

COL 865

Deep Learning

August 17, 2017

Logistics:-

Monday August 14, 2017

CNNs (advanced)
By Chetan Arora

4.5 weeks:- Basics / Foundation

Thursday August 7, 2017

CNN: Basics

5.5 weeks

1.5 weeks:- CNNs

1.5 weeks:- RNNs

1.5 weeks:- GANs / Autoencoders

1 week:- Deep RL

4 weeks

Student Presentation

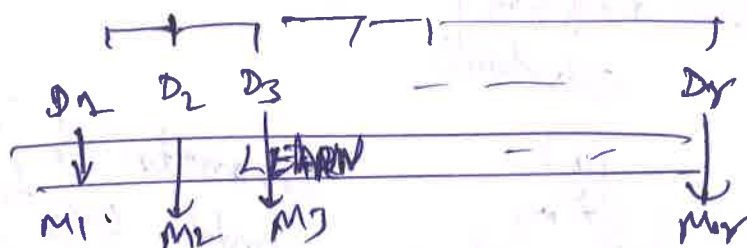
Ensemble Methods:-

Combine multiple different models & their outputs.

Example:-

Bagging

Data: $\{x_i, y_i\}_{i=1}^m$



M_i :- $P(y|x_i; \theta_i)$

Then, final prediction:-

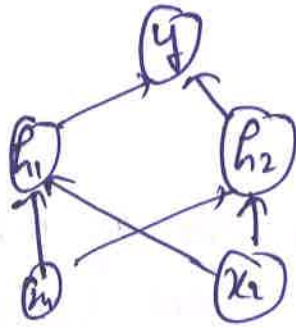
$$\frac{1}{n} \sum_{i=1}^n P(y|x_i; \theta_i)$$

In general can be:-

- ① Bagging
- ② Random initialization
- ③ Random selection of mini-batches

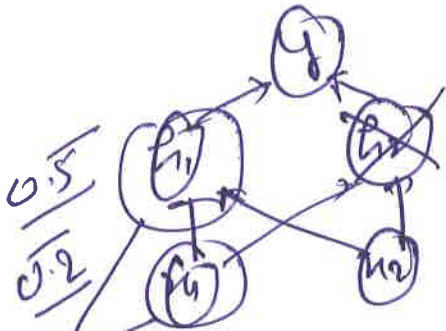
④ Dropout

Next



4 non-output units -
for each such unit, consider dropping it out independently.
16 possible networks

example -



Different probabilities of being dropped out

Note: explicitly constructing all such sub-networks may be expensive.

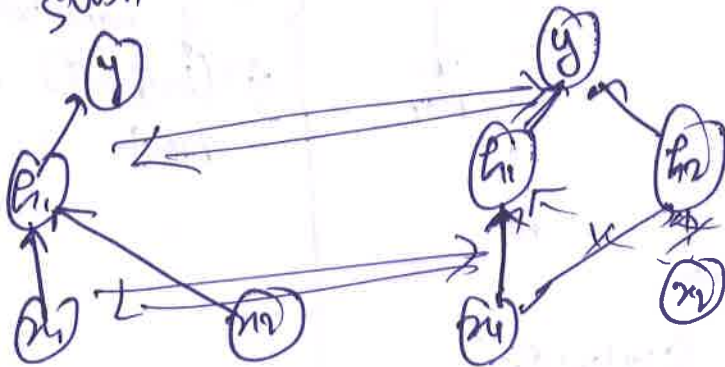
Idea: Use in conjunction with mini batch algorithm

$(x^{(i)}, y^{(i)}) \sim x^{(i)}$
size minibatch (independent examples)

$(x^{(m)}, y^{(m)})$
Batch is independent of each other

- ① construct a sub network using random sampling
 - ② compute gradient, \uparrow [Forward Pass, Backward Pass]
 - ③ moved to next minibatch
- stochastic gradient descent

Difference - (from Bagging) subnetworks



parameters of each model are tied with each other.

Therefore, running single iteration (each minibatch) with a subnetwork is also ok.

②

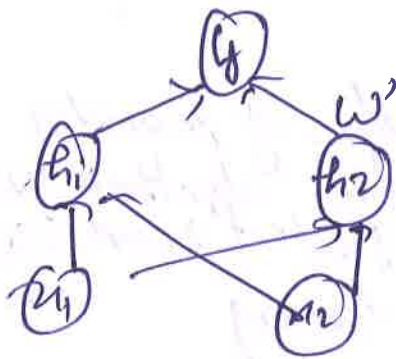
Why does it / should it work?

What is the form of regularization that it is introducing:- it is reducing over dependence on a specific architecture.

~~State~~

Inference:- Need to average over model
(implicitly represented)
↳ A very large # of them

Idea:-



$$w' = p \times w$$

p:- probability of dropping w .

Weight scaling inference rule

$$\tilde{P}_{\text{dropout}}(y|x) \doteq \left(\prod_{u \in U} p(y|x; w_u) \right)^{1/|U|}$$

possible subnetworks

weight scaling inference.

units to be dropped

Why does it work?:-

Robust set of features

Drop inputs w work as hidden units

Provide additional enhancements compared to independently trained ensemble.

Adversarial Training

Read from Book

Optimization:-

(Done with Modeling + Regularization)

$$J(\theta) + \lambda L(\theta)$$

\Downarrow
cost function
(-ve log-likelihood)

\Downarrow regularizer

- ↳ ① way different from pure optimization
- ② challenges
- ③ specific Algorithms

- ↳ ① Initialization of parameters
- ② Optimization algorithm (various)
↳ gradient based / second order
- ③ Adaptive learning rate

$$J(\theta) = E_{(x,y) \sim \text{Data}} [L(f(x;\theta); y)]$$

\Downarrow per-example loss

$\frac{1}{n} \sum_{i=1}^n$

proxy for actual performance metric that you care about (generalization error)

$$[J^*(\theta) = E_{(x,y) \sim \text{Data}} [L(f(x;\theta); y)]]$$

\Downarrow

if we know this, then pure optimization problem

③

$$\Rightarrow \underset{0}{\text{argmin}} \rightarrow \text{Training Loss} = \frac{1}{m} \sum_{i=1}^m L(f(x_i; \theta), y_i)$$

Empirical
Risk
Minimization

How is it related
to true risk
(generalization
error)

PAC Learning

↳ Problem of Overfitting

Issues:- May care about \mathcal{L}_0 loss
↳ Not differentiable

Optimize Surrogate Loss :-

Negative log-likelihood

↳ Another idea:-

Early Stopping

[convergence criteria
↳ Based on a validation set]

Batch / Mini Batch Algorithms :-

Batch gradient Descent

Minibatch / Stochastic
gradient Descent

$$\frac{1}{m} \sum_{i=1}^m L(f(x_i; \theta), y_i)$$

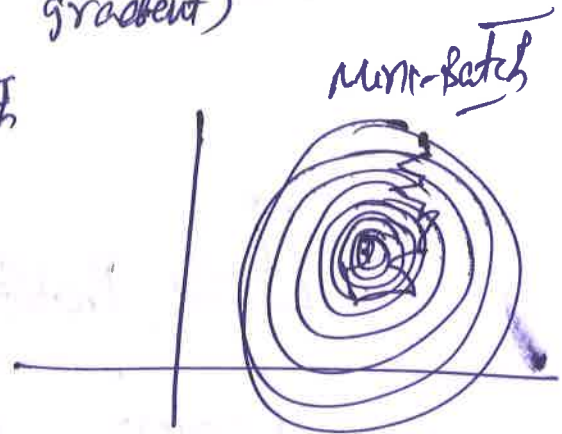
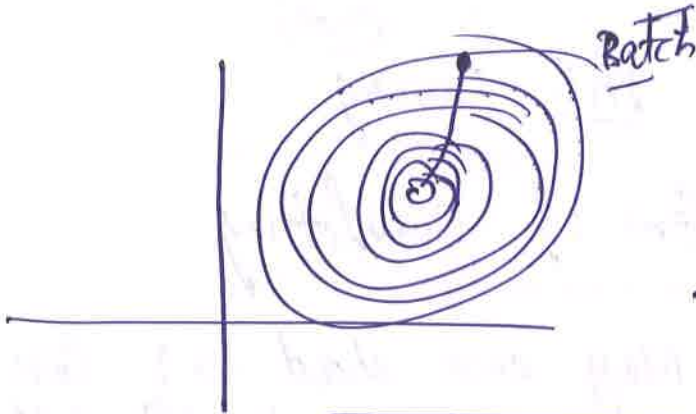
⇔

$$\frac{1}{B} \sum_{i=1}^B L(f(x_i; \theta), y_i)$$

① Divide into m/B batches
② could be randomly chosen

Significant computational advantages

- ↳ ① Each iteration $O(m)$ time
(Batch gradient Descent)
- ② Each iteration takes $O(r)$ time
(approximate gradient)



↳ Offers regularization

Not practical for most applications

Most Algorithms
 $1 \leq r \leq m$
Mini Batch

Gradient: $m \rightarrow$ estimate $E[x]$
 $\hat{\mu} = \frac{1}{m} \sum_{i=1}^m x^{(i)}$ i.i.d.

$$\hat{\mu} = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$E[\hat{\mu}] = \mu$$

$$\sigma[\hat{\mu}] =$$

$$\text{var}[\hat{\mu}] = \frac{1}{m} \sum_{i=1}^m \text{var}(x^{(i)})$$

$$= \frac{1}{m^2} \sigma^2 m = \frac{\sigma^2}{m}$$

\Rightarrow Standard Deviation in estimate of μ $= \frac{\sigma}{\sqrt{m}} \rightarrow$ goes down as m increases
 But m times more computation stochastic

④

Considerations:-

- ① Larger Batch size \Rightarrow more accurate estimate of gradient
- ② Parallelization within batch?
 - \hookrightarrow prefer larger batches
 - \hookrightarrow certain batch sizes may be preferable (powers of two)
- ③ Regularizing effect: small batch size
- ④ Kind of algorithm:-

① Gradient Based \checkmark

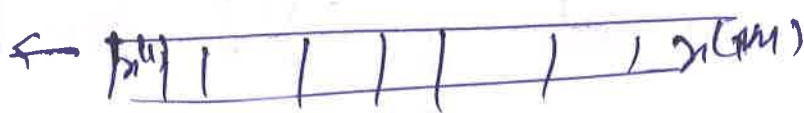
② Second order $H^{-1}g$
 \hookrightarrow small errors in estimating g can magnify \Rightarrow larger batch

⑤ Example in each batch should be uncorrelated (independent):-

- ① Similarly examples across batches be independent from each other.
- ② When examples are not repeated, stochastic gradient descent is a unbiased estimate of true gradient why?

$$J(\theta) = (x, y) \sim p_{\text{data}}[I \neq (x, y)]$$

Variance may be high



Challenges in Neural Network Optimization :-

Non-Convex Optimization

① ill-conditioning of Hessian Matrix :-

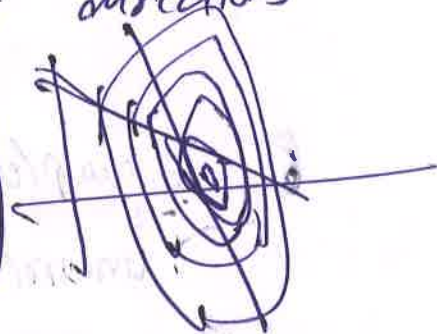
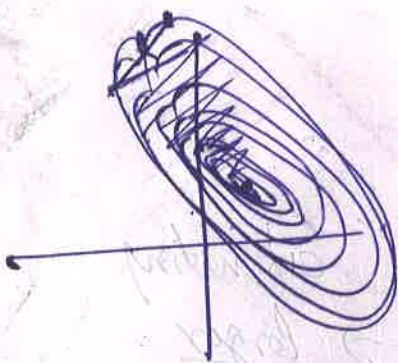
Conditioning Number :-

Matrix :- H (Hessian Matrix)

Condition Number :- $\max_{i,j} \left| \frac{\lambda_i}{\lambda_j} \right|$

ratio of the largest & the smallest eigenvalues

\Rightarrow High Condition number - \Rightarrow varying curvature in different directions



$$J(\theta) = J(\theta^{(0)}) + (\theta - \theta^{(0)})^T \underset{g}{\nabla J(\theta)} + (\theta - \theta^{(0)})^T H (\theta - \theta^{(0)})$$

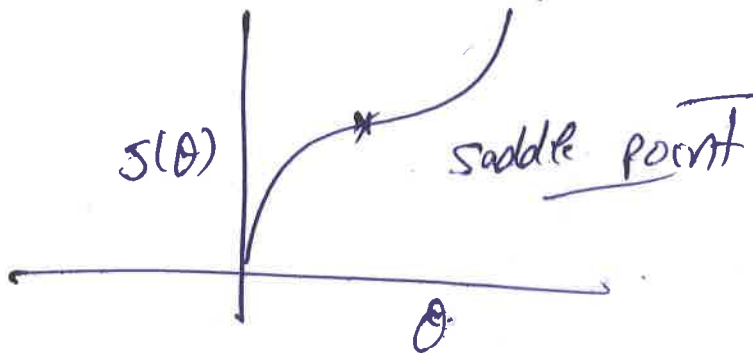
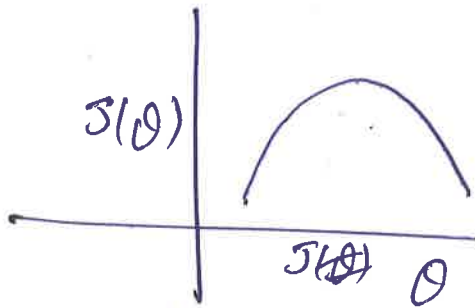
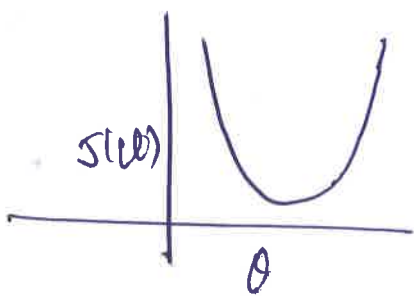
Learning rate :- η

$$\theta \leftarrow \theta^{(0)} - \epsilon g$$

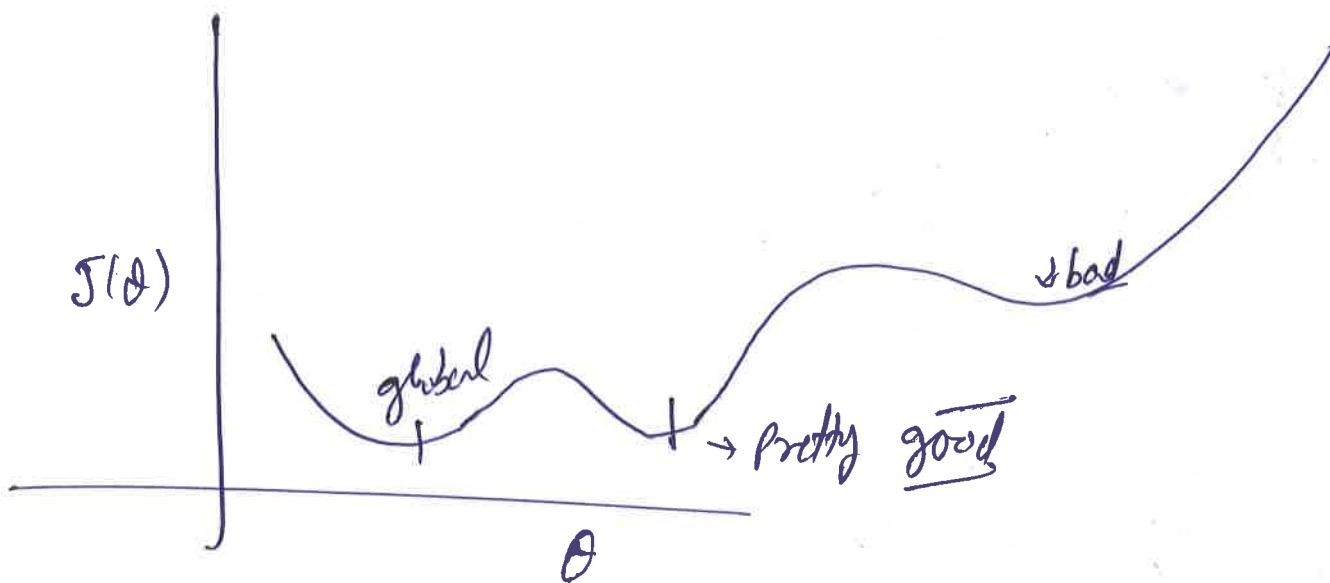
$$J(\theta^{(0)} - \epsilon g) = J(\theta^{(0)}) - \epsilon g^T g + \frac{1}{2} \epsilon^2 g^T H g$$

$$\text{Decrease in cost :- } \left[\epsilon g^T g - \frac{1}{2} \epsilon^2 g^T H g \right]$$

5)



H has both
+ve & -ve
eigenvalues



Dropout proof:-

Ensemble prediction:-

Can work
better than
sampling

Average prediction $\left[\sum_n \frac{1}{n} P(y|x; \theta_n) \right] \frac{1}{2d}$

multiply
each weight
by $\frac{1}{2d}$
Scaling
No non-linearity

Geometric
Average

$\left[\prod_n P(y|x; \theta_n) \right]^{1/2d} \rightarrow \# \text{ of configurations}$
Klinton et al., 2012

1. 1. 1. 1. 1.
 2. 2. 2. 2. 2.
 3. 3. 3. 3. 3.



1. 1. 1. 1. 1.
 2. 2. 2. 2. 2.
 3. 3. 3. 3. 3.

1. 1. 1. 1. 1.
 2. 2. 2. 2. 2.
 3. 3. 3. 3. 3.

1. 1. 1. 1. 1.
 2. 2. 2. 2. 2.
 3. 3. 3. 3. 3.

1. 1. 1. 1. 1.
 2. 2. 2. 2. 2.
 3. 3. 3. 3. 3.