

Lecture 7 Recap

Regression Losses: L2 vs L1

• L2 Loss:

$$-L^{2} = \sum_{i=1}^{n} (y_{i} - f(x_{i}))^{2}$$

- Sum of squared differences (SSD)
- Prone to outliers
- Compute-efficient (optimization)
- Optimum is the mean

• L1 Loss:

$$-L^{1} = \sum_{i=1}^{n} |y_{i} - f(x_{i})|$$

- Sum of absolute differences
- Robust
- Costly to compute
- Optimum is the median



Softmax Formulation

• What if we have multiple classes?



Example: Hinge vs Cross-Entropy

Hinge Loss:
$$L_i = \sum_{k \neq y_i} \max(0, s_k - s_{y_i} + 1)$$

Cross Entropy : $L_i = -\log(\frac{e^{s_{y_i}}}{\sum_k e^{s_k}})$

Hinge loss: Cross Entropy loss: Given the following scores for x_i : $\max(0, -3 - 5 + 1) +$ $-\ln\left(\frac{e^5}{e^{5}+e^{-3}+e^2}\right) = 0.05$ s = [5, -3, 2]Model 1 $\max(0, 2 - 5 + 1) = 0$ $\max(0, 10 - 5 + 1) +$ $-\ln\left(\frac{e^5}{e^{5}+e^{10}+e^{10}}\right) = 5.70$ s = [5, 10, 10]Model 2 $\max(0, 10 - 5 + 1) = 12$ s = [5, -20, -20] $-\ln\left(\frac{e^5}{e^{5}+e^{-20}+e^{-20}}\right)$ $\max(0, -20 - 5 + 1) +$ Model 3 $\max(0, -20 - 5 + 1) = 0$ $= 2 * 10^{-11}$ $y_i = 0$ Cross Entropy *always* wants to improve! (loss never 0) Hinge Loss saturates. I2DL: Prof. Dai 5

Sigmoid Activation



TanH Activation



[LeCun et al. 1991] Improving Generalization Performance in Character Recognition

Rectified Linear Units (ReLU)



Quick Guide

• Sigmoid is not really used.

• ReLU is the standard choice.

• Second choice are the variants of ReLU or Maxout.

• Recurrent nets will require TanH or similar.



Xavier/Kaiming Initialization

• How to ensure the variance of the output is the same as the input?

$$(nVar(w)Var(x)) = 1$$

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

ReLU Kills Half of the Data





Lecture 8



Data Augmentation

Data Pre-Processing



For images subtract the mean image (AlexNet) or per-channel mean (VGG-Net)

Data Augmentation

• A classifier has to be invariant to a wide variety of transformations



cat



Videos

Shopping More

Settings Tools

J Q

0



SafeSearch *



Cute



And Kittens









White Cats And Kittens





News















Illumination

I2DL: Prof. Dai



Appearance

17

Data Augmentation

• A classifier has to be invariant to a wide variety of transformations

• Helping the classifier: synthesize data simulating plausible transformations

Data Augmentation

a. No augmentation (= 1 image)



224x224



b. Flip augmentation (= 2 images)



224x224



c. Crop+Flip augmentation (= 10 images)



224x224



+ flips

Data Augmentation: Brightness

• Random brightness and contrast changes



Data Augmentation: Random Crops

- Training: random crops
 - Pick a random L in [256,480]
 - Resize training image, short side L
 - Randomly sample crops of 224x224



- Testing: fixed set of crops
 Resize image at N scales
 - 10 fixed crops of 224x224: (4 corners + 1 center) × 2 flips

21

Data Augmentation: Advanced

Magnitude: 9

ShearX

Magnitude: 17







AutoContrast





Original





Original



ShearX



AutoContrast

AutoContrast









Sample strength



Sample augmentation and apply it

Muller et al., Trivial Augment, ICCV 2021



I2DI : Prof. Dai

Data Augmentation

• When comparing two networks make sure to use the same data augmentation!

Consider data augmentation a part of your network
 design



Advanced Regularization

L2 regularization, also (wrongly) called weight decay

• L2 regularization



- Penalizes large weights
- Improves generalization



L2 regularization, also (wrongly) called weight decay

• Weight decay regularization



• Equivalent to L2 regularization in GD, but not in Adam.

Loshchilov and Hutter, Decoupled Weight Decay Regularization, ICLR 2019

Early Stopping



Bagging and Ensemble Methods

• Train multiple models and average their results

• E.g., use a different algorithm for optimization or change the objective function / loss function.

• If errors are uncorrelated, the expected combined error will decrease linearly with the ensemble size

Bagging and Ensemble Methods

• Bagging: uses k different datasets (or SGD/init noise)





Dropout

Dropout

• Disable a random set of neurons (typically 50%)



(a) Standard Neural Net



(b) After applying dropout.

• Using half the network = half capacity



(b) After applying dropout.



- Using half the network = half capacity
 - Redundant representations
 - Base your scores on more features

• Consider it as a model ensemble

• Two models in one



(b) After applying dropout.









[Srivastava et al., JMLR'14] Dropout





- Using half the network = half capacity
 - Redundant representations
 - Base your scores on more features
- Consider it as two models in one
 - Training a large ensemble of models, each on different set of data (mini-batch) and with SHARED parameters

Reducing co-adaptation between neurons

Dropout: Test Time

• All neurons are "turned on" – no dropout



Conditions at train and test time are not the same

PyTorch: model.train() and model.eval()

Dropout: Test Time Dropout probability $z = (\theta_1 x_1 + \theta_2 x_2) \cdot p$ p = 0.5• Test: $E[z] = \frac{1}{4} (\theta_1 0 + \theta_2 0) \\ + \theta_1 x_1 + \theta_2 0 \\ + \theta_1 0 + \theta_2 x_2 \\ + \theta_1 x_1 + \theta_2 x_2) \\ = \frac{1}{2} (\theta_1 x_1 + \theta_2 x_2)$ • Train: Z θ_1 θ_2 x_2 x_1 Weight scaling inference rule

Dropout: Before

• Efficient bagging method with parameter sharing

- Try it!
- Dropout reduces the effective capacity of a model, but needs more training time

• Efficient regularization method, can be used with L2

Dropout: Nowadays

- Usually does not work well when combined with batch-norm.
- Training takes a bit longer, usually 1.5x
- But, can be used for uncertainty estimation.
- Monte Carlo dropout (Yarin Gal and Zoubin Ghahramani series of papers).

Monte Carlo Dropout

- Neural networks are massively overconfident.
- We can use dropout to make the softmax probabilities more calibrated.
- Training: use dropout with a low p (0.1 or 0.2).
- Inference, run the same image multiple times (25-100), and average the results.

Gal et al., Bayesian Convolutional Neural Networks with Bernoulli Approximate Variational Inference, ICLRW 2015 Gal and Ghahramani, Dropout as a Bayesian approximation, ICML 2016 Gal et al., Deep Bayesian Active Learning with Image Data, ICML 2017 Gal, Uncertainty in Deep Learning, PhD thesis 2017

ТΠ

Batch Normalization: Reducing Internal Covariate Shift

ТΠ

Batch Normalization: Reducing Internal Covariate Shift

What is internal covariate shift, by the way?

Our Goal

• All we want is that our activations do not die out



I2DL: Prof. Dai

- Wish: Unit Gaussian activations (in our example)
- Solution: let's do it



Mean of your mini-batch examples over feature k $\widehat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{Var[x^{(k)}]}}$

• In each dimension of the features, you have a unit gaussian (in our example)



• In each dimension of the features, you have a unit gaussian (in our example)

• For NN in general, BN normalizes the mean and variance of the inputs to your activation functions

BN Layer

• A layer to be applied after Fully Connected (or Convolutional) layers and before non-linear activation functions



• 1. Normalize

$$\widehat{\boldsymbol{x}}^{(k)} = \frac{\boldsymbol{x}^{(k)} - E[\boldsymbol{x}^{(k)}]}{\sqrt{Var[\boldsymbol{x}^{(k)}]}} \qquad \qquad \text{Differentiable function so we} \\ \text{can backprop through it}...$$

• 2. Allow the network to change the range

$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}$$
 These parameters will be optimized during backprop

• 1. Normalize

$$\widehat{\boldsymbol{x}}^{(k)} = \frac{\boldsymbol{x}^{(k)} - E[\boldsymbol{x}^{(k)}]}{\sqrt{Var[\boldsymbol{x}^{(k)}]}}$$

• 2. Allow the network to change the range

$$y^{(k)} = \gamma^{(k)} \widehat{x}^{(k)} + \beta^{(k)}$$

backprop

The network *can* learn to undo the normalization

$$\gamma^{(k)} = \sqrt{Var[\mathbf{x}^{(k)}]}$$
$$\beta^{(k)} = E[\mathbf{x}^{(k)}]$$

• Ok to treat dimensions separately? Shown empirically that even if features are not correlated, convergence is still faster with this method

BN: Train vs Test

 Train time: mean and variance is taken over the minibatch

$$\widehat{\boldsymbol{x}}^{(k)} = \frac{\boldsymbol{x}^{(k)} - \boldsymbol{E}[\boldsymbol{x}^{(k)}]}{\sqrt{Var[\boldsymbol{x}^{(k)}]}}$$

- Test-time: what happens if we can just process one image at a time?
 - No chance to compute a meaningful mean and variance

BN: Train vs Test

Training: Compute mean and variance from mini-batch 1,2,3 ...

Testing: Compute mean and variance by running an exponentially weighted averaged across training minibatches. Use them as σ_{test}^2 and μ_{test} .

 $\begin{aligned} Var_{running} &= \beta_m * Var_{running} + (1 - \beta_m) * Var_{minibatch} \\ \mu_{running} &= \beta_m * \mu_{running} + (1 - \beta_m) * \mu_{minibatch} \\ \beta_m : \text{momentum (hyperparameter)} \end{aligned}$

I2DL: Prof. Dai

BN: What do you get?

• Very deep nets are much easier to train, more stable gradients

• A much larger range of hyperparameters works similarly when using BN

BN: A Milestone



BN: Drawbacks



Other Normalizations



Other Normalizations





What We Know



Concept of a 'Neuron'



Activation Functions (non-linearities)



Backpropagation



SGD Variations (Momentum, etc.)



Data Augmentation

a. No augmentation (= 1 image)





b. Flip augmentation (= 2 images)







Weight Regularization e.g., L^2 -reg: $R^2(W) = \sum_{i=1}^N w_i^2$ Batch-Norm

$$\widehat{\boldsymbol{x}}^{(k)} = \frac{\boldsymbol{x}^{(k)} - E[\boldsymbol{x}^{(k)}]}{\sqrt{Var[\boldsymbol{x}^{(k)}]}}$$

Dropout



(b) After applying dropout.

Why not simply more layers?

- Neural nets with at least one hidden layer are universal function approximators.
- But generalization is another issue.
- Why not just go deeper and get better?
 - No structure!!
 - It is just brute force!
 - Optimization becomes hard
 - Performance plateaus / drops!
- We need more! More means CNNs, RNNs and eventually Transformers.



See you next week!

References

- Goodfellow et al. "Deep Learning" (2016),
 Chapter 6: Deep Feedforward Networks
- Bishop "Pattern Recognition and Machine Learning" (2006),
 Chapter 5.5: Regularization in Network Nets
- http://cs231n.github.io/neural-networks-1/
- http://cs231n.github.io/neural-networks-2/
- http://cs231n.github.io/neural-networks-3/