

Lecture 2:

Bayes Rule:

Propensity for
deserving a certain
value of x given
a certain value of y .

$$P(y|x) = \frac{\text{Posterior}}{\text{Likelihood} \downarrow \text{Prior: what we know of } y \text{ before seeing } x} = \frac{P(x|y) P(y)}{P(x)}$$

what we know about "y after seeing x"

$$\mathbb{E}[f(x)] = \sum_x f[x] \Pr(x)$$

$$\mathbb{E}[f(x)] = \int f[x] \Pr(x) dx.$$

Generalization Error: Error of the hypothesis (model measured w.r.t. probability dist. from which training samples are drawn).

Basically how well we will do on future data, but we don't know the future data " x_i " and also their labels y_i . We only know the "range" of all possible (x, y) .

$$E_{\text{gen}} = \int$$

Underfitting | Overfitting: Learning of an algorithm depends on

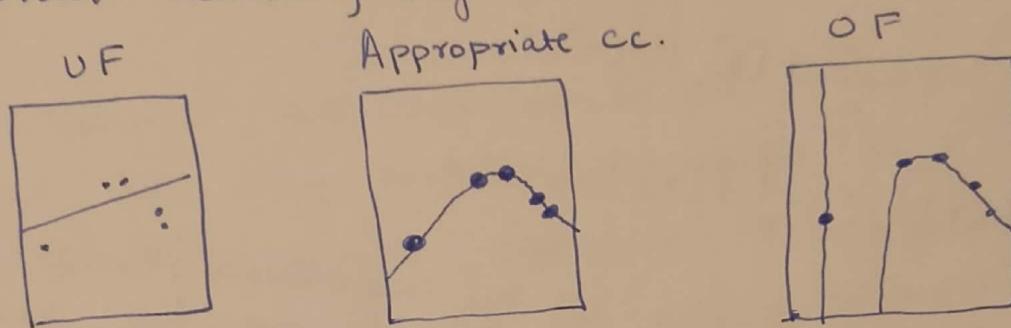
- ① training error small
 - ② making gap b/w training & test error small
- generalization gap.

In the case of underfitting: ① not fulfilled.
Overfitting: ② F ① but ~~not~~ not ②.

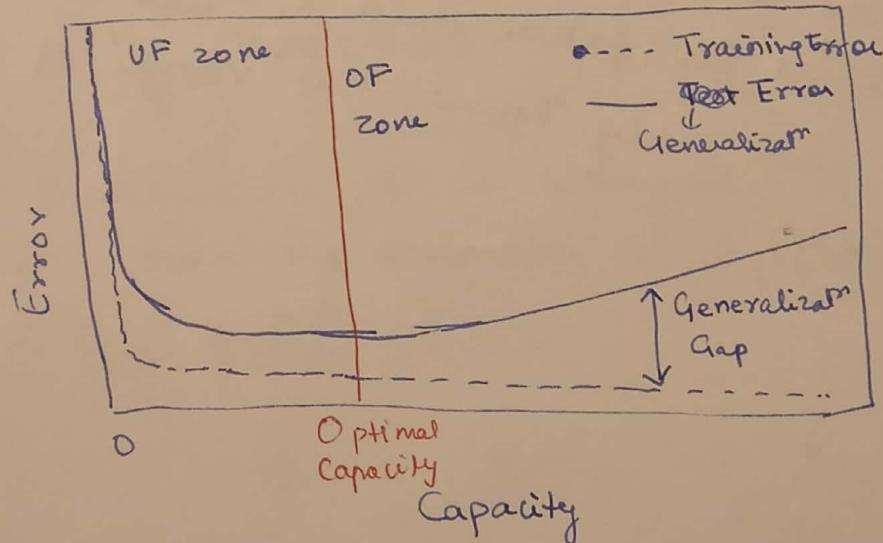
②

Model Capacity: OF and UF controlled by model capacity.
If UF then ↑ the model capacity, while in OF ↓ the model capacity. (because in UF model less expressive while in OF, model memorizes the data).

→ Capacity controlled by hypothesis space. (i.e. set of solutⁿ from which learning algorithm chooses from).



Capacity vs. Error



Hyperparameters: Not set during learning, set in advance or determined through a nested learning procedure. It overall controls the behaviour of algorithm.

Point Estimation: Provides a single best prediction for a parameter or a function, denoted by $\hat{\theta}$ and \hat{f} . The point estimation function doesn't need to guarantee or even have something close to true values or it might even not have same range. (3)

For our sake, we assume θ be fixed and $\hat{\theta}$ as r.v.

Function Approx.: $\exists f(\cdot)$ which can approx. relationship from x to y i.e. $y = f(x) + \epsilon$, ϵ error/noise.

→ In functⁿ estimⁿ, we approximate true f with \hat{f} based on some model. Basically, the class of model (i.e. y) is already known, so functⁿ approx. simplifies to estimating the functⁿ parameter θ . (Can also be seen as point estimⁿ).

Bias: Bias of estimator

$$\text{Bias}(\hat{\theta}_m) = \underbrace{E(\hat{\theta}_m)}_{\text{estimator}} - \underbrace{\theta}_{\text{true value}}$$

When $\text{Bias}(\hat{\theta}_m) = 0$, $\hat{\theta}$ is unbiased. If bias is zero, only if $m \rightarrow \infty$, then the estimator is asymptotically unbiased.

$$\lim_{m \rightarrow \infty} E(\hat{\theta}_m) - \theta = 0$$

$$\Rightarrow \lim_{m \rightarrow \infty} E(\hat{\theta}_m) = \theta.$$

estimator property is

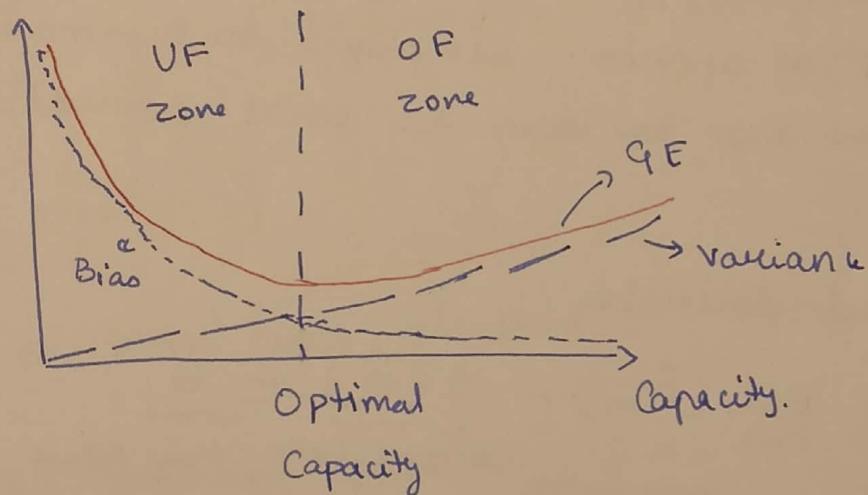
Variance: Another variance.

$$\begin{aligned} \text{Var}(\hat{\theta}_m) &= E[(\hat{\theta} - E(\hat{\theta}))^2] \\ &= E(\hat{\theta}^2) - (E(\hat{\theta}))^2 \end{aligned}$$

$$\text{Standard error} = SE(\hat{\theta}_m) = \sqrt{\text{Var}(\hat{\theta}_m)}$$

Mean Squared Error. (MSE). : Good measure of generalization error (4)

$$\begin{aligned}
 \text{MSE}(\hat{\theta}_m) &= \mathbb{E}[(\hat{\theta} - \theta)^2] \\
 &= \mathbb{E}(\hat{\theta}^2) - 2\mathbb{E}(\hat{\theta})\mathbb{E}(\theta) + \mathbb{E}(\theta^2) \quad (\because \theta \text{ & } \hat{\theta} \text{ are independent}) \\
 &= \mathbb{E}(\hat{\theta}^2) - 2\theta\mathbb{E}(\hat{\theta}) + \mathbb{E}\theta^2 + \mathbb{E}^2(\hat{\theta}) - \mathbb{E}^2(\theta) \\
 &\quad (\because \theta \text{ is fixed}). \\
 &= \underbrace{(\mathbb{E}(\hat{\theta}) - \theta)^2}_{\text{Bias}} + \underbrace{\mathbb{E}(\hat{\theta}^2) - \mathbb{E}^2(\hat{\theta})}_{\text{Var}} \\
 &= (\text{Bias}(\hat{\theta}_m))^2 + \text{Var}(\hat{\theta}_m).
 \end{aligned}$$



Maximum Likelihood Estimator (MLE). : We define estimator as a function of data i.e. $\hat{\theta}_m = g(x^{(1)}, \dots, x^{(m)})$ where $x^{(1)}, \dots, x^{(m)}$ are i.i.d.

But where does $g(\cdot)$ come from?

→ Commonly, it comes from MLE, defined as

$$\hat{\theta}_{ML} = \underset{\theta}{\operatorname{argmax}} P_{\text{model}}(\mathbf{x}, \theta) = \underset{\theta}{\operatorname{argmax}} \log P_{\text{model}}(\mathbf{x}, \theta)$$

$$= \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^m P_{\text{model}}(x^{(i)}, \theta)$$

$$= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log P_{\text{model}}(x^{(i)}, \theta)$$

- Activation functions should be non-linear and differentiable
- Regression: For regression, $\sigma(\cdot)$ is simply an identity functⁿ i.e. $y_p = a_p$.
- Classification: logistic sigmoid for binary problems, or softmax for multiple classes.

$$\sigma_{\text{binary}}(a) = \frac{e^{(a)}}{1 + e^a}$$

$$\sigma_{\text{multi-class}}(a) = \frac{e^{(a_i)}}{\sum_k e^{a_k}}$$

$$\rightarrow \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad \frac{d}{dx} \tanh(x) = 1 - \tanh^2(x)$$

$$\rightarrow \text{ReLU} = \max\{0, a\}.$$

Advantages of tanh and sigmoid : Bounded, smooth, differentiable

Disadvantages: Vanishing gradients, gives only dense representation as no way can be zero.

Advantages of ReLU: enforces sparsity, easy to implement.

Disadvantages: unbounded, exploding gradients.

- Why is it important to enforce that the activation functⁿs should be non-linear?
- Because, if they are linear, then essentially the multi-layer feed forward network can be scaled down to effectively to one layer network!!

→ Suppose we have network with N hidden layers and "h" activations. (6)

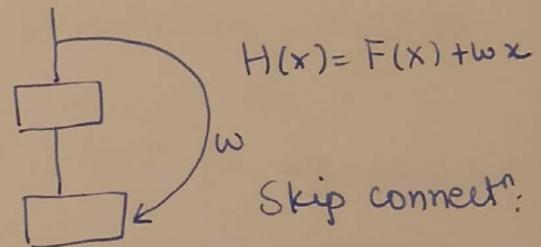
$$\begin{aligned}
 y_k(x, w) &= \sigma \left(w_k^{(N)} \circ h^{(N-1)} \circ \dots \circ h^{(2)} \left(w_k^{(2)} \cdot h^{(1)}(w_j^{(1)} \cdot x) \right) \right) \\
 &= \sigma \left(w_n^{(N)} \circ h^{(N-1)} \circ \dots \circ h^{(2)} \left(w_k^{(2)} \cdot w_j^{(1)} \cdot h^{(1)}(x) \right) \right) \\
 &= \underbrace{\sigma \left(w_n^{(N)} \dots w_k^{(2)} \right)}_{w'} \cdot \underbrace{h^{(N-1)} \circ h^{(2)} \circ h^{(1)}(x)}_g \\
 &= \sigma(w' g(x)) \quad \text{Equivalent collapsed network.}
 \end{aligned}$$

→ Depth of NN: # of activations or weight matrices!!

Skip Connections:

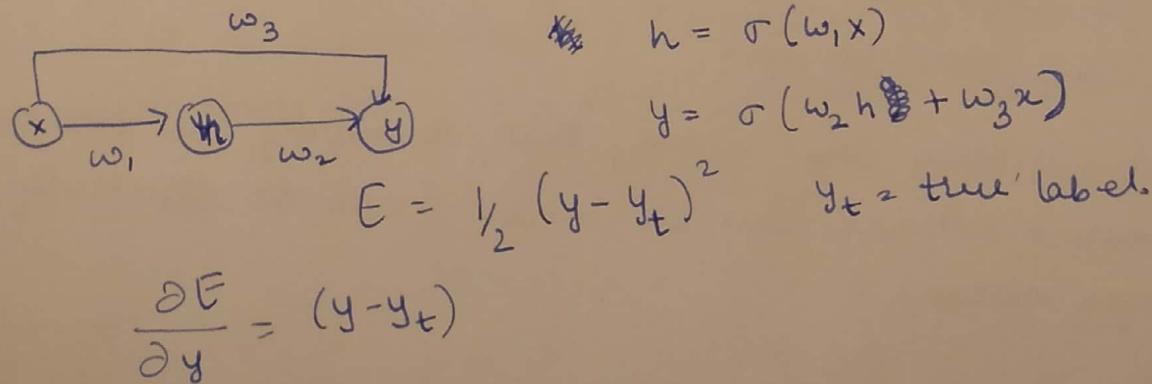


$H(x) = F(x) + x$
 "Residual Network"
 where "weights" are set to identity, only connections are made.



where we connect all the neurons of i^{th} layer by skipping some layers ~~and~~ to j^{th} layer with ^{some} weights associated to the connection.

Skip Connect Ex. in two layered NN.



$$\frac{\partial E}{\partial w_2} = \frac{\partial E}{\partial y} \times \frac{\partial y}{\partial w_2} = (y - y_t) \sigma(w_2 h + w_3 x) (1 - \sigma(w_2 h + w_3 x)) h.$$

~~$$\frac{\partial E}{\partial w_3} = \frac{\partial E}{\partial y} \times \frac{\partial y}{\partial w_3} = (y - y_t) \sigma(w_2 h + w_3 x) (1 - \sigma(w_2 h + w_3 x)) x.$$~~

$$\frac{\partial E}{\partial w_1} = \frac{\partial E}{\partial y} \times \frac{\partial y}{\partial h} \times \frac{\partial h}{\partial w_1}$$

$$= (y - y_t) \times \sigma(w_2 h + w_3 x) [1 - \sigma(w_2 h + w_3 x)] w_2 \times \sigma(w_1 x) (1 - \sigma(w_1 x)) x$$

Weight Space Symmetry:

Odd Activation Functn's: Odd functn's such as $\tanh(a)$, since $\tanh(-a) = -\tanh(a)$. If consider a fully connected two layered NN, then setting any consecutive wt's w_{ji} and w_{kj} negative cancels. M hidden units leads to M symmetries, so any "w" given wt. vector will be one of the 2^M leading to same O/P.

Non-unique: ^{Ordering} Also can swap hidden units' wts and i/p's, which leaves network unchanged. $1 \Rightarrow M!$ orderings.

→ Total $M! 2^M$ symmetries → This is not limited ~~to~~ only for tanh functn.

Universal Approximators: Basically given enough ^{model} capacity can approximate any functn (i.e. Turing Complete).

Why would Learning fail?

- ① Insufficient # HVs.
- ② Optimizat^m alg. cannot find optimal parameter values
- ③ training algo. chooses wrong function ~~param~~^{by} overfitting

Training Feed forward NN.

→ Gradient Descent Optimizat^m.

$$\omega^{(t+1)} = \omega^{(t)} - \eta \nabla E(\omega^{(t)})$$

learning rate

→ Stochastic Gradient Descent. (SGD).

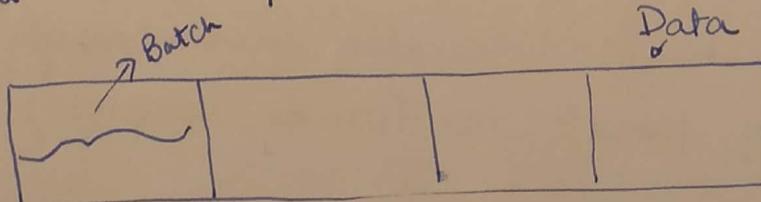
For iid's based on ML, we can decompose the error function into sum of terms, one for each training sample.

$$\rightarrow E(\omega) = \sum_n E_n(\omega).$$

→ Online Learning : Updates are made to "w" based on the error gradient one data pt. at a time i.e

$$\omega^{(t+1)} = \omega^t - \eta \nabla E_n(\omega^{(t)})$$

→ Minibatches are an intermediate when we update based on small batches of data pts. (SGD).



epochs: # times you iterate over training. the whole data

Minibatch | Batch size: Divide your data accordingly to do back propagation after seeing that subset of data.

iterations: # times you do BP to complete one epoch.

Initializing wts : If w is near zero, so is w.x. Then
if we apply activatⁿ as sigmoid, then operating sigmoid becomes
roughly linear. \Rightarrow Network collapses into an approximate
linear model. ⑨

still

- Always advisable to start wts near zero randomly in the initial layers so that model starts linearly and then wts ↑ and in a way become non-linear. (Circulum learning)
- Large wts lead to bad results. (Next Chp May be).
- Generally advisable to scale the i/p's s.t. they have zero-mean and std. deviatⁿ 1. (Standardizatⁿ of Data)
 \Rightarrow Ensures all dim. are handled equally in the regularizatⁿ.

How to standardize data?
$$\frac{x^{(i)} - \mu}{\sigma}$$
, where $\mu = \text{mean}$ &
 $\sigma = \text{STD.}$

Directional derivative: of $f(x)$ is in the direction of unit vector "u" is the instantaneous rate of change of f in dirⁿ of "u":

$$\nabla_u f(x) = \lim_{\alpha \rightarrow 0} \frac{f(x+du) - f(x)}{\alpha} = u^T \nabla f(x).$$

To minimize " f ", want to find dirⁿ"u" in which it ↓'s

$$\min_{u, u^T u=1} u^T \nabla_x f(x) = \|u\|_2 \| \nabla_x f(x) \|_2 \cos \theta$$

where θ is \angle b/w " x " and " u ".

→ This comes down to $\min_u \cos \theta$, since $\nabla_x f(x)$ ~~depends~~ does not depend upon ~~u~~ , i.e. in the opposite dirⁿ of $\nabla_x f(x)$.

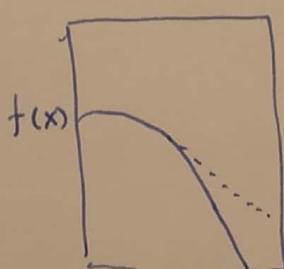
i.e. $x' = x - \epsilon \nabla_x f(x)$, where ϵ = learning rate
 \uparrow

how big a step to take in the opposite dirⁿ of
desirative gradient.

First Order Derivatives → Jacobian

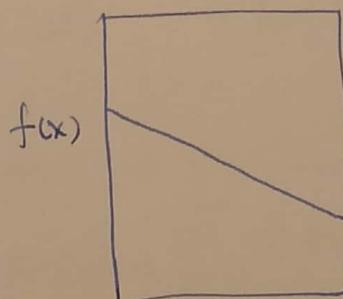
Second Order " " → Help us understand the properties or the nature of the curvature and ~~estimates~~ estimates how much improvement does the gradient step has.

-ve curvature



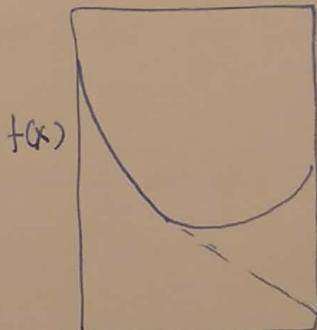
Cost functn ↓
faster than gradients
predicts.

No curvature



Gradient predicts
↓ correctly.

the curvature



f(x) is
slower than expected
and eventually ↑'s,
so taking large steps
would ↑ the functn
again.

THE HESSIAN MATRIX.

$$H(f(x))_{ij} = \frac{\partial^2}{\partial x_i \partial x_j} f(x).$$

Directional second derivatⁿ in dirⁿ of unit vector d.
is given by $d^T H d$.

→ GD in NN.

→ In most of cases NN, can be thought of as a parametric model that defines cond. distribution $p(y|x, \theta)$. It assigns the cost functⁿ based on ML.

$$\rightarrow \theta_{ML} = \operatorname{argmin}_{\theta} - \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}, \theta)$$

→ For some given mapping $\hat{y} = f(x, \theta)$, we do not want a functⁿ "f" that estimates a single \hat{y} , but rather a distribution (i.e. $p(y|x_0)$)

→ Suppose we use Gaussian to model this conditional i.e. $N(y; \hat{y} = f(x, \theta), \sigma^2)$, the conditional log-likelihood is

$$\sum_{i=1}^m \log p(y^{(i)} | x^{(i)}, \theta) = -m \log \sigma - \frac{m}{2} \log(2\pi) - \sum_{i=1}^m \frac{\|f(x_i, \theta) - y^{(i)}\|^2}{2\sigma^2}$$

→ We are only interested in the terms depending on θ ,

$$\Rightarrow \theta_{ML} = \operatorname{argmin}_{\theta} \sum_{i=1}^m \|f(x_i, \theta) - y^{(i)}\|^2$$

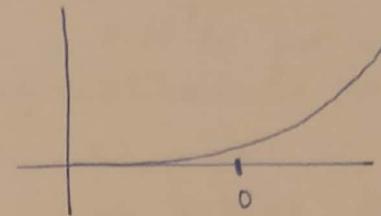
i.e. in some way equivalent to MSE.

→ In the cost functⁿ when $p_{model}(y|x)$ is a gaussian $N(y, f(x, \theta), I)$, then we get $J(\theta) = \mathbb{E}_{x, y \sim P_{data}} \|f(x, \theta) - y\|^2$

Activation function's.

(12)

→ Softplus function: $\text{softplus}(x) = \log(1 + e^x) = \log(1 + e^{\max(0, x)})$.
softened version of $x^+ = \max(0, x)$.



Binary Classification: If use sigmoid functn at the output layer, then use log loss functn (ML) and not least squares error.

(final) i/p into O/p layer

$$\log \tilde{P}(y) \rightarrow yz \Rightarrow \tilde{P}(y) = e^{yz}$$

$$P(y) = \frac{e^{yz}}{\sum_{y \in \{0,1\}} e^{yz}} \stackrel{=} \sigma((2y-1)z)$$

Parameterized by Bernoulli

Multi label Classification with Softmax: Generalizes sigmoid (σ).

$$\text{Softmax}(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

use feed

if "z" into softmax as the log of unnormalized prob. $z_i = \log \tilde{P}(y=i|x)$

In training, want to maximize log-likelihood i.e

$$\max \log P(y=i|z) = \max \log \text{softmax}(z_i).$$

$$\text{where } \log \text{softmax}(z_i) = z_i - \log \sum_j e^{z_j}$$

\Rightarrow i/p z_i will always have direct contributn to the cost function, and in principle this term cannot saturate (since it is linear w.r.t. of previous layers) learning can always progress !!

→ The ^{more} stable version of softmax is $\text{softmax}(z) = \text{softmax}(z - \max_i z_i)$.
 This is ∵ softmax is invariant to addition of some scalar
 i.e. $\text{softmax}(z) = \text{softmax}(z + c)$. This allows us to evaluate softmax with only small numerical errors when z contains extremely large #'s.

→ Sigmoid → Use the negative log-likelihood as functⁿ of z (last input) to output rather than $\hat{y} = p(y=i|x)$. This is so that sigmoid underflows to zero, and can avoid taking log of \hat{y}_i , which results to $-\infty$.

→ ReLU = $\max\{0, x\}$.

Generalizing ReLU.

$$h_i(a_i, \alpha_i) = \max(0, a_i) + \alpha_i \min(0, a_i).$$

\downarrow
non-zero slope for unit i "
when $a_i < 0$.

Three variations include:

① Absolute Value Rectification: fixes $\alpha_i = -1$ and get $z = |a_i|$.
 Used for object recognition from images, where it makes sense to learn features that are invariant under polarity reversal to input illuminatn's.

② Leaky ReLU: fixes α_i to very small value eg. 0.01

③ Parametric ReLU: treats α_i as learnable parameter.

→ During initializatⁿ, generally advisable to initialize wts with value 0.1, so that they can pass through the activatⁿ layer.

→ Max Out Units: Generalize ReLU's by dividing "a" into "k" ^{groups of} values. Each max out unit then o/p's max of one of these groups.

$$z = h_i(a_i) = \max_{j \in [1, k]} a_{ij}$$

→ Since max-out unit is parameterized by "K" w_i vectors instead of one, learning such units needs more regularization than std. ReLU.
Can work well in large data & # pieces per unit is low.
+ training

→ Since max-out units have multiple filters, it is difficult for them to forget how to perform tasks based on previous tasks, ~~is~~ \Rightarrow catastrophic forgetting.

Cost function: Usually avg. over training set.

$$J(\theta) = \mathbb{E}_{(x,y) \sim P_{\text{data}}} [L(f(x, \theta), y)],$$

where L is per sample loss based predicted output $f(x, \theta)$ and target y .

$$\mathbb{E}_{(x,y) \sim P_{\text{data}}} [L(f(x, \theta), y)] = \frac{1}{m} \sum_{i=1}^m L(f(x^{(i)}, \theta), y^{(i)})$$

→ Mini batch algorithms converge faster using approx.^m gradient rather than exact slower gradients.

Ill-conditioning of Hessian: When this is the case, SGD gets stuck, so that even a small step ↑'s the cost functn.

$$\begin{aligned} x &= x^0 - \varepsilon g. \\ (\text{2}^{\text{nd}} \text{ Order} \atop \text{Taylor series} \atop \text{approx}^m) \quad f(x) &\approx f(x^0) + (x - x^0)^T g + \frac{1}{2} (x - x^0)^T H (x - x^0) \\ &= \underbrace{f(x^0)}_{\text{orig. cost}} - \underbrace{\varepsilon g^T g}_{\text{expected improvement}} + \underbrace{\frac{1}{2} \varepsilon^2 g^T H g}_{\text{correct to account for curvature of functn.}}. \end{aligned}$$

→ The last term when too large can cause gradient to move uphill. i.e. 3rd term $> \varepsilon g^T g$. (2nd term)

→ When ill-conditioning, then learning becomes very slow despite of strong gradient, and one needs to ↓ the learning rate in order to ~~be~~ nullify the strong ↗ curvature.

Exploding Gradients: NN have steep regions in cost funct! Such a structure comes from multiplying many large wts.

→ When parameters get close to the cliff, a small update can throw the parameters very far, thus undoing all the optimization done before.

Gradient Clipping: Exploding gradients can be overcome by GC heuristic. In an update gradient does not specify step size, only the dirⁱⁿ, while SGD step size \propto gradient, but in clipping one can "cap" the step size to prevent from parameter moving away far.

Stochastic Gradient Descent

Require: learning rate ε_k , initial parameter θ .

while stopping criterion not met do

sample mini-batch of m examples from training data

Compute gradient $\hat{g} \leftarrow + \frac{1}{m} \nabla_{\theta} \sum_i L(\cdot)$

Update $\theta \leftarrow \theta - \varepsilon \hat{g}$.

→ In practice, it is common to decay the learning rate with some schedule. e.g. linear \downarrow until t iteratⁿ

$$\varepsilon_k = (1-\alpha)\varepsilon_0 + \alpha\varepsilon_t \quad \alpha = \frac{k}{t}.$$

→ Only use validatⁿ for hyperparameter tuning, no wts get update when you train on validatⁿ (i.e. no BP). If you update wts, then you most~~ly~~ probably invite overfitting!!

→ Monitoring validation error over the progress of learning and when we face error plateaus, \downarrow the learning rate, see what happens.

Momentum: Using it, one can accelerate learning when there is high curvature, small consistent gradients, or noisy gradients.

→ Aggregates exp. decay, gradients from past iterat's to continue stepping in that dirⁿ.

$$v \leftarrow \alpha v - \epsilon \nabla_{\theta} \left(\underbrace{\frac{1}{m} \sum_{i=1}^m L(\cdot)}_{\text{current gradient elements}} \right)$$

$$\theta \leftarrow \theta + v.$$

$\alpha \in [0, 1]$.

→ Momentum tries to solve the ① poor conditioning of Hessian
② variance in the ~~stochastic~~ stochastic gradient.

→ SGD with Momentum

I/P: G , α momentum parameter, θ , initial velocity v .
which stopping criteria not met do

sample "m" examples

$$\text{Compute GD: } g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\cdot), y)$$

$$\text{Compute velocity update: } v \leftarrow \alpha v - Gg$$

$$\text{Apply update: } \theta \leftarrow \theta + v$$

end while

→ Nesterov Momentum: Same as Momentum, but gradients are evaluated after applying current velocity i.e.

$$v \leftarrow \alpha v - \epsilon \nabla_{\theta} \left(\frac{1}{m} \sum_{i=1}^m L(f(\overset{*}{\theta} + \alpha v), y) \right)$$

difference b/w Momentum and NM.

$$\theta \leftarrow \theta + v$$

→ Can speed up convergence for convex gradients.

Approximating 2nd Order Methods.

$$J(\theta) \approx J(\theta_0) + (\theta - \theta_0)^T \nabla_{\theta} J(\theta_0) + \frac{1}{2} (\theta - \theta_0)^T H(\theta - \theta_0)$$

Solving for critical pt.

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0). \quad (\text{Newton Update})$$

→ The update is more accurate as Hessian captures/determines the curvature of the loss
local

Initialization's

→ Should do random initialization ∵ we don't know which deterministic algorithm can be used to update the wts of two hidden units (having same wts and activatn funct) and update will take place in same manner.

→ Large wt. Initialization's : ^(A) Helps in breaking symmetry and avoids losing signal during fwd propagatn.

^(D) Exploding gradients, both in FP & BP

^(D) Makes the activatn function saturated.

e.g. wt. decay,

→ Small wt. Initialization's : ^(A) Good from regularization perspective
^(A) they induce prior for final parameters being close to initial parameters

^(D) If we initialize with large values, then our prior specifies which units should interact with each other and how they should interact.

→ Wt. Initialization as a hyperparameter if you have good amount of computational power, then train on one minibatch and see activations and gradients and appropriately scale the wts.

→ Bias Initialization: Usually zero initialization works but can also do

① to match op's to marginal statistics

② avoid saturation eg. ReLU units bias should be initialized as 0.1 instead of 0, avoiding saturation at initializatⁿ.

Adaptive Learning: Learning rate should ↓ periodically.

① Adagrad - scales a fixed learning rate inversely proportional to root of sum of all historical values of gradients.

② RMSProp - changes gradient accumulatⁿ of Adagrad into an exp. wtd moving avg.

③ Adam - adds momentum to RMSProp.

① Adagrad : $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}, \theta), y^{(i)})$
gradient: Accumulate squared, $r \leftarrow r + g \odot g$ element wise

$$\text{Update} : \Delta \theta \leftarrow \frac{-\varepsilon}{\delta + \sqrt{r}} \odot g$$

Apply update: $\theta \leftarrow \theta + \Delta \theta$.

② RMS Prop: $g \leftarrow \text{same}$

Accumulate sqd gradient: $r \leftarrow \delta r + (1-\delta) g \odot g$

$$\text{Update} : \Delta \theta \leftarrow \frac{-\varepsilon}{\sqrt{\delta r + r}} \odot g$$

Apply update: $\theta \leftarrow \theta + \Delta \theta$.

③ ~~Adagrad~~ RMS Prop with Momentum

$$\tilde{\theta} \leftarrow \theta + \alpha v$$

$$g \leftarrow \frac{1}{m} \nabla_{\tilde{\theta}} \sum_i L(f(x^{(i)}, \tilde{\theta}), y^{(i)})$$

$$r \leftarrow \delta r + (1-\delta) g \odot g$$

$$\text{Velocity update} v \leftarrow \alpha v - \frac{\varepsilon}{\sqrt{r}} \odot g$$

$$\text{Update} \quad \theta \leftarrow \tilde{\theta} + v$$

Batch Normalizatⁿ : Method of adaptive-reparameterizatⁿ 20

→ Can be applied to any input or hidden layer and reduces the problem of coordinating updates across many layers.

→ H' is a minibatch of activations, each column represents the activation of a given unit layer while a row represents activation of an example.

$$\underbrace{H'}_{\text{normalized}} = \frac{H - \mu}{\sigma}, \text{ where } \mu \text{ & } \sigma \text{ are column wise means & std. deviat.}$$

$$\mu = \frac{1}{m} \sum_{i=1}^m H_i$$

$$\sigma = \sqrt{\cancel{\lambda}(s) + \frac{1}{m} \sum_{i=1}^m (H_i - \mu)^2}$$

usually 10^{-8} to avoid undefined gradients.

→ Can also be considered as whitening of data, since each dimension (unit of a layer) is de-correlated and all have a unit variance learning alg.

→ Normalizatⁿ makes BP efficient: Doesn't waste time to adjust mean and std deviat. Also it prevents scaling errors from propagating through the network since o/p's are normalized after normalizatⁿ layer.

→ For a layer of form $\phi(WX + b)$, where ϕ is the non-linear activatⁿ functⁿ, it is recommended to apply batch normalization on XW and remove "b", since bias is redundant with B of BN.

Usually $\gamma H' + \beta$, where γ & β are learned parameters.

Lecture 7 CNN I

(Regularization less skipped)

- Building invariance directly into the network is the basis for convolutional neural networks, widely applied to images
 - System should remain invariant to translations, scaling and small rotations and also invariant to some elastic deformations.
small nonlinear transform.
- Invariance: Invariant to translations, scaling and small rotations and also invariant to some elastic deformations.
small nonlinear transform.
- 3 Pillars of convolution in CNN!
- The convolution operation allows us to do this for i/p's of varying sizes and
- ① local receptive fields (sparse connectivity in range vs fully connected)
- ② parameter sharing (allows for detectn of same features over image)
- ③ equivariant representation and subsampling (reduces information hierarchy).
(i.e. downsampling !!?)
- ① local receptive fields imply sparse connectivity and interact. Though, the connectivity is sparse, the neurons still interact with each other after some depth and eventually some neurons which did not interact, will interact when we apply / add fc-layer at the end before o/p.
- ② Parameter sharing -

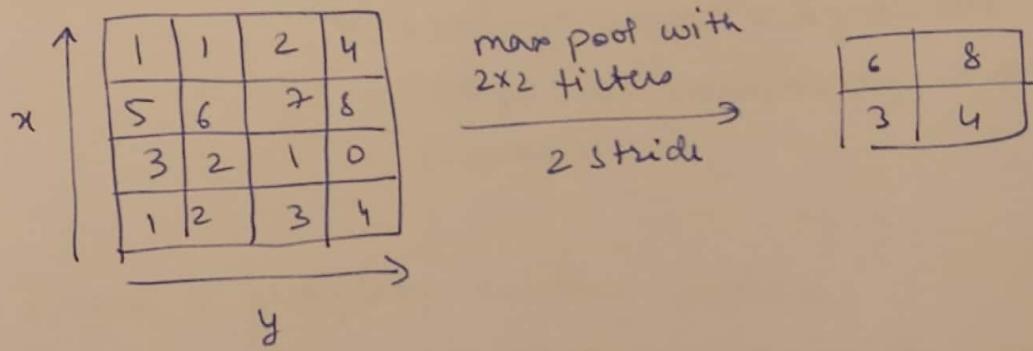
③ Translatⁿ equivariance - Change in I/P, \otimes changes O/P in the same manner. i.e. $f(g(x)) = g(f(x))$, f is equivariant to g . (21)

- Through Parameter sharing, convolutⁿ becomes spatially invariant.
- Suppose now that remapping is shift of pixel locatⁿ's i.e. $I' = g(I) = I(x-1, y)$. Translatⁿ equivariance means that applying convolution to the remapped I' is equivalent to first applying convolutⁿ to I and then applying the remapping $g(\cdot)$.
- Why is equivariance useful & and when is it not?
- Also prove 'the shift equivariance property of convolution'!!

- # Convolution variants: On each image / layer we generally apply multiple kernels in parallel. We can think of a kernel as a feature extractor.
- # Multiple Image Channel: Color image is 3-d tensor and thus $W \times H \times C$ channel. (R, G, B).
- # Stride: For saving computational expense skip some pixels. Can be thought of as downsampling in the O/P of fully connected convolutional funct'.
- # Zero Padding: Depends upon the size of kernel.

CNN Layer: 3 Stages

- ① Convolut' operat' ② non-linear activation eg. ReLU
- ③ pooling funct' (max, avg.) \rightarrow some summary statistic.
↓
Can also define strides..



- Pooling over spatial regions results in invariance to translation. Can also lead to other types of ~~eg~~ invariance when parameterizat' is changed.
- Pooling also allows handling input images of varying sizes. By changing the stride, the size of feature map ~~can be~~ will remain the same for images of varying sizes.
- Do an example of counting # parameters in the CNN architecture !!

Lecture 8. CNN II

(23)

→ Alex Net , VGGNet , Google LeNet , MSRA

Dense Net .

CNN Visualizatⁿ.

- ① Visualize filters
- ② Visualizing last layers (K-NN)
- ③ Multiple^{layer} Activations - Feature map visualizatⁿ.
- ④ Saliency Maps: How to tell which pixels matter for classification?
of (unnormalized) class

$\frac{\partial E}{\partial x}$ → Compute gradient score w.r.t.
image pixels, take absolute
value & max over RGB channels

⑤ Gradient Ascent

$$I^* = \text{argmax}_I f(I) + R(I)$$

Neuron value

Natural image regularizer

Gradient ascent: Synthes~~e~~^c
image that maximally activates
a neuron.

⑥ Deep Dream: amplify existing features.

Rather than synthesizing an image to maximize a specific neuron, instead try to amplify the neuron activations at some layer in a network.

Choose an image & a layer in CNN: repeat

- ① Fwd: Compute activation's at chosen layer
- ② Set gradient of chosen layer equal to its activation
- ③ Backward: Compute gradient on image
- ④ Update image

Equivalent to $I^* = \text{argmax}_I \sum_i f_i(I)^2$

CNN Applications

① Transfer Learning (Just like we did our project on synthesizing images).

- Take pretrained model.
- Remove the last fc-layer
- Use pretrained model as a fixed feature extractor for the new dataset
- Train a new classifier for the new dataset i.e linear SVM.

Also finetune all the layers of pretrained model. Also keep some earlier layers fixed and only fine-tune some higher-level layers as their impact in the loss function is more.

② R-CNN

③ Faster R-CNN.

④ YOLO - You only Learn Once

} Detecting objects
in the video.

Regularization → Loosely defined as modification to a learning algorithm which reduces its generalization error but not training error.

How & what to do?

- add constraints to the model parameters (θ)
- extra ~~cost~~ terms in the objective function (i.e. loss function)

→ Most strategies focus on regularizing the estimator which trades off increased bias for reduced variance

→ 3 Bias Variance Situations

(1) Model excludes true data generating data and underfits
∴ excluding bias

(2) model matches true data generating process

(3) model includes true data generating process also with other possible generating processes & overfits, thus allowing variance to dominate estimation error (instead of bias).

Most Probable, But try to shift to (2)

Parameter Norm Penalties: Add norm penalty $\mathcal{J}_2(\theta)$ to objective functn. i.e.

$$\mathcal{J}(\theta, X, y) = \mathcal{J}(\theta, X, y) + \alpha \mathcal{J}_2(\theta). \quad \begin{matrix} \text{hyper parameter (ideally diff. } \\ \text{for every layer)} \\ \text{but expensive)} \end{matrix}$$

Larger α means more generalization, smaller less.

- \mathcal{J}_2 only penalizes wts of affine transform & not bias (b).
- Regularizing bias (b) may induce ~~WF~~ → each wt. specifies how two variables interact, so fitting a wt. well requires both in variety of conditions, ∵ bias only requires single variable observing hence needs less data.

L²- Regularizer (Ridge regression). : $J_2(\theta) = \frac{1}{2} \|w\|^2$.

→ Also called wt. decay as it tries to decrease the wt. and in turn makes gradients smoother.

→ How does wt. decay work?

$$\tilde{J}(w, x, y) = J(w, x, y) + \frac{\alpha}{2} w^T w.$$

$$\nabla_w \tilde{J}(w, x, y) = \nabla_w J(w, x, y) + \alpha \cdot w$$

$$w \leftarrow w - \varepsilon (\nabla_w J(w, x, y) + \alpha w)$$

$$\equiv \underbrace{w(1-\alpha)}_{\text{This implies that at every update, the wt-decay}} - \varepsilon \nabla_w J(w, x, y)$$

regularization performs a multiplicative shrink.

→ Simple wt. decay is inconsistent with certain scaling properties

of network.

$$\xrightarrow{\text{1/p}} \xrightarrow{\text{O/p}}$$

→ 2 layer perceptron, $\{x_i\}$ $\{y_k\}$

$$z_j = h\left(\sum_i w_{ij} x_i + w_{j0}\right) \quad \text{and} \quad y_k = \sum_j w_{kj} z_j + w_{k0}.$$

→ Suppose we add linear term $x_i \rightarrow x'_i = ax_i + b$.

$$w_{ji} \rightarrow w'_{ji} = \frac{1}{a} w_{ji} \quad \text{and} \quad w_{j0} \rightarrow w'_{j0} = w_{j0} - \frac{b}{a} \sum_i w_{ji}$$

Why for O/p's $y_k \rightarrow y'_k = c y_k + d$, we can adjust wt's of 2nd layer

$$w_{kj} \rightarrow w'_{kj} = c w_{kj} \quad \& \quad w_{k0} \rightarrow w'_{k0} = c w_{k0} + d.$$

→ Network Consistency: Model is ^{equi}~~invariant~~ (i.e. if networks are consistent, then the wt's should differ only by linear transform).

→ Simple wt. decay gives same importance to bias and wt's and prefers "smaller weights".

→ Consistent way to avoid inconsistency of L_2 -regulariz⁽²⁻⁸⁾

$$\frac{\lambda_1}{2} \sum_{w \in W_1} w^2 + \frac{\lambda_2}{2} \sum_{w \in W_2} w^2, \text{ where } W_1, W_2 \text{ are wts of layers 1 \& 2 (excluding bias)}$$

→ By doing this, wts are invariant to scaling.

$$\lambda_1 \rightarrow \lambda'_1 = \alpha^{-1} \lambda_1 \quad \lambda_2 \rightarrow \lambda'_2 = \alpha^{-1} \lambda_2$$

$$p(w | \alpha_1, \alpha_2) \propto e^{(\frac{\alpha_1}{2} \sum_{w \in W_1} w^2 + \frac{\alpha_2}{2} \sum_{w \in W_2} w^2)}. \text{ where } \alpha_i's \text{ are HP.}$$

New regularizer corresponds to a prior.

→ Priors of this form cannot be normalized since bias is unconstrained. To break the symmetry, we do shift-invariance property.

L_1 -regulariz $\therefore J(\theta) = \|w\|_1 = \sum_i |w_i|$

$$\tilde{J}(w, x, y) = \alpha \|w\|_1 + J(w, x, y)$$

$$\nabla_w \tilde{J}(w, x, y) = \alpha \cdot \text{sign}(w) + \nabla_w J(w, x, y).$$

In general L_1 -regulariz gives sparse wts, where some wts have values close to zero.

→ ~~If θ follows Gaussian prior~~ L^2 regulariz is equivalent to MAP with a Gaussian prior of wts. $\underbrace{\log p(\theta | x)}_{\text{cost.}} + \underbrace{\log p(\theta)}_{\text{regulariz term}}$

→ L^1 is equivalent to using an isotropic Laplace distribution over "w" for prior.

Tangent Propagation: (can also regularize model opp's (make invariant to transformations), by tangent propagation. continuous)

→ If we a transformation on input vector x_n such as (rotation, translation), then the transformed pattern follows a manifold ~~in~~ M with D-dim. input space

→ The transformⁿ is parameterized by ξ . The manifold swept by x_n is a 1-D space with parameter ξ_s .

→ $s(x_n, \xi) =$ resultant vector with $s(x, 0) = x$. The tangent to M is directional derivative $\tau = \frac{\partial s}{\partial \xi_s}$

$$\tau_n = \left. \frac{\partial s(x_n, \xi)}{\partial \xi_s} \right|_{\xi_s=0}$$

→ The O/p vector will also change due to transformⁿ of i/p.

$$\left. \frac{\partial y_k}{\partial \xi_s} \right|_{\xi_s=0} = \sum_{i=1}^D \left. \frac{\partial y_k}{\partial x_i} \times \frac{\partial x_i}{\partial \xi_s} \right|_{\xi_s=0} = \sum_{i=0}^D J_{ki} T_i$$

↓
Jacobian

$$\tilde{E} = E + \lambda \mathcal{R} \quad (\text{this encourages local invariance on neighborhood of pt } i)$$

$$\mathcal{R} = \frac{1}{2} \sum_n \sum_k \left(\left. \frac{\partial y_k}{\partial \xi_s} \right|_{\xi_s=0} \right)^2 = \frac{1}{2} \sum_n \sum_k (J_{ki} T_i)^2$$

where λ determines the trade-off b/w fitting & learning the invariance.

Reducing Model Capacity.

→ Early Stopping: To prevent overfitting. Good form of hyperparameter selection. Disadvantage is that it is computationally expensive and you also take validation data also as a hyperparameter. Also retraining the model all parameters get change and in the end as we do random initialization, so one can tell what will happen when if you train your model again.

→ One can also think it as restricting optimization procedure to a relatively small vol^m of parameter space

→ Parameter Sharing: CNN.

Adjusting Training Data.

→ Data Augmentation: Synthesize more data and train NN on this data.

→ Injecting Noise:

① Noisy I/P: Train your NN with noise-injected I/P:
i.e. add random noise to I/P's.

② Noisy wts: Noise added to the wts (RNN's).
Consider model wts as R.V.'s associated with some probability dist.

Noise as a Regularizer.

$$J = \mathbb{E}_{P(x,y)} [(\hat{y}(x) - y)^2]$$

Now add noise $\epsilon_w \sim N(\epsilon, 0, \eta I)$.

$$\begin{aligned} \tilde{J}_w &= \mathbb{E}_{P(x,y, \epsilon_w)} [(\hat{y}_{\epsilon_w}(x) - y)^2] \\ &= \mathbb{E}_{P(x,y, \epsilon_w)} [\hat{y}_{\epsilon_w}^2(x) - 2\hat{y}_{\epsilon_w}(x)y + y^2]. \end{aligned}$$

For η = small, $\min J = \min \tilde{J}$, with added regularization term $\eta \cdot \mathbb{E}_{P(x,y, \epsilon_w)} [\|\nabla_w \hat{y}(x)\|]$.

③ Noisy Output. → Hard classification problem becomes soft classification problem with correct being $1-\epsilon$ and wrong with ϵ .

↓ One hot coding					$(1, 0)$
<u>1 0 0 0 0</u>					$(1-\epsilon, \epsilon)$
<u>$1-\epsilon \epsilon_4 \epsilon_{14} \epsilon_{14} \epsilon_6$</u>					$(\epsilon, 1-\epsilon)$

- Bagging - can reduce generalization error by combining ensemble of models. i.e. multiple models trained separately and then votes on O/P for a given test sample
- The ensemble performs at least as well as any member on an average and if members make errors independently, then the ensemble performs better! Error correlatn of members cannot be helped by ensemble.
- Train different NN on different sampled dataset's sampled from training data.
with exp. large # models.
- # DropOut : Approximates the process of bagging (i.e. learn K-different models from K-different datasets).
- For each sample, create a mask and apply it to a layer.

$$\text{Dropout } 0.6 \rightarrow \frac{6}{10} \rightarrow 0.6$$

$$\text{Mask } 0.6 \rightarrow \frac{6}{10} \rightarrow 1.] \underbrace{\text{Masking value 1.}}_{\text{Hiding the wt.}}$$

DROPOUT

- Models share parameter
- Parameter sharing makes the representation of an exp. # models easy (less memory)
- Most ^{not} ~~other~~ models are trained explicitly as parameter sharing helps.

BAGGING

- Models are independent
- Parameters are more due to models being independent. (Hence more memory)
- Each model is trained to its convergence
- Both follow bagging algorithm !!

→ Inferenceing in Dropout.

$$\text{P}_{\text{bagging}} = \frac{1}{K} \sum_{i=1}^K p^i(y|x).$$

↗ dist. given mask vector u and data

$$\text{P}_{\text{dropout}} = \sum_u p(u) \underbrace{p(y|x,u)}_{\substack{\downarrow \\ \text{prob dist used} \\ \text{to sample } u \text{ for training}}}.$$

↙ This term is intractable, but 20-30 masks suffice.

② Use geometric mean as dropout. (Don't know, don't care)

③ wt. Scaling Inference Rule. at layer i

- ④ Basically, dropout = 0.5, is equivalent to doubling the wts of that layer and doing traditional ~~no~~ training!! No need of masking!!.

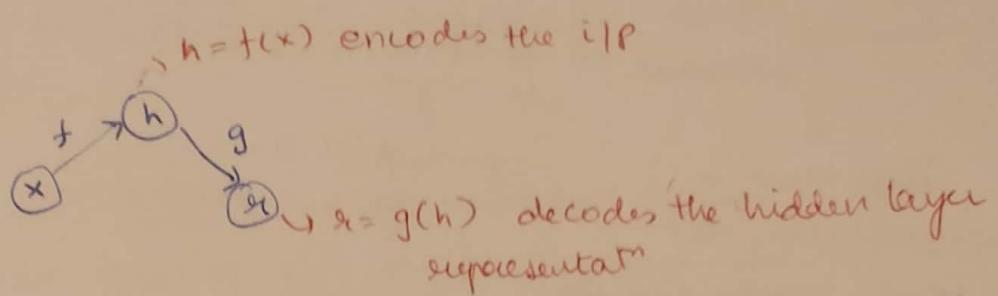
→ Effectiveness of Dropout.

- Computationally cheap ($O(n)$) computat' per ex. per update.
- Model Agnostic.
- Dropout reduces the effective c.c. of model, hence one should ↑ the capacity of model to compensate.
- Alternative view as a regularizer: Regularizes each HU to be not merely a good feature but a feature that is good in many contexts.

→ Dropout as injecting noise at HU: adding noise at feature extraction helps to enforce use of multiple features for prediction rather than relying on handful of features.

Skipped (RNN-I, II). Lecture -II Autoencoders.

- Special types of NN which are trained to "copy" the input to output.



- Based on def", it is simply 'a series of identity functors', but is a trivial setup, but we restrict h e.g. ^{its} dimensionality, so that cannot copy the i/p.
- We are interested in learning the hidden layer representation of i/p.
- Used largely for dimensionality reduction and feature learning, but with their relation with latent variables, also used for generative modelling.

Discriminative Model : x-data, y-label.
 ↴ MAP Inference task to find y^* maximizing $P(y|x)$.
 Models posterior directly & determines class label accordingly

Generative Models: Models the class conditional prob. $P(x|y)$ and uses prior $P(y)$ to find posterior $P(y|x)$.
 ↴ to generate new data via $P(x) = \sum_y P(x|y)P(y)$

(Class label determined by) $y^* = \operatorname{argmax}_y \frac{P(x|y)P(y)}{P(x)}$
 ↴ y what we want to learn and it's more difficult to come with

→ Almost any generative model with a latent variable and the ability to inference on latent variable given i/p can be inferred as an autoencoder

→ Encoder & decoder as stochastic mappings with $P_{\text{encoder}}^{(h|x)}$ and $P_{\text{decoder}}(x|h)$.

Undercomplete AE - constrain the hidden layer to be of lower dimension than that of x , try to capture salient features of x .
 $L(x, f(g(x))) \rightarrow \text{loss functn}$

→ If decoder "g" were linear and L was the MSE, then an undercomplete AE would span the same subspace as PCA !!

Complete & Overcomplete AE -

Complete = dim of i/p = dim of hidden layer

Overcomplete = dim of i/p < dim of HL. ($\frac{\text{hidden}}{\text{code}}$)

Even with little capacity can learn to copy data

Regularized AE : Modify the loss functn which encourages the model to learn other properties than copying. like

- ① Sparsity of representations
- ② Smallness of the derivative of representations
- ③ Robustness of noise or to missing i/p's.

Sparse AE → Adding the sparsity regularization term $\Omega(h)$ on code layer "h" to loss.

$$L(x, g(f(x))) + \Omega(h).$$

→ Used for learning features of other tasks eg. Classification.

→ In other regularizations such as wt. decay, the regularizer depends upon the prior on model parameters. i.e. maximizing log-likelihood $\log P(x|\theta) + P(\theta)$.

→ While there is no interpretation to regularization in AE, as the term depends upon the data.

→ $\mathcal{S}(h)$ is not a penalty for copying task, rather it is an approx^m to ML training of a generative model with latent variables.

$$P_{\text{model}}(x, h) = P_{\text{model}}(x|h) P_{\text{model}}^{(h)} \quad \begin{matrix} \rightarrow \text{prior over latent variables} \\ \text{model} \end{matrix}$$

λ represents beliefs before seeing "x".

$$\Rightarrow \log P_{\text{model}}(x) = \log \sum_h P_{\text{model}}(x|h).$$

→ The term $\log P_{\text{model}}^{(h)}$ can be sparsity inducing.

For ex., with Laplacian prior

$$P_{\text{model}}^{(h)} = \frac{\lambda}{2} e^{-\lambda \|h\|_1}$$

$$\begin{aligned} \text{If } \mathcal{S}(h) = \lambda \sum_i |h_i| \text{ we get } -\log P_{\text{model}}^{(h)} &= \sum_i (\lambda |h_i| - \log \frac{\lambda}{2}) \\ &= \mathcal{S}(h) + \text{const.} \end{aligned}$$

→ The result of applying $P_{\text{model}}^{(h)}$ on approx^m ML learning, is a result of model's distribution over its latent variables

→ This explains why features learned by the AE are useful, since they describe latent impacts on i/p.

Representation of Depth → Training with one HL is shallow and we cannot add arbitrary constraints while preserving sufficient representation. Adding depth helps ∵ it decreases the computational cost of representing some complex function.

↳ Also ↓'s the amount of training data needed to learn some representation of a function.

Stochastic Encoders & Decoders.

→ Minimizing $-\log P_{\text{model}}(x|h)$ at decode.

→ Minimizing negative log likelihood is equivalent to using MSE.

$$P_{\text{encoder}}(h|x) = P_{\text{model}}(h|x)$$

$$P_{\text{decoder}}(x|h) = P_{\text{model}}(x|h).$$

→ Binary classification \Rightarrow Bernoulli dist. from Lec-2
 Multiclass classification \Rightarrow Softmax

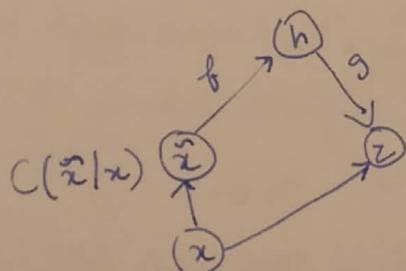
Regularize by Denoising AE.

→ Can also regularize by changing the reconstruction error.

→ To prevent eventual identity funct learning.
 the denoising AE minimizes

$$L(x, g(\tilde{f}(x))), \quad \tilde{x} = \text{noisy version of } x.$$

→ The AE will undo the noise instead of copying ip to op.



- ① Sample training ex: "x" from training data
- ② Sample \tilde{x} from $C(\tilde{x}|x)$
- ③ Learn f & g, using (x, \tilde{x}) where $P(x|\tilde{x}) = P_{\text{decoder}}(x|h=f(\tilde{x}))$, where P_{decoder} is determined by $g(h)$.

→ Now, the injected noise in data can be thought learning the manifold of the underlying probability distribn $p(x)$.

- Core assumption of AE is that data is concentrated around low-dim. manifold. Goal is to learn the structure of manifold.
- Tangent plane in manifold - For some pt x in d -dim M , the tangent plane is given by "d" basis vectors spanning the local dirⁿ's of variations allowed on manifold.
- If data is concentrated in low-dim manifold, then only variations tangent to the manifold around " x " need to correspond the changes in $h=f(x)$.
- Hence, encoder is mapping sensitive changes along the manifold dirⁿ's and insensitive to orthogonal dirⁿ's.

How to do Manifold Learning?

- Embedding: a low-dimensional vector with less dim than input space, of which the manifold is a low-dim. subset.
- Many manifold learning alg., learn the embedding of each training example.
 - The fundamental problem with non-parametric approaches to manifold learning is that manifold themselves are not smooth, hence need a lot of training data to cover variations
 - Because such method estimates the shape of manifold by interpolating b/w neighbouring ex., there is no chance to generalize ^{for} unseen data variations.
 - Alternatively, learning a mapping which can map any pt. from i/p to embedding in hope that it would generalize better for unseen variations than by doing local interpolat.

Contractive Autoencoders: Regularize the code $h = f(x)$ with $\mathcal{R}(h)$ to encourage derivatives of " f " to be as small as possible. (38)

$$\mathcal{R}(h) = \lambda \left\| \frac{\partial f}{\partial x} \right\|_F^2$$

$F = \text{Frobenius norm}$

→ $\mathcal{R}(h)$ encourages feature encoder to resist small change to the input, which in turn maps a neighbourhood of i/p pts to a smaller o/p neighborhood - Contractive AE.

$$L(x, g(f(x))) + \lambda \left\| \frac{\partial f}{\partial x} \right\|_F^2$$

$\underbrace{\quad}_{\text{reconstruction error -}} \quad \underbrace{\quad}_{\text{contractive penalty -}}$

$\underbrace{\quad}_{\text{Only learn identity functr}}$

features learned would be constant w.r.t. x

Applicatⁿs: AE

- ① Dimensionality Reduction
- ② Information Retrieval
 - ↳ Semantic Hashing
 - ↳ Embed data and query in low dimension.
- ③ As an "initializⁿ", AE can also be used for "pre-training"
 - ⓐ Pretrain a sequence of shallow AE, greedily, stacking one layer at a time with unsupervised ~~all~~ data
 - ⓑ train last layer supervised, keep others fixed
 - ⓒ BP & fine-tune the entire network supervised.

89

Lecture 13 GANS
Deep Generative Models

Why Study Generative Modeling

- Training and sampling from such models tests out ability to represent and handle high-dim. dist.
- helps dealing with multi-modal data learning.
- Simulate possible outcomes to be incorporated to reinforcement learning.
- Loosely define generative modeling as capturing $P(x)$ for some complex x , but the more correlated data, it becomes difficult to learn.
- Hence latent variables (unobserved var.) can capture the variations of one or more observed variable
- Ex. last lecture, each dim. of the manifold can be considered as latent variables.

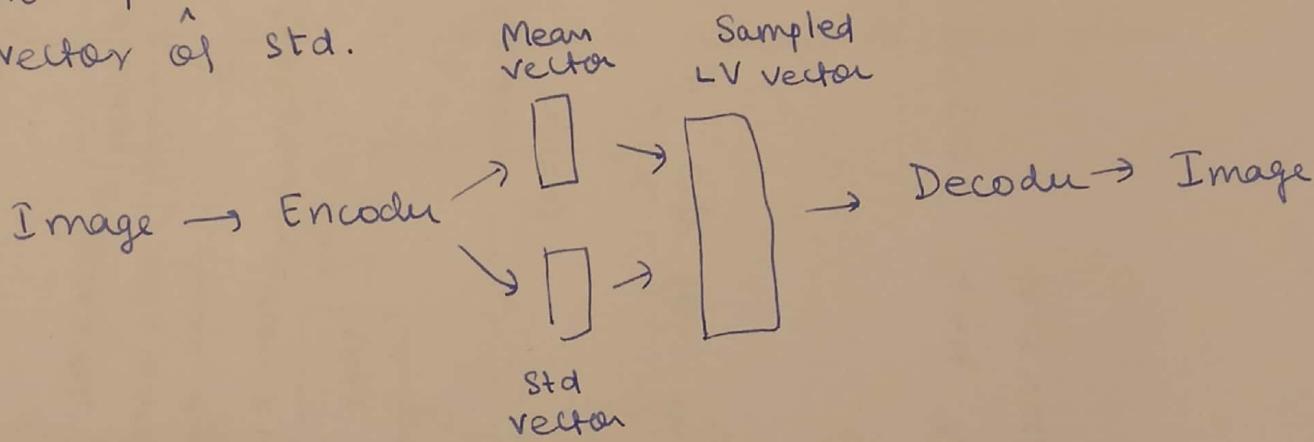
Latent Variable Generative Model : learn mapping from some latent variable to a complicated dist. on x

$$p(x) = \int p(x, z) dz \quad \text{where} \quad p(x, z) = \underbrace{p(x|z)}_{\text{difficult}} \underbrace{p(z)}_{\text{simple}}$$

Variational AE: Take initial vector as ^avector on real #, which can be seen has latent variable and the code (h) is some kind of a representation of the image which can synthesize them.

Image \rightarrow Encoder \rightarrow Latent Variable \rightarrow Decode \rightarrow Image

- But we don't know how to create latent vector. For this we add a constraint s.t. latent variable's match a unit Gaussian dist.
- This differentiates variational AE from std. AE.
- For our loss term, we sum two separate losses: ① generative loss ② latent loss
- ① MSE of reconstruction error
 - ② KL-divergence how closely to LV match a unit Gaussian.
- In order to optimize KL-divergence, we need to apply simple reparameterization trick: instead of the encoder generating a vector of real values, it will generate a vector of means & a vector of std.



GAN: Generative Adversarial Networks

- Based on game-theory scenario, where two networks compete
- Generator - Generates samples $x = g(z, \theta^{(g)})$ with random noise z .
- Discriminator - differentiates b/w real and fake,
 $d(x, \theta^{(d)})$
- Generator Network: $G(z, \theta^{(g)})$.
- Must be differentiable
 - No invertibility requirement
 - Trainable for any size "z"
 - Can make "x" conditionally Gaussian given "z".

Training a GAN: GAN has a binary entropy objective. (41)

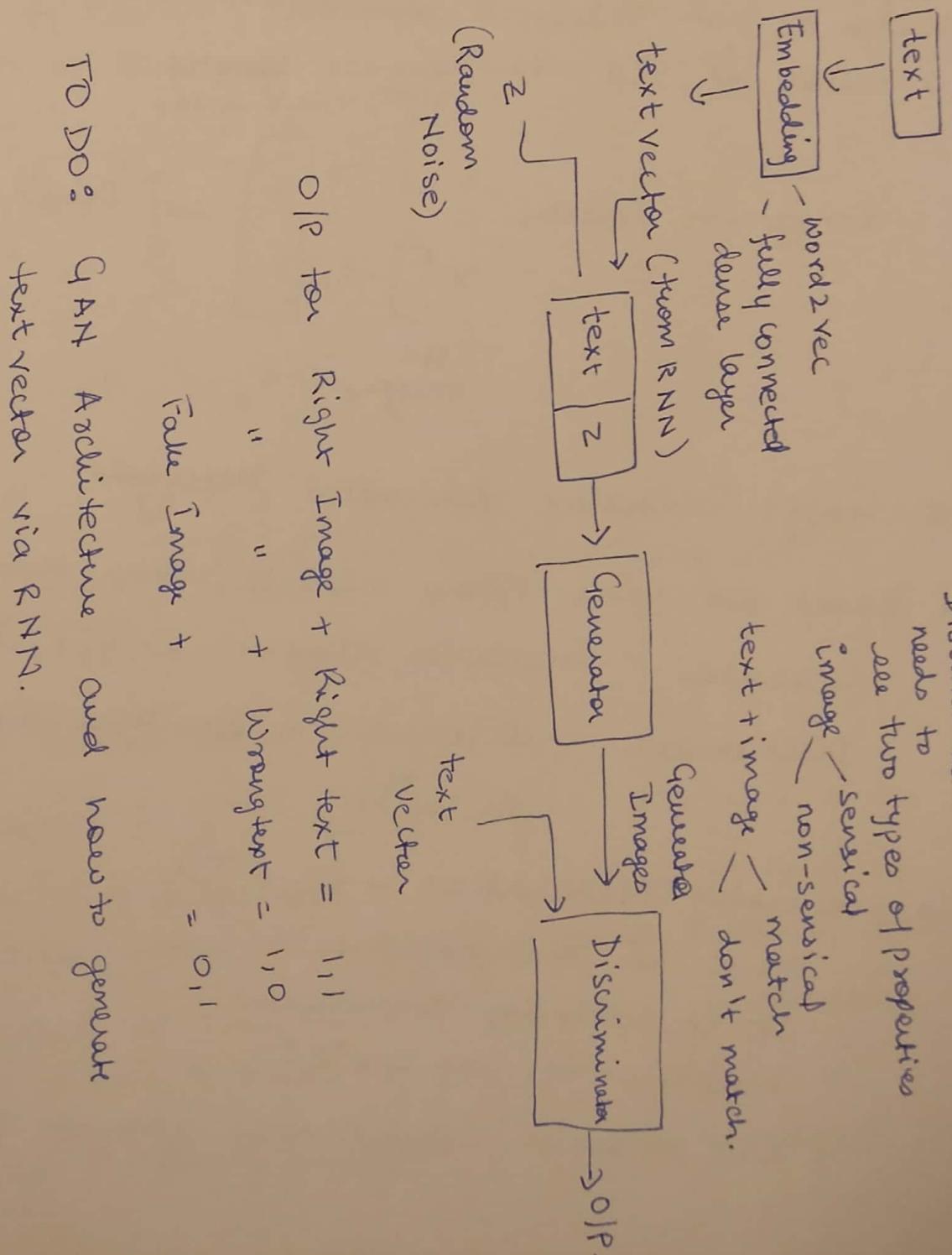
$$J_{GAN} = \log d(x, \theta^{(a)}) + \log [1 - d(g(x, \theta^{(g)}), \theta^{(a)})].$$

Training alternates b/w

$$\theta^{(g)*} = \operatorname{argmin}_{\theta^{(g)}} J_{GAN} \quad \text{and} \quad \theta^{(d)*} = \operatorname{argmax}_{\theta^{(d)}} J_{GAN}.$$

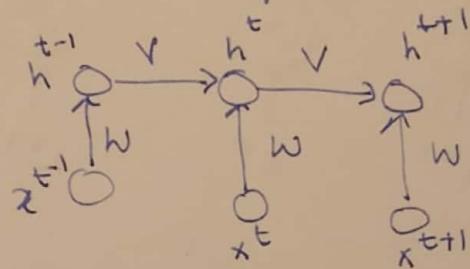
→ Convergence when generator samples are no longer distinguishable from real data. i.e. d outputs = 1/2.

ICML (2016): Generative Adversarial for Image Synthesis.

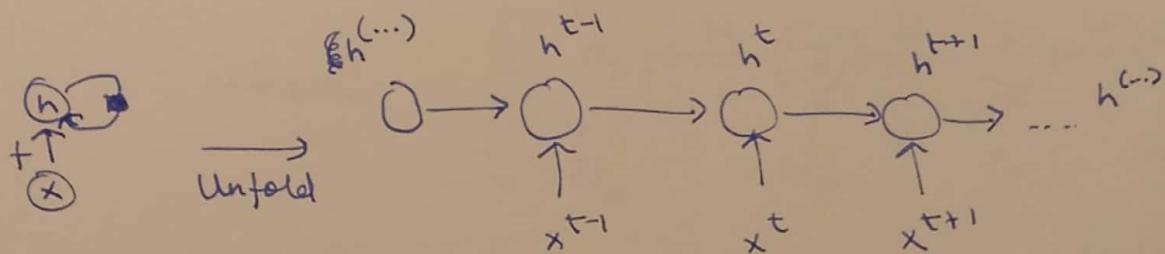


- The relaxation of the allowing of cycles in the network gives rise to Recurrent Neural Nets (RNN).
- The MLP can only form a map from i/p to o/p layers (vectors). But RNN can map from entire history to each o/p. This gives RNN a memory of previous i/p's.

Parameter Sharing is essential for RNN, because if the parameters were time specific, models could not generate the sequence of lengths not seen in training, nor share any statistical strength across diff. length of sequences across different times.



$$\rightarrow \text{Recursion: } s^{(t)} = f(s^{(t-1)}, \theta).$$



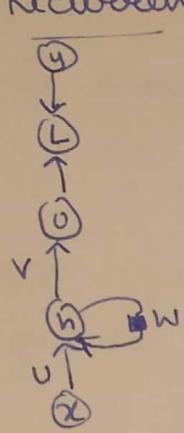
→ Recursiveness in hidden units

$$h(t) = f(h^{(t-1)}, x^{(t)}, \theta) = g(\dots).$$

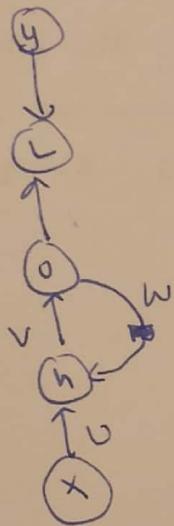
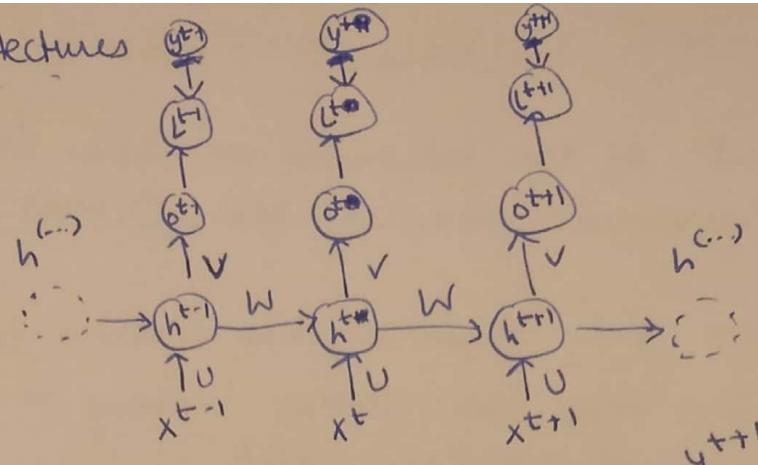
- (A) t is sized, \therefore it is specified in terms of transition from one state to another, rather than variable-length history of states of g
- (A) f can be used at any time "t", with same parameters, which "g" is time specific.

⇒ Needs less # training examples to estimate the model.

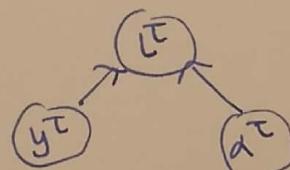
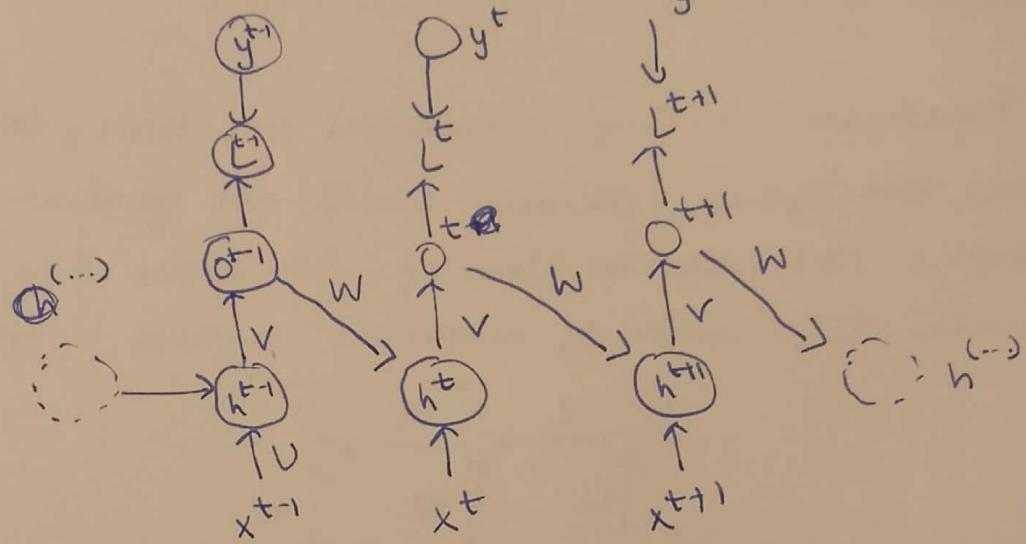
→ Recurrent Architectures



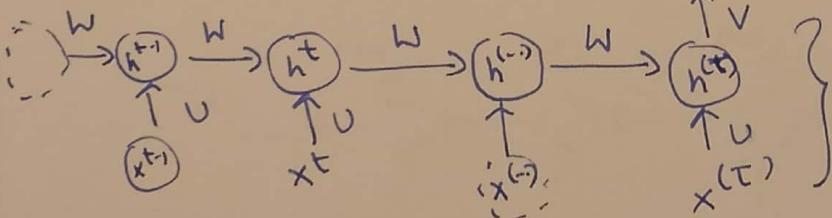
→



→



Encoder

Update:

$$a^{(t)} = b + W h^{(t-1)} + U x^{(t)}$$

$$h^{(t)} = \tanh(a^{(t)})$$

$$o^{(t)} = c + V h^{(t)}$$

$$\hat{y}^{(t)} = \text{softmax}(o^{(t)})$$

→ $L^{(t)}$ as negative log likelihood of $y^{(t)}$ given $x^{(1)}, \dots, x^{(t)}$, then

$$L(\{x^{(1)}, \dots, x^{(T)}\}, \{y^{(1)}, \dots, y^{(T)}\}) = \sum_{t=1}^T L^{(t)}.$$

$$= - \sum_t \log P_{\text{model}}(y^{(t)} | \{x^{(1)}, \dots, x^{(t)}\})$$

given by $y^{(1)}$ from $\hat{y}^{(1)}$

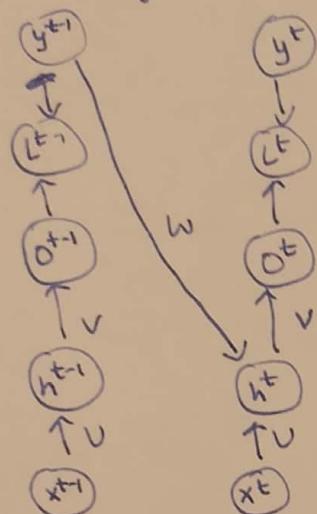
Back Propogat^m through Time (BPTT) Derivation.

(44)

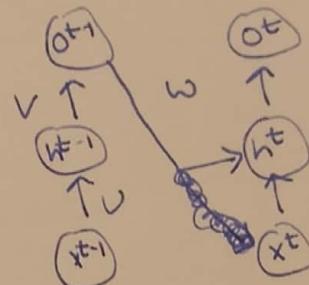
Teacher Forcing: Consider max likelihood criterion for training, where the model receives the true label $y(t)$ as i/p at $t+1$.

$$\log P(y^{(1)}, y^{(2)} | x^{(1)}, x^{(2)})$$

$$= \log P(y^{(2)} | x^{(1)}, x^{(2)}, y^{(1)}) + \log P(y^{(1)} | x^{(1)}, x^{(2)}).$$



Train time



Test time

→ Due to this, training can be done in parallel as we know the labels first hand and they don't depend upon time.

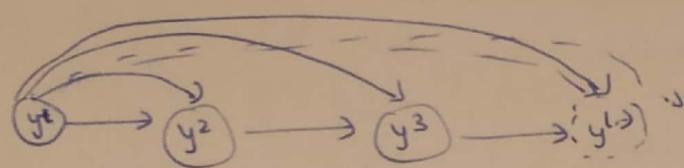
RNN's as Directed Graphs

→ If we consider the loss as max. log likelihood, and consider an RNN as a probabilistic graphical model, the graph would have no edges from any past $y^{(i)}$ to current $y^{(t)}$. Furthermore, y are all conditionally independent, given $x^{(i)}$.

→ But if the network has connections b/w previous o/p & current o/p, then log likelihood

$$P(y_0^{(t)}, x^{(1)}, \dots, x^{(t)}, y^{(1)}, \dots, y^{(t-1)}).$$

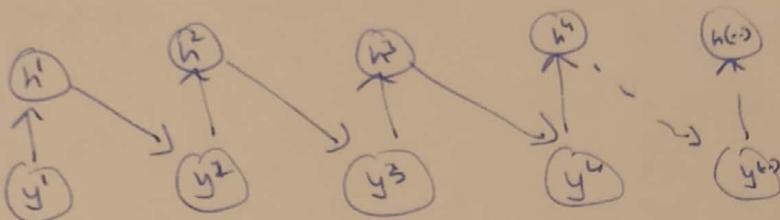
In such a case $y^{(i)}$ does not have connection with $y^{(j)}$, $\forall i < j$.



More Parameters

- What would happen if we introduce a hidden units into such an RNN.

$$\Rightarrow h^{(t+1)} = f(h^{(t-1)}, y^{(t)})$$



Less Parameters.
Lossy Summary.

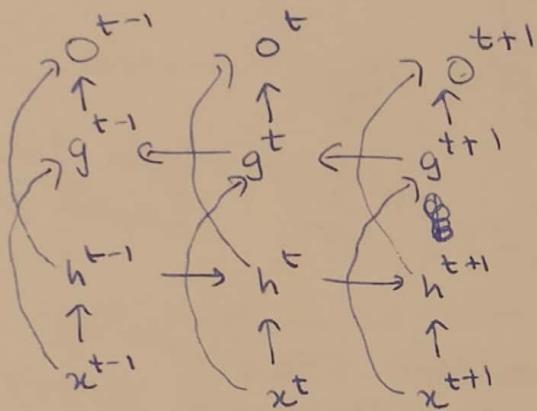
- Good parametrization doesn't give you free lunch!!
That is, it takes more time to train. (Because of BPTT)
↓
cumbersome nature.

makes an

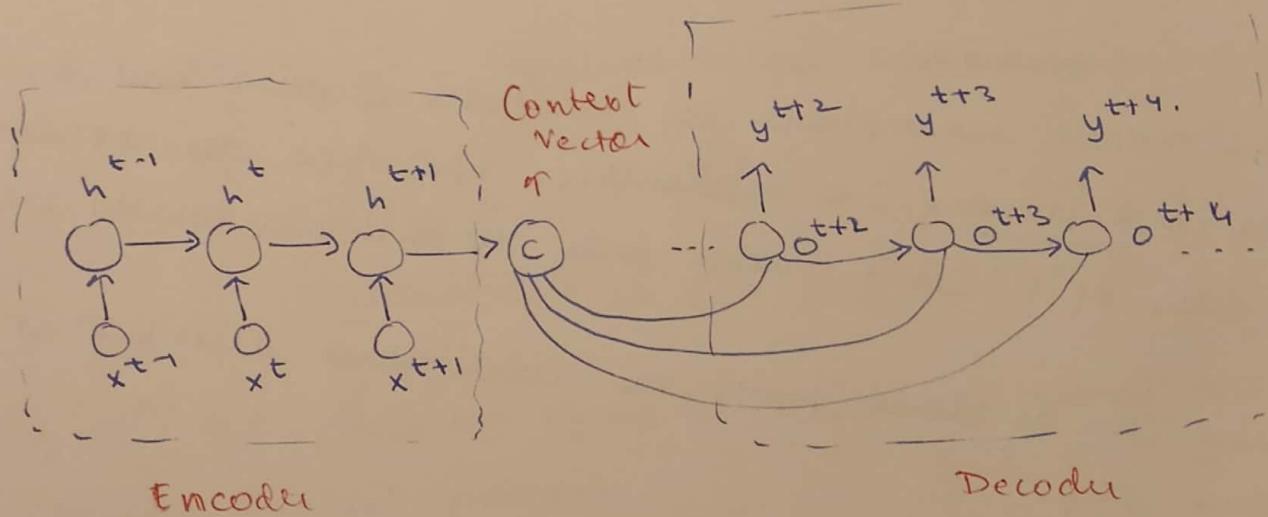
- Also the model's assumption by parameter sharing that the sequence is stationary i.e. relationships b/w steps does not depend on the value of "t" itself, rather only the difference b/w the steps.
- Sampling can be done from the conditional distribution at each time step.

But when do you know when to stop.

- Always beneficial to know the length of the sequence.
- Can add, extra "stop" dimension in output (i.e. y)
- introduce extra Bernoulli op at hidden layer/context layer which decides to continue / halt generative process
- add an extra op in model to predict T, which samples the value for T and then sample for T steps.

Bi-directional RNN.

Seq-to-Seq via Encoder-Decoder

# Long term Dependencies (LSTM, GRU, Peephole)

Refer to Colah blog for better understanding !! :)

→ Learning long-term dependencies is non-trivial:

- ① gradients tend to vanish or explode over many time steps
- ② Even if we have stable network, long term dependencies / interactions, tend to get exponentially less wt.'s (as Jacobians are multiplied repeatedly) in comparison to short term memories

Leaky Units

$$h^{(t)} = + (h^{(t-1)}, x^{(t)}, \theta).$$

$$h_L^{(t)} = \alpha h_L^{(t-1)} + (1-\alpha) h^{(t)} + (h_L^{(t-1)}, x^{(t)}, \theta). \text{ where } \alpha \in [0, 1].$$

Gated Recurrence Unit

→ GRU computes reset gate $r^{(t)}$ and update gate $z^{(t)}$.

$$h^{(t)} = g_R(h^{t-1}, x^t, \theta_R)$$

$$z^{(t)} = g_z(h^{t-1}, x^t, \theta_z).$$

Reset gate determines how much of the old memory should propagate (or kept or memorized)

$$\tilde{a}^t = b + r^t W h^{t-1} + U x^t \quad \tilde{h}^{(t)} = \tanh(\tilde{a}^t).$$

Now, update gate is applied to the memory content & the previous step.

$$h_q^{(t)} = \underbrace{z^t}_{\text{decides how much past state should matter}} h_q^{(t-1)} + (1-z^t) \tilde{h}^{(t)}$$

decides how much past state should matter

→ If $r^{(t)} \approx 0$, then we ignore the previous information and only consider current memory.

→ Has one advantage of forgetting information than the Leaky units.

LSTM (Long-Short Term Memory) Unit

3 Gates: an input gate $i^{(t)}$, a forgetting gate $f^{(t)}$, o/p gate $o^{(t)}$.

$$i^{(t)} = g_i(h^{(t-1)}, x^{(t)}, \theta_i)$$

$$f^{(t)} = g_f(h^{(t-1)}, x^{(t)}, \theta_f)$$

$$o^{(t)} = g_o(h^{(t-1)}, x^{(t)}, \theta_o)$$

→ Now, we define memory cell $c^{(t)}$ by gating previous memory and new content.

$$a^{(t)} = b + Wh^{(t-1)} + Ux^{(t)} \quad c^{(t)} = f^{(t)}c^{(t-1)} + i^{(t)}\tanh(a^{(t)}).$$

→ The final hidden state is a gated activation of memory cell:

$$h_{\text{LSTM}}^{(t)} = o^{(t)}\tanh(c^{(t)}).$$

Peephole Connections

→ Allow gate layers to look at cell state.

$$i^{(t)} = g_i(h^{(t-1)}, c^{(t-1)}, x^{(t)}, \theta_i).$$

$$f^{(t)} = g_f(h^{(t-1)}, c^{(t-1)}, x^{(t)}, \theta_f)$$

$$o^{(t)} = g_o(h^{(t-1)}, c^{(t-1)}, x^{(t)}, \theta_o).$$

→ Intuitively, taking into account what is in the state, when deciding which values to forget/update.