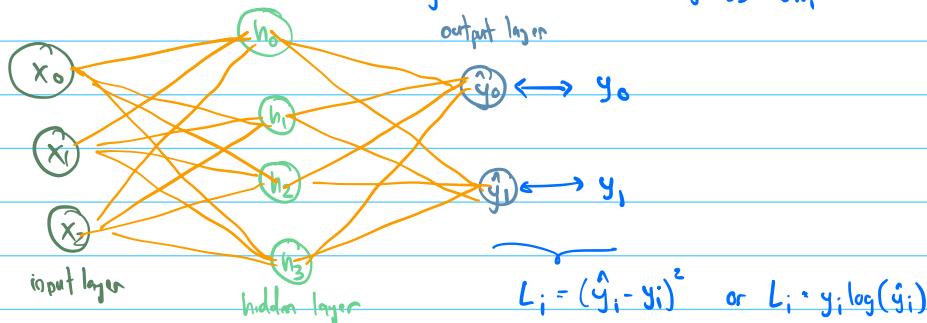
$w_{l,m,n}$ - weight of n -th neuron in l -th layer to m -th neuron in the next layer

\hat{y}_i - output of i -th row y_i - ground truth targets L - loss function



$$\text{forward} \Rightarrow f(x, w) = \tanh \left(\underbrace{w_2 \sigma(xw)}_s \right) \Rightarrow L(y, f(x, w)) = CE(\tanh \dots)$$

$$\text{backward} \Rightarrow \text{for each weight } w_i \rightarrow \text{for } w_2 \rightarrow \frac{\partial L}{\partial g} \cdot \frac{\partial g}{\partial w_2} \quad \text{for } w_1 \rightarrow \frac{\partial L}{\partial g} \cdot \frac{\partial g}{\partial s} \cdot \frac{\partial s}{\partial w_1}$$



$$y_i = A(b_{i,j} + \sum_j h_j w_{j,i,j})$$

let $A(x) = \max(0, x)$

$$\frac{\partial L}{\partial w_{j,i,j}} = \frac{\partial L}{\partial y_i} \cdot \frac{\partial y_i}{\partial w_{j,i,j}} = 2(\hat{y}_i - y_i) \cdot \frac{\partial \hat{y}_i}{\partial w_{j,i,j}}$$

$$= 2(\hat{y}_i - y_i) \cdot h_j$$

$$\frac{\partial \hat{y}_i}{\partial w_{j,i,j}} \rightarrow \text{use } \max(0, x) \quad \frac{\partial \hat{y}_i}{\partial w_{j,i,j}} = \begin{cases} 0 & \text{if } x \leq 0 \\ h_j & \text{if } x > 0 \end{cases}$$

Calculate # of weight in each layer

Layer = (# of neurons • # input neuron) + # of biases (usually just # of neuron in layer)
(Input layer doesn't have any weight, it starts in 1st hidden layer)

$$\text{if } A = \text{Sigmoid} = \frac{1}{1+e^{-s}} \text{ and } L = -\sum_{i=1}^{n_{\text{output}}} [y_i \log(g_i) + (1-y_i) \log(1-g_i)] \Rightarrow \text{BCE}$$
$$s = \sum_{j=1}^n b_j w_{ji}$$

$$\frac{\partial L}{\partial w_{ji}} = \frac{\partial L}{\partial g_i} \cdot \frac{\partial g_i}{\partial s} \cdot \frac{\partial s}{\partial w_{ji}}$$
$$= \left(-\frac{y_i}{g_i} + \frac{1-y_i}{1-g_i} \right) \cdot (g_i(1-g_i)) \cdot (b_j)$$
$$\frac{\partial L}{\partial s} = -y_i(1-g_i) + g_i(1-y_i)$$
$$= \hat{g}_i - y_i$$
$$= \frac{-y_i + \hat{g}_i}{g_i(1-g_i)} \cdot (g_i(1-g_i))(b_j)$$
$$= (-y_i + \hat{g}_i)b_j$$

Regularization

→ Use to reduce overfitting → Aim to lower validation error and increase training error

$$\text{Loss function} = \sum_{i=1}^n (\hat{y}_i - y_i)^2 + \lambda R(\theta) \Rightarrow \lambda R(\theta)$$

↳ some constant λ

Regularization techniques

→ L2 Regularization

↳ if λ is large, we care more about regularization

→ L1 Regularization

↳ if λ is small, care more about predicting

→ Max Norm Regularization

as close as possible

→ Dropout (Drop some of the neuron)

→ Early Stopping

$$\theta_1 = [0, 0.75, 0]$$

$$\theta_2 = [0.25, 0.5, 0.25]$$

L2 Regularization → $R(\theta) = \sum_{i=1}^n \theta_i^2$ → penalize larger weight → want weight to have similar values
→ Take all feature into account

$$R(\theta_{11}, \theta_{12}, \theta_{13}) = 0^2 + 0.075^2 + 0^2 = 0.5625 \quad R(\theta_{21}, \theta_{22}, \theta_{23}) = 0.25^2 + 0.5^2 + 0.25^2 = 0.375 \rightarrow \text{Use this one}$$

L1 - Regularization → $R(\theta) = \sum_{i=1}^n |\theta_i|$ → Enforce Sparsity → will focus on a few key feature that have high θ

$$R(\theta_{11}, \theta_{12}, \theta_{13}) = 0 + 0.75 + 0 = 0.75 \quad R(\theta_{21}, \theta_{22}, \theta_{23}) = 0.25 + 0.5 + 0.25 = 1$$

↳ use this one



Deep NN (multiple layer)

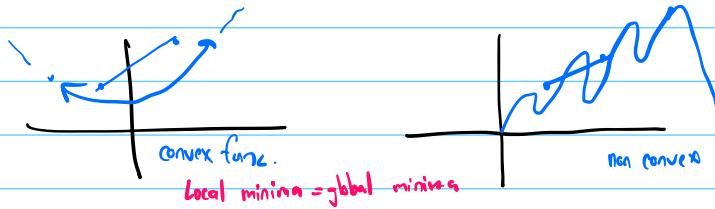
Adv → better at generalization

Problem → might cause vanishing gradient, memory

Lecture 5 - Scaling Optimization

Gradient Descent

- Gradient steps in direction of negative descent $x' = x - \alpha \nabla_x f(x)$
- Not guarantee to reach global optimum; only local minimum
- NN is a non-convex function with many local minima



- Learning Rate → Big Learning Rate might cause weight to bounce off local minima
→ Small Learning Rate might cause learning to be slow esp. if it's in plateau part

$$\theta^{k+1} = \theta^k - \alpha \nabla_{\theta} L_i(\theta^k, x_i, y_i)$$

Iterate until convergence $|\theta^{k+1} - \theta^k| < \epsilon$

$$L = \frac{1}{n} \sum_{i=1}^n L_i(\theta, x_i, y_i) \quad \text{where } \theta = \arg \min L$$

For Multiple Training sample

→ the gradient is average of all data (or batch)

$$\nabla_{\theta} L = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} L(\theta^k, x_i, y_i)$$

↑
no. of data

For big data, Gradient Descent might take a very long time to update weight since it only update weight after 1 epoch

Stochastic Gradient Descent (SGD)

- We divide data into mini-batch and then iterate through each minibatch.
Each time it finish the minibatch, recalculate the gradient and update the weight
- We can express total loss over the training data as expectation of all samples.
The expectation can be approximate with a small subset of data

$$\mathbb{E}_{i \in \{1, \dots, n\}} [L_i(\theta, x_i, y_i)] \approx \frac{1}{|S|} \sum_{j \in S} L_j(\theta, x_j, y_j) \text{ with } S \text{ as set of minibatch}$$

- if there's n data and each minibatch has m data, then #of minibatch $\propto \frac{n}{m}$
 - Smaller batch size means greater variance in the gradients \rightarrow noisy update
- For each iteration (batch)
1. Pick minibatch
 2. Feed it to NN
 3. Calculate the mean gradient of the minibatch
 4. Update the weights with the mean gradient
 5. Iterate Step 1-4 for n minibatch \rightarrow When we go through all minibatch, this is 1 epoch

Convergence of SGD

we want to minimize $F(\theta)$ (Loss for each minibatch)

$$E(H(\theta^k, x)) = \nabla F(\theta^k) \quad \theta^{k+1} = \theta^k - \alpha_k \nabla_{\theta} H(\theta^k, x) \quad \text{where } \alpha_1, \dots, \alpha_n \text{ is a sequence of positive step-sizes}$$

This converge to local (global) minimum if this is all true

1. $\alpha_n \geq 0, \forall n \geq 0$
2. $\sum_{i=1}^{\infty} \alpha_n = \infty$
3. $\sum_{i=1}^{\infty} \alpha_n^2 < \infty$
4. $F(\theta)$ is strictly convex

Problem w/ SGD \rightarrow it scales weight equally across all dimension (1 learning rate for all dim)

Gradient Descent with Momentum → Look at the direction of gradient (average gradient)

$$v^{k+1} = \beta v^k - \alpha \nabla_{\theta} L(\theta^k)$$

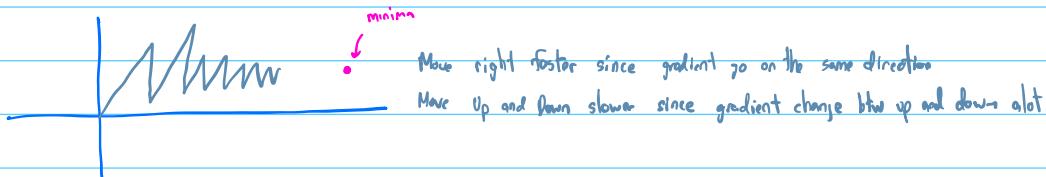
↑ momentum
velocity (in vector-value)
so direction important

↓ current weight

$$\theta^{k+1} = \theta^k + v^{k+1}$$

↑ updated weight
velocity

- If the gradient point in the same direct, build up velocity (v^{k+1}) and weight will change more
 - If gradient keep changing direction, then we want to slow down (take smaller step)
- v^{k+1} will stay the same cuz it get add and subtract



Hyper parameter = β, α

Root Mean Squared Prop (RMSProp) → look at the size of gradient (square of average gradient)

- Dampening direction with high variance (gradient) by divide learning rate by exponential average of squared gradients
- This mean high gradient will has less effect on changing of θ
- This also mean that we can increase learning rate α and it will
1. speed up the direction with less variance/gradient (like momentum)
 2. while not effect the direction of high variance/gradient as much like of damping

$$\theta^{k+1} = \theta^k - \alpha \cdot \frac{\nabla_{\theta} L}{\sqrt{s^{k+1} + \epsilon}}$$

where $s^{k+1} = \beta \cdot s^k + (1-\beta)[\nabla_{\theta} L \cdot \nabla_{\theta} L]$

↑
exponential
average sum of
squared

$$\nabla_{\theta} L^2 = \left(\sum_{i=1}^n \frac{\partial L}{\partial \theta_i} \right)^2 = \text{size of gradient}$$

Adam (Adaptive Moment Estimation)

Combine Momentum and RMSProp idea

where Momentum = θ increase faster when keep going in the same direction and θ stay still when keep going in opposite direction

and RMSProp = dampening direction with high gradient/variance which mean we can increase the α to help speed up

Combine together, Adam can

1. speed up when go in same direction
2. increase learning rate to speed up direction with less variance while not effecting direction with high variance too much.

$$\begin{aligned} \text{first momentum} \quad m^{k+1} &= \beta_1 m^k + (1-\beta_1) \nabla_{\theta} L(\theta^k) \rightarrow \hat{m}^{k+1} = \frac{m^{k+1}}{1-\beta_1^{k+1}} \\ \text{second momentum} \quad v^{k+1} &= \beta_2 v^k + (1-\beta_2) \nabla_{\theta}^2 L(\theta^k)^2 \rightarrow \hat{v}^{k+1} = \frac{v^{k+1}}{1-\beta_2^{k+1}} \end{aligned} \quad \left. \begin{array}{l} \theta^{k+1} = \theta^k - \alpha \left(\frac{\hat{m}^{k+1}}{\sqrt{\hat{v}^{k+1}} + \epsilon} \right) \\ \text{for bias toward 0} \\ \text{correction} \end{array} \right\} \beta_2 \text{ to the power of } k+1$$

Hyperparameter = $\beta_1, \beta_2, \alpha, \epsilon$

Second Derivative

- Don't work well with minibatches but...
- Converge to local minima fast (in term of # of iteration)
- Use curvature of gradient to take a more direct route

Conclusion

Gradient Descent → calculate gradient of all data and adjust the weight

SGD → calculate gradient and adjust weight per batch

GD with Momentum → Gradient go in some direction make θ change faster or stay still if direction is fluctuated

RMSProp → Dampen change of gradient for the direction with high variance/gradient. So we can increase learning speed

Adam → change of θ faster if gradient go in the same direction and dampen gradient if it's high variance. So we can increase learning rate

$$GD \rightarrow \theta^{k+1} = \theta^k - \alpha \nabla_{\theta} L \quad SGD \rightarrow \theta^{k+1} = \theta^k - \alpha \nabla_{\theta} L_{\text{minibatch}} \quad GD \text{ w/ Momentum} \rightarrow \theta^{k+1} = \beta \theta^k - \alpha \nabla_{\theta} L \quad \theta^{k+1} = \theta^k + v^{k+1}$$

$$RMSProp \rightarrow s^{k+1} = \beta s^k + (1-\beta) \nabla_{\theta} L^2 \quad \theta^{k+1} = \theta^k - \alpha \frac{\nabla_{\theta} L}{\sqrt{s^{k+1}} + \epsilon}$$

$$\begin{aligned} \text{Adam} \rightarrow m^{k+1} &= \beta_1 m^k + (1-\beta_1) \nabla_{\theta} L \rightarrow \hat{m}^{k+1} = \frac{m^{k+1}}{1-\beta_1^{k+1}} \\ v^{k+1} &= \beta_2 v^k + (1-\beta_2) \nabla_{\theta}^2 L^2 \rightarrow \hat{v}^{k+1} = \frac{v^{k+1}}{1-\beta_2^{k+1}} \end{aligned} \quad \left. \begin{array}{l} \theta^{k+1} = \theta^k - \alpha \frac{\hat{m}^{k+1}}{\sqrt{\hat{v}^{k+1}} + \epsilon} \end{array} \right\}$$

Lecture 6 - Training NN Part 1

Learning Rate Decay

Many options → 1. $\alpha = \frac{\alpha_0}{1 + (\text{# of epochs} \times t)}$
(t is decay rate)

→ 2. Step decay $\Rightarrow \alpha = \alpha - t\alpha$ every n epochs

→ 3. Exponential decay $\Rightarrow \alpha = t^{\text{epoch}} \cdot \alpha_0$

→ 4. $\alpha = \frac{+}{\sqrt{\text{epoch}}} \cdot \alpha_0$

Learning means generalization to unknown dataset \rightarrow test optimize parameters w/ unknown dataset

Training set \rightarrow train NN

Validation set \rightarrow optimizing hyperparameter, check generalization progress

Test set \rightarrow test at the very end

if Training error $\xrightarrow{\text{high}}$ have bigger model
train longer (underfitting)
diff architecture
diff hyperparameter / LR
 \downarrow
low

if Validation error $\xrightarrow{\text{high}}$ More data
(overfitting) Regularization
New model architecture
 \downarrow
low

if diff btw training error and validation $\xrightarrow{\text{high}}$ Make training data more similar to test data
Data augmentation
New model architecture
 \downarrow
low

Test error high $\xrightarrow{\text{high}}$ More dev set data
 \downarrow
Done

Underfitting → training and validation loss/error can be further decrease at the end of training

Overfitting → training loss is decrease while validation stop decrease or increase
(so it does not generalise well anymore) ↑
not as much

Bad sign if validation loss is less than training loss

Hyperparameter Tuning

1. Manual Search → You randomly select value for hyperparameter
2. Grid Search → Define a list or a range with step for all hyperparameter
(coarse) Then iteration over all value of hyperparameter
so if there's 3 hyperparameters, 4 values each
Then the number of iteration = 4^3

After getting the best hyperparameter → Define a list or range around that value → Refine Grid Search

3. Random Search → Define a range for all hyperparameter and a number of iteration (combination) that you want to random.

How to start Training

- Start with Single training sample
→ check that model has capable of overfitting (accuracy should be 100%)
- Increase to handful of sample
- As you increase, model should generalise more

Fix Underfitting by - 1. Increase model's capacity → bigger model size → more layers/nodes/weight

Fix Overfitting by - 1. Reduce capacity → smaller model size → less layer/node/weight
- 2. Regularization → weight decay, Dropout, Batch Norm, L1/L2 Reg, Early Stop
- 3. Data augmentation → More training data

Lecture 7 - Training NN Part 2

Regression Losses

$$\text{L2 Loss} \rightarrow L_2 = \sum_{i=1}^n |y_i - g_i|^2$$

- Sum of squared difference
- Prone to outlier
- Efficient to optimize
- Optimum = mean

$$\text{L1 Losses} \rightarrow L_1 = \sum_{i=1}^n |y_i - \hat{g}_i|$$

- sum of absolute diff
- Robust (cost of outlier is linear)
- Costly to optimize
- Optimum = median

Sigmoid

Score of class ($y_i = 1$) $y_i \in \{0, 1\}$ 2 classes

$$p(y_i = 1 | x_i, \theta) = \sigma(s) = \frac{1}{1 + e^{-s}}$$

Softmax (use as activation function for Multiclass Classification Problem)

→ Output probability for each class

$y_i \in \{1, 2, \dots, C\}$ C Classes

score of class k

$$p(y_i = k | x_i, \theta) = \text{softmax}(s_k)$$

weight of class k
(each class has its own weight
for each input)

$$p(y_i = 1 | x_i, \theta) = \frac{e^{x_i \theta_1}}{e^{x_i \theta_1} + e^{x_i \theta_2} + e^{x_i \theta_3}}$$

$$p(y_i = 2 | x_i, \theta) = \frac{e^{x_i \theta_2}}{e^{x_i \theta_1} + e^{x_i \theta_2} + e^{x_i \theta_3}}$$

$$p(y_i = 3 | x_i, \theta) = \frac{e^{x_i \theta_3}}{e^{x_i \theta_1} + e^{x_i \theta_2} + e^{x_i \theta_3}}$$

weight that is compared to that
output class (score)

$$p(y_i | x_i, \theta) = \frac{e^{x_i \theta_{y_i}}}{\sum_{k=1}^C e^{x_i \theta_k}}$$

sum of all the softmax output of class per 1 input row

For Numerical Stability = $\frac{e^{s_{y_i} - s_{\max}}}{\sum_{k=1}^C e^{s_k - s_{\max}}}$ to reduce computation in the case that
each $e^{s_{y_i}}$ is very big

Cross Entropy Loss (Maximum Likelihood Estimate)

$$L_i = -\log \left(\frac{e^{s_i}}{\sum_{k=1}^C e^{s_k}} \right)$$

normally have y_{ik} here but b/c of one-hot bits encoding
it's 1 for value of predict label = true label

Hinge Loss (SVM Loss)

$$L_i = \sum_{k=1, k \neq y_i}^C \max(0, s_k - s_{y_i} + 1) \rightarrow \text{stop learning when } L_i < 0$$

if there's several model with $L_i < 0$, then we cannot diff. which model is better
(the score for the correct predicted label is clearly higher) cuz all the loss is set to 0.

CE can show better which model is better (closer to 0), CE always to improve

Hinge Loss saturates (no longer optimize at some point)

Regularization

- Change with respect to weight θ

$$L = \frac{1}{N} \sum_{i=1}^N L_i + \lambda R^2(w)$$

If use L-2 Reg = $\sum w_i^2$

$$L1-\text{Reg} = \sum_{i=1}^D |w_i|$$

$$L2-\text{Reg} = \sum_{i=1}^D w_i^2$$

Activation Function

x and w is matrix

On linear layer ($s = xw + b$), the deriv are

dout = upstream derivative

$$\frac{\partial L}{\partial x} = \text{dout} \cdot w^T \quad \frac{\partial L}{\partial w} = x^T \cdot \text{dout} \quad \frac{\partial L}{\partial b} = \text{dout}$$

when deriv matrix-matrix mult \rightarrow get transpose, vector-vector mult \rightarrow normal or matrix-vector

$$\text{Sigmoid} \rightarrow y' = \frac{1}{1+e^{-s}} \quad \frac{\partial y'}{\partial w} = y'(1-y')$$

Disadvantage \rightarrow If s is very high or low, the gradient is almost 0 (saturate neuron cause vanishing gradient)

\rightarrow Output is always positive (not zero-center)

\rightarrow all of the weight update will either be positive or negative $\left(\frac{\partial L}{\partial w} = \frac{\partial L}{\partial s} \cdot \frac{\partial s}{\partial w} \cdot \frac{\partial s}{\partial w} = y'(1-y') \cdot x \right)$

\rightarrow Create a zig-zag path toward optimal w

so it does not depend on the weight, so all weight will always be increase or decrease based on $y'(1-y') > 0$ or < 0

TanH \rightarrow Advantage \rightarrow zero center

Disadvantage \rightarrow still saturate (gradient move slow when score are very high or low)

ReLU $\rightarrow \sigma(x) = \max(0, x)$

\rightarrow Advantage \rightarrow Fast convergence

\rightarrow Does not saturate (large and consistent gradient)

\rightarrow Disadvantage \rightarrow If score < 0, then ReLU = 0 then you cannot improve and get out (Dead ReLU)

\rightarrow give some slightly positive biases (0.01) to make ReLU stay active for most inputs

called Leaky ReLU if $\alpha = 0.01$ always

Parametric ReLU $\rightarrow \sigma(x) = \max(\alpha x, x)$

\rightarrow not result in dead ReLU

α = hyperparameter that can be trained via backpropagation to have diff number

Advantage \rightarrow Gradient do not die

Maxout Units $\rightarrow = \max(w_1^T x + b_1, w_2^T x + b_2)$

Advantage \rightarrow can create any approximate of an activation function
with enough combination of Maxout Unit

\rightarrow Generalization of ReLU, Linear, do not die, not saturate

Disadvantage \rightarrow Each Maxout Unit require new set of weight \rightarrow Increase # of parameters

Weight Initialization \rightarrow good initialise can result in better local minimum

If all weight initial with 0 / same value \rightarrow Get same score, same gradient \rightarrow No symmetry breaking, all neuron behave the same

If initialise with Gaussian with small Random Number ($\mu=0, \sigma^2=0.01$)

\hookrightarrow small weight cause small output \rightarrow output approach 0 very fast for each layer in forward pass

\hookrightarrow b/c the input for the next layer is small, the gradient (which uses the input) will be very small as well
 \hookrightarrow or output of last layer for calculation

and approach 0 with the deeper layer. \Rightarrow Vanishing gradient cause by small output

If initialise w/ Gaussian with Big Number ($\mu=0$, $\sigma^2=1$)

↳ b/c the weight is large, the score is large as well

↳ if we use tanh, the large score will cause saturate neuron and the gradient will move very little (close to 0) → Vanishing gradient caused by saturated activation function (for tanh)

Xavier Initialisation ($\mu=0$, $\sigma^2 = \frac{1}{n}$)

$$\text{calculate for each layer} \rightarrow \text{Var}(s) = \text{Var}\left(\sum_{i=1}^n x_i w_i\right) = \sum_{i=1}^n \text{Var}(x_i w_i) = n(\text{Var}(x) \text{Var}(w))$$

n = # of input (based on each layer)

Goal : we want $\text{Var}(s) = \text{Var}(x) \rightarrow \text{Var}(x) = n(\text{Var}(x) \text{Var}(w))$

(variance of output = input)

$$\frac{1}{n} = \text{Var}(w)$$

n = # of input neuron for layer l

So Xavier w/ tanh will make the output stay Gaussian distribution

But Xavier do not work with ReLU cuz ReLU kill half the data



For ReLU, we have to use Kaiming Initialisation = $\text{Var}(w) = \frac{2}{n}$ (ReLU kill half the data, so we use half the input, $\frac{n}{2} \rightarrow \frac{1}{\frac{n}{2}} = \frac{2}{n}$)

Exercise : Autoencoder

- unsupervised learning

- NN that is train to copy its input to its output

- consist of encoder and decoder

- encoder learn feature

- layer between input and output has smaller size (latent space)

Lecture 8 - Training Neural Network Part 3

Data Augmentation → One way to better generalise, type of regularization

You always crop the image first, to have image dimension as input of what you want

1. Flipping

2. Crop + Flip → Do multiple crop on diff region of the image (crop to the dim of input you want)
and then flip all of them.

→ More robust to partial exclusion

3. Random Brightness and Contrast change

4. Random Crop → Random crop on training

→ Fixed crop on testing

Use same data from data augmentation when compare NN

Advanced Regularization → Another way to improve generalise

1. Add weight decay $\Rightarrow \theta_{k+1} = \theta_k - \alpha \nabla_{\theta_k} L - \frac{\lambda}{2} \theta_k \rightarrow$ so $\theta_{k+1} = \theta_k - \alpha (\nabla L(\theta_k) + \nabla R(\theta_k))$

If this is L2 Reg → gradient of regularization = $\frac{\partial}{\partial w} \sum_{i=1}^n \|w_i\|^2$

2. Early Stopping \Rightarrow stop after validation loss increase n times in a row.

3. Ensemble \Rightarrow Train multiple model in diff way (diff activation, augmentation, loss, etc.)

that lead to diff local optima and average it.

Expected combined error should decrease linearly w/ ensemble size \Rightarrow provided errors are uncorrelated

2 type of ensemble

→ 1. Bagging \Rightarrow build diff model like ensemble \rightarrow divide training data into k parts (can be overlap)

inputs

→ 2. Dropout \Rightarrow disable a random set of neurons (do it during the forward pass) (random row set for each minibatch)

→ reduce capacity of NN

Advantage \rightarrow prevent one node from being too specialised, make other node learn other feature as well \Rightarrow More Robust

Disadvantage \rightarrow another hyperparameter (λ of dropout) that need to be optimize

→ Need larger model + more training time

Consider as ensemble technique
We're training large ensemble of
models, each new Model
w/ diff minibatches
w/ shared parameter

✓ bagging method w/ shared param.

On test time \Rightarrow turn on all neuron \rightarrow no dropout for validation & test data

→ Need to do weight scaling inference rule during training

→ Adapt weight to match the condition at training time (multiply weight with dropout prob.)

$$E(S) = p \left(\sum \theta_i x_i \right)$$

\uparrow
 $i=1$
keeping probability (or $1-p$ if $p = \text{dropping probability}$)

Batch Normalization

Want to do Normalization after finishing the layer to not have saturate neuron

1. Normalize mean and variance of the input to the activation function w/ the mean & var of minibatch

$$\hat{x}^k = \frac{x^k - E(x^k)}{\sqrt{\text{Var}(x^k)}} \quad \begin{matrix} \text{mean of minibatch over feature } k \\ \text{variance of } x^k \end{matrix}$$

Batch Normalization layer → apply after Fully Connected Layer (FC) or CNN
and before non-linear activation function

2. Allow Network to change the range of x^k during backprop

$$y^k = \gamma^k \hat{x}^k + \beta^k$$

parameter that could learn to undo normalization $\hat{x}^k \rightarrow x^k$

$$\text{if } \gamma^k = \sqrt{\text{Var}(x^k)} \quad \beta^k = E(x^k)$$

Advantage → converge to optima faster

→ got much deeper network → more stable gradient

→ Reduce internal covariate shift → Change in input distribution

At test time → Compute μ and σ^2 by running an exponentially weighted averaged across training minibatches

$$\text{Var test} \rightarrow \text{Var}_{\text{running}} = \beta \text{Var}_{\text{running}} + (1-\beta) \text{Var}_{\text{minibatch}}$$

$$\mu_{\text{test}} \rightarrow \mu_{\text{running}} = \beta \mu_{\text{running}} + (1-\beta) \mu_{\text{minibatch}}$$

} use this μ and σ^2 to normalise test data instead

BN get worse performance when batch size is smaller

Other Normalisation → Layer Norm, Instance Norm, Group Norm

Lecture 9 : CNN

FC use too much hyperparameter from weight, so we need to
 restrict degree of freedom (apply some inductive bias \rightarrow look for specific feature/thing)
 \hookrightarrow Weight sharing, extract same features independent of location

Convolutions

Filter Kernel \rightarrow a matrix that's slide across input matrix

\rightarrow Equal to dot product b/w filter weight w and input x_i -th chunk (output for each chunk = 1x1)

$$z_i(\text{Output}) = w^T x_i + b$$

\rightarrow Need to have same depth as input (if input is RGB, then filter need to have depth = 3, result depth = 1)

\rightarrow Output dim is depend on dimension

\hookrightarrow multiply weight in each depth to each input in chunk and sum it all up

1. Size of kernel (the bigger, the smaller output dim)

$$(x_1w_{11} + x_2w_{12} \dots) + (x_1w_{21} + x_2w_{22} \dots) + (x_1w_{31} + x_2w_{32} \dots)$$

2. # of stride (reduce output dim)

$$\boxed{}_{\text{depth dim. 1}} + \boxed{}_{\text{depth dim. 2}} + \boxed{}_{\text{depth dim. 3}} = \boxed{}_{\text{output}}$$

3. padding

* The weight inside filter kernel is a hyperparameter that need to be learn

so kernel w/ 3 depth and size 5x5 has # of weight/params = $3 \times 5 \times 5 = 75$

Convolution Layer

\rightarrow Apply different filter with diff. weights

\rightarrow Resulting layer from all the filter kernel to input = Activation maps ($D \times W \times H$ of output)

\rightarrow Number of depth in Activation maps = Number of filter kernel apply to input (w/ same dim)

\rightarrow Filters that apply to input form a Convolution layer (n filter \times $W \times H$ of filter)

together (depth is implicitly given by input depth)

Stride

\rightarrow place filter S pixel from previous input

\rightarrow can have illegal stride (filter does not fit, there's a remainder)

$$\text{Output dim} = \left(\frac{N-F}{S} + 1 \right) \times \left(\frac{N-F}{S} + 1 \right) \quad S = \text{Stride} \\ N = H \text{ or } W \text{ dim of input} \quad F = \text{H or W of filter}$$

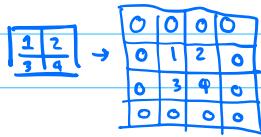
If there's a fraction, then the stride is illegal \rightarrow Can be solved by padding

round the
other one

Shrinking dim too fast is not a good idea \rightarrow so we apply padding to not reduce dim

Padding

- Use to stop dimension from shrinking too quickly
- Also allow corner pixel to be used in calculation more than once
- Most common = zero padding (add 0 pixels to all four sides)



$$\text{Output Size} = \left(\frac{N+2P-F}{S} + 1 \right) \times \left(\frac{N+2P-F}{S} + 1 \right)$$

if fraction, then do floor operation

Type of Convolution → 1. Valid Convolution = no padding

2. Same Convolution = Input dim = output dim → Set padding $P = \frac{F-1}{2}$

Example - Input image = $32 \times 32 \times 3$

10 filter 5×5
Stride = 1
Padding = 2

$$\left. \begin{array}{l} \text{Output} = \frac{N+2P-F}{S} + 1 = \frac{32+2(2)-5}{1} + 1 = 31 + 1 = 32 \\ \text{Output} = 32 \times 32 \times 10 \end{array} \right\}$$

No. of params per filter = $W \times H \times D + 1$ bias = $5 \times 5 \times 3 + 1 = 76$ (each filter has 1 bias)

Total params for this Convolution layer = # of filters × params per filter = $10 \times 76 = 760$ params

Convolution Neural Network

Input → Convolution layer → Activation function → Convolution layer (use output of Activation function as input → Activation → → FC Layer
(do convolution with n filters, then apply activation function to it))

w/ dim depth of # of filter on last
Convolution layer

Pooling Layer → reduce size of input (downsampling), no learning (no weight)

→ Come after the Convolution layer (Conv layer = Feature Extraction, Pooling layer = Feature Selection)

→ Max pooling → A filter that only pick max value of the chunk as the output

select the strongest
feature in the region

→ Average pooling → A filter that average value from the chunk

Pooling layer has 2 hyperparams

1. Filter dim F (Fxp)

2. Stride S

Output of Pooling layer

$$W_{out} = \frac{W_{in} - F}{s} + 1 \quad H_{out} = \frac{H_{in} - F}{s} + 1 \quad D_{out} = D_{in}$$

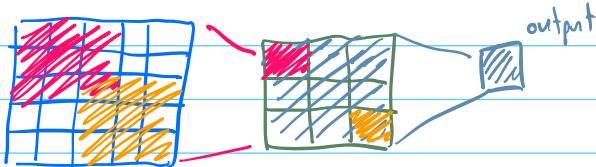
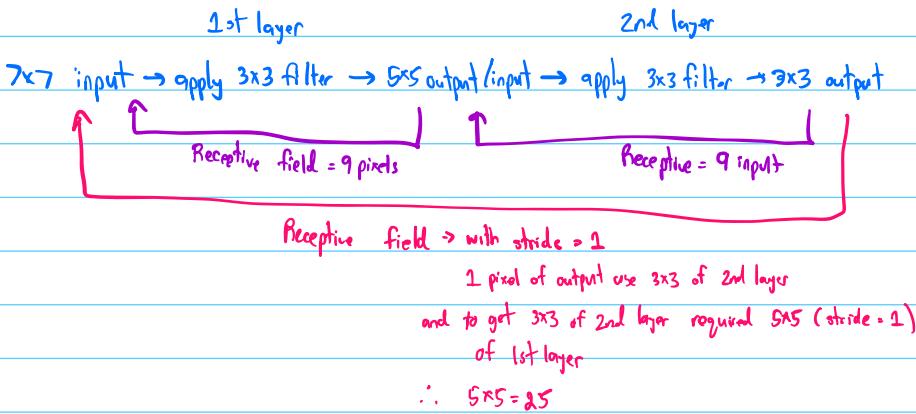
The Final FC Layer → Make decision w/ the extracted features from convolutions

Receptive Field (Spatial extend of the connectivity of a convolutional filter)

How many input space were used to get this output pixel

If filter 3×3 apply on input, result in 1 pixel. Then receptive field = 9 pixels

When we chain Convolution filter



* Note, it's more efficient to use many kernel than 1 big filter kernel (have less weight + More nonlinearity)

Spatial Batch Normalization

When using BN on Convolution layer/NN

input = NxCxHxW

We compute μ and var for each channel for all data in minibatch. (normally we normalise for each feature)

$$N \times C \times H \times W \rightarrow 1 \times C \times 1 \times 1$$

If there's C channel, each with HxW dim, and there's N minibatch

then we get μ and var for each C using stats for $H \times W \times N$ (entire weight in the channel \times all the data in minibatch)

C_i, use data from entire weight of channel i of all the data in minibatch

$$N \times C \times H \times W \rightarrow N \times H \times W \times C \rightarrow (N \times H \times W) \times C$$

find μ and σ^2

Lecture 10 - CNN Part 2

Classic Architecture

• LeNet

- Digit Recognition (10 classes)
- use sigmoid/tanh } not common anymore
- use average pool }

As NN is deeper, $W \times H \downarrow$, no. of filter \uparrow

• Alex Net (first architecture to use CNN with Imagenet)

- use large size filter \rightarrow not common

As NN is deeper, $W \times H \downarrow$, no. of filter \uparrow

- use Max pool ($3 \times 3, s=2$)

- use ReLU

- 1000 classes

• VGG Net

- more simplicity

always
- use 3×3 Convolution filter, stride = 1 that result same dim for output (add padding)

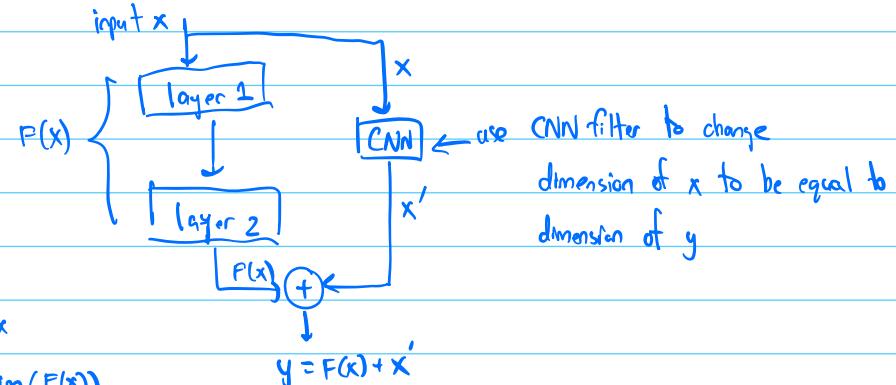
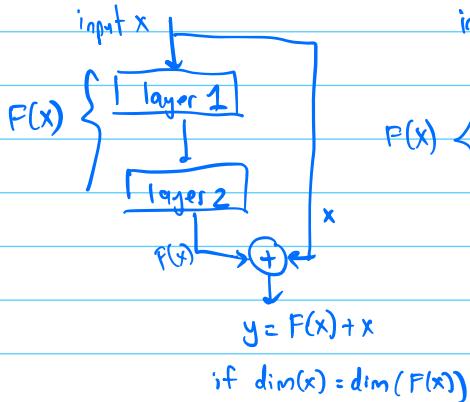
- Max pool ($2 \times 2, s=2$)

As NN is deeper, $W \times H \downarrow$, no. of filter \uparrow

As we add more layer, training is harder + vanishing gradient + degradation

Skip Connection \leftarrow solution

Residual Block (ResNet) \rightarrow With this technique, you can add abt more layer



So when there's a vanishing gradient (either from input or backprop) in the $F(x)$ ($F(x)$ close to 0), then we can still use value of x for further layer and ignore ("skip") output $F(x)$
look like we're skipping the output

$$x^{L+1} = f(y) = f(F(x) + x) = f(x) \text{ instead of } f(0) \text{ w/ no residual}$$

if $F(x) = 0$

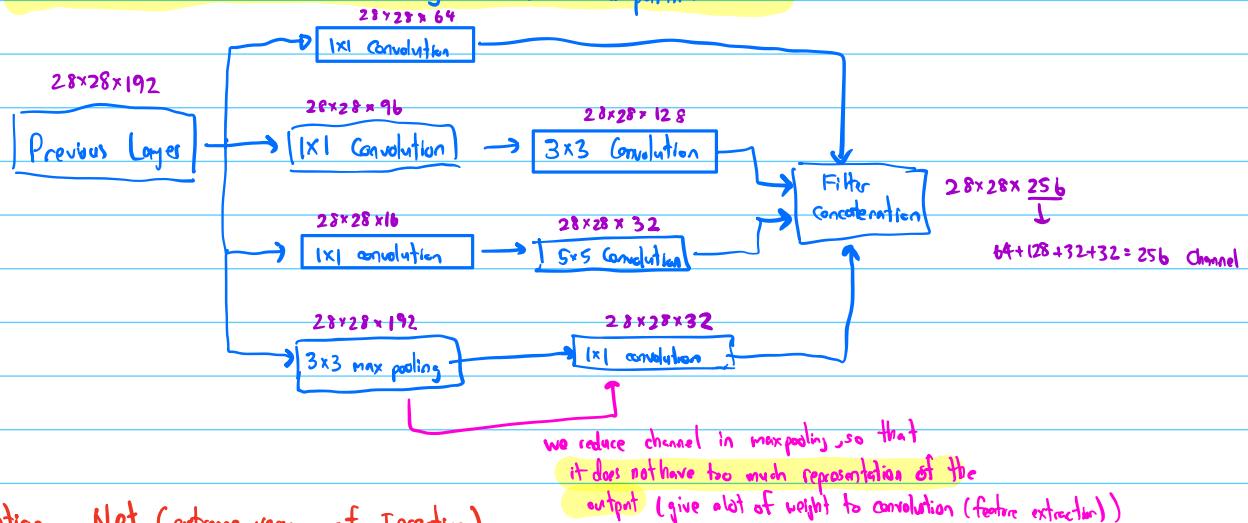
- Identity (input x) is easy for the residual block to learn ($x^{l+1} = x$)
- Not hurt the performance, can only improve

1x1 Convolutions

- Reduce or Increase the channel
- Using n filter of 1×1 convolution is same as using n FC layer
- Add non-linearity (not by itself but w/ activation func afterward)
- When use in last layer instead of FC to get output \rightarrow # of output = # of filter

Inception Layer (Google Net)

- Learn best weight when training by using multiple convolutions (kernel) size with multiple filters and pooling layers in parallel and concat all layer together. This help NN find more useful feature
- Use 1×1 convolution to reduce dimensionality (channel) and computational cost



Xception Net (extreme version of Inception)

- Use depthwise Separable Convolutions
- 36 Conv layer w/ skip connection
- Depthwise Separable Convolution = Filter apply only on certain depth

If input has D channel, then use D filter of $W \times H \times 1$ (1 channel) $\Rightarrow D \times W \times H \times 1$

Output has same # of channel as input

Then use N filter of $1 \times 1 \times D$ to turn into N Channel $\rightarrow N \times 1 \times 1 \times D$

* normal conv would use N filter of $W \times H \times D \rightarrow N \times W \times H \times D$

(help reduce computational cost)

Fully Convolutional Network

- We can convert Fully Connected layer into Convolutional Layer
- $1 \times 1 \times D$ Convolution w/ N filter = Fully Connected layer w/ N output ($\# \text{of filter} = \# \text{of output}$)
- When use FC Layer \rightarrow size of input is fixed (take same # of inputs)
- When use Convolution layer $\rightarrow |X|$ does not care about size, so input is flexible

FCN : Semantic Segmentation

- same output size as input (but # of channel of output = # of class want to predict)
- go from ^{from} little input size (from convolution) to big output size (Upsampling)
- Give score of classes for each pixel \rightarrow "void class" for unlabelled pixels

Type of Upsampling

First, double the size of input w/ blank space in between each input, then ...

1. Interpolation - Increase the size using
 1. Nearest Neighbor (look at value of neighbor)
 2. Bilinear (weighted average of pixel)
 3. Bicubic

2. Transposed conv - Increase size w/ convolution - upscaling

- Convolution filter (learned) (up-convolution)
↑ add value to blank space using filter

U-Net (Autoencoder using Convolution, pooling, up-convolution, and skip connection)

Use skip connection to add details back to decoder so that it has more detail to work with

Encoder (Contraction Path)

- Height & Width \downarrow , Depth \uparrow
- Use Conv & max pooling to down sampling
- ReLU

Decoder (Expansion Path)

- Height & Width \uparrow , Depth \downarrow
- 2×2 Up conv \rightarrow double input
- Skip Connection \rightarrow Concatenate feature maps of encoder to decoder \rightarrow then 2 3x3 Convs \rightarrow final layer
txt Conv to map
channel to # classes

→ input and output width has to be the same

Lecture 11 - Recurrent Neural Network and Transform

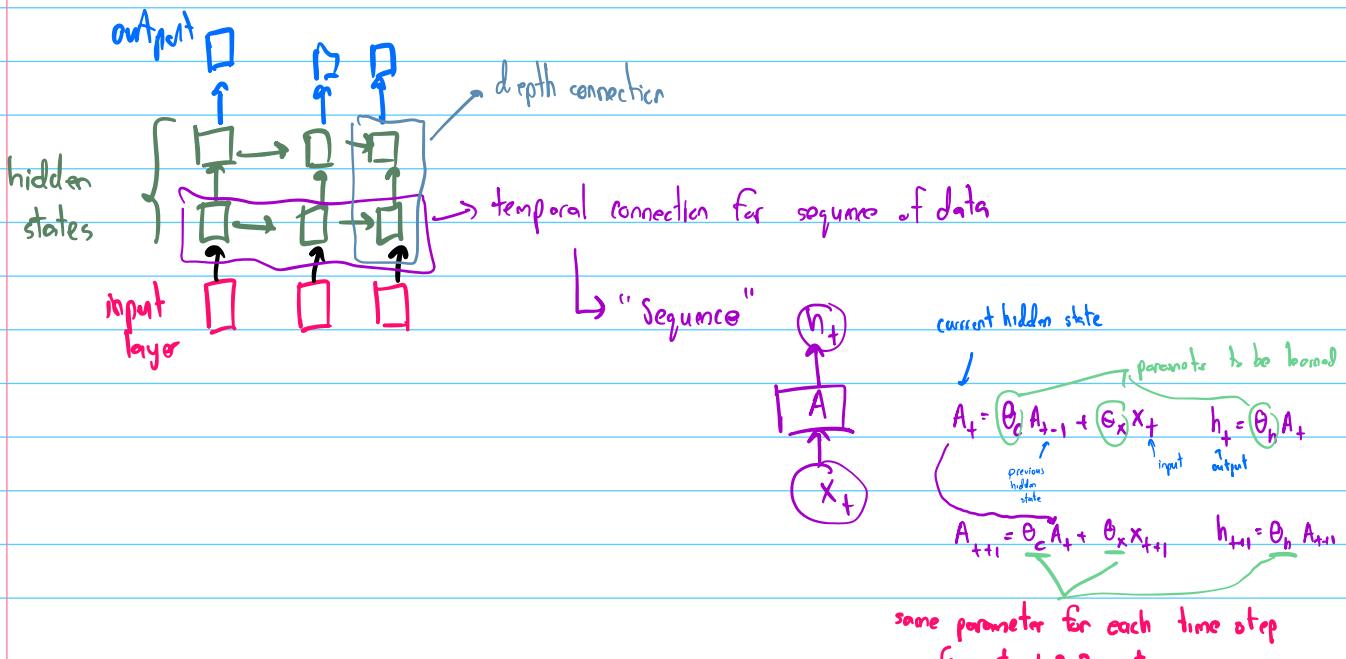
Transfer Learning

- We don't have enough resources or data → use Pretrain model instead
- Pretrain should have similar low level feature or same input
- Use the pretrain feature extraction (CNN layers + FC layer)
- Replace last FC layer to have correct output classes (If have large dataset, we can train more FC layer w/ low learning rate)

Recurrent Neural Network

- Input/Output can be sequence (set of inputs = dependent to each other)
- Flexible
 - One to one (classic NN for image classification)
 - One to many (image captioning - one image input, output as a sentence)
 - many to one (language recognition - what language is this sentence)
 - many to many (Machine translation - translate language to other language)
 - start produce output after all the input are processed (require all info first)
 - many to many (Video classification)
 - can produce first output right away with first input

Multi-layer RNN



RNN has long term dependency (A_t is depend on A_0) and since weight is always the same, it become exponentially big

$$A_t = (\theta_c)^t A_0 \quad \text{b/c } A_1 = \theta_c A_0 \quad A_2 = \theta_c A_1 = \theta_c^2 A_0 \quad A_3 = \theta_c A_2 = \theta_c^3 A_0 \dots$$

Problem of RNN

if θ_c is ^{very} small (< 1) $\rightarrow \theta^t \rightarrow 0$ ∵ vanishing gradient (also can vanish from activation function sigmoid or tanh)
 if θ_c is ^{large} (> 1) $\rightarrow \theta^t \rightarrow \infty$ ∵ exploding gradient → fix by using gradient clipping (if value $> x$, use x)

θ_c is matrix → eigen decomposition $\theta_c^t = Q \Lambda Q^{-1}$

Make eigenvalue = 1 → Allow cell to maintain state → Long Short Term Memory (LSTM)

Long Short Term Memory (LSTM)

Setting eigenvalue = 1 to avoid vanishing gradient problem

- use tanh
- Cell - transport info through the Unit
- Gate - remove or add info to the cell state

3 Type of Gate

1) Forget Gate \rightarrow Use sigmoid = output between 0 (forget) and 1 (keep)

- use x_t (input of current) and h_{t-1} (past hidden state)

- multiply output to cell state (output = 0.7 → keep 70% of info)

$$f_t = \text{sig}(\theta_{xf} x_t + \theta_{hf} h_{t-1} + b_f)$$

2) Input Gate \rightarrow Decide which value will be update

- Use sigmoid to check if output (tanh) will be add to cell state or not

$$i_t \text{ (input)} = \text{sig}(x_t \theta_{xi} + h_{t-1} \theta_{hi} + b_i)$$

- Output (g_t) tanh will be between (-1, 1)

$$C_t = f_t \odot C_{t-1} + i_t \odot g_t$$

↑ forget gate ↑ previous cell state
 ↑ input gate ↑ output from tanh → $g_t = \tanh(x_t \theta_{xg} + h_{t-1} \theta_{hg} + b_g)$

3) Output Gate \rightarrow Decide which values will be output (and use in next cell as a hidden state)

$\rightarrow o_t = \text{sig}(x_t \theta_{xo} + h_{t-1} \theta_{ho} + b_o) \rightarrow$ decide if we will send hidden state h_t to the next cell

$$h_t = o_t \odot \tanh(C_t)$$

→ Produce Actual Output and pass to next time step

All weight and bias are the same for each time step and are learned through backprop

How LSTM solve vanishing gradient

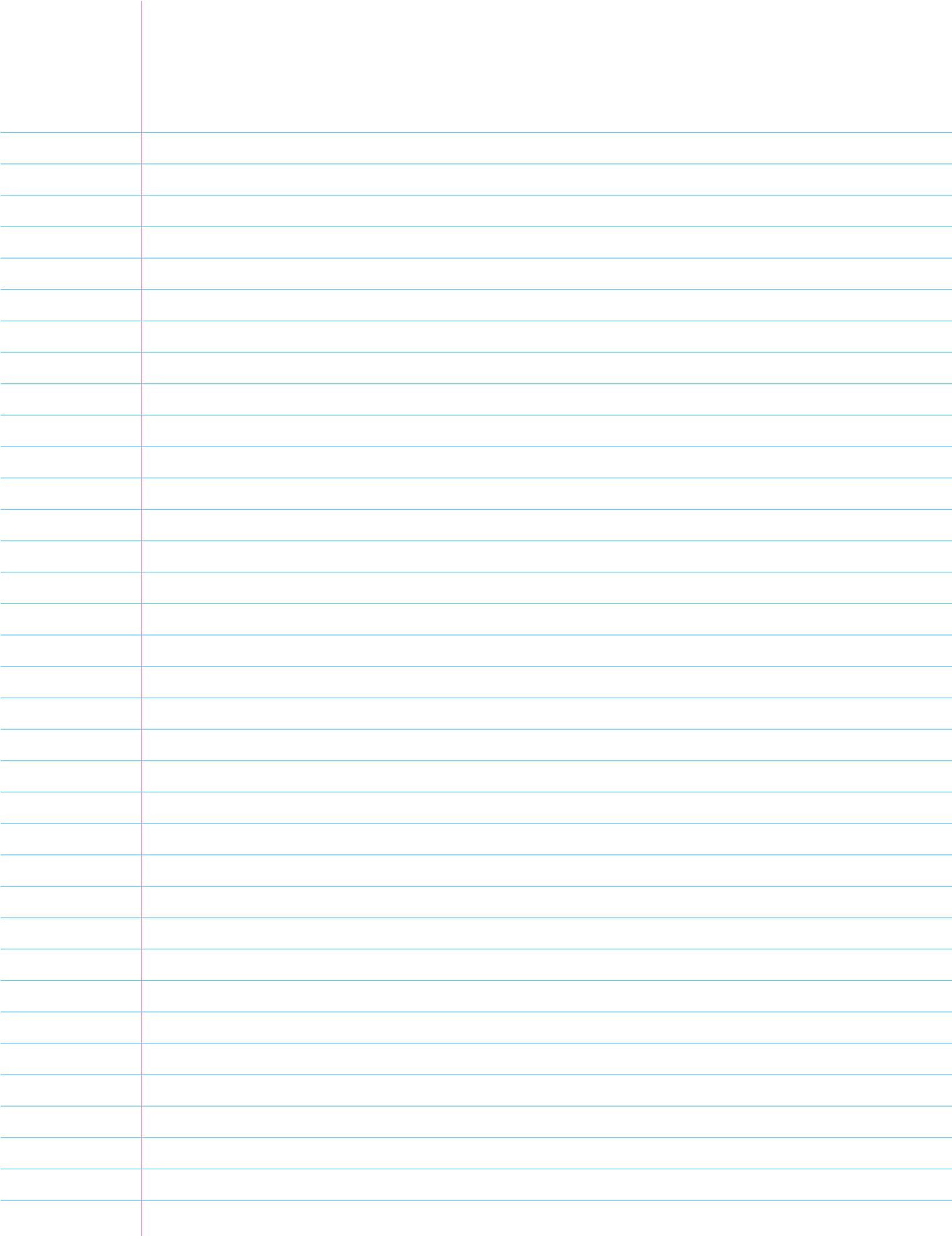
$$C_t = f_t \odot C_{t-1} + i_t \odot g_t$$

1. Vanishing from weight (f_t)

f_t use sigmoid = 1 for important info, so when you get important info in the past, there will be no exploding/vanishing gradient from exponential

2. Vanishing from activation function ($\tanh = g_t$)

→ if g_t has vanishing gradient, we still has $f_t \odot C_{t-1}$ to propagate backward further
(similar to residual connection)



Knowledge from past exam

21ss

1. Maximum Likelihood Estimate \rightarrow supervised learning
2. \uparrow Batch size \uparrow Learning rate, \downarrow Batch size \downarrow LR (more confidence)
3. Batch Norm \rightarrow formula add non-linearity
4. Non suitable for activation function \rightarrow change value of input to function
5. Batch Norm is not regularization it's regularizations
6. Convolution - translation equivariant ($f(g(x)) = g(f(x))$)
 - Not rotation equivariant \rightarrow need data augmentation to learn picture w/ rotation
 - Convolution layer is not translation invariant but
Other part of CNN might cause output to be translation invariant such as pooling layer

BCE & CE not use $\frac{1}{N}$

only $\sum_{i=1}^n \sum_{j=1}^K y_{ij} \log \hat{y}_{ij}$

or

$$\sum_{i=1}^n y_i \log \hat{y}_i + (1-y_i) \log (1-\hat{y}_i)$$