CPSC 340: Machine Learning and Data Mining

Discrete Labels Fall 2015

Admin

- Assignment 5 is posted.
 - Due Friday of next week.
 - A2.2 update: use k = 10.

Last Time: Convolutional Neural Networks

- Convolutional neural networks:
 - 1. Convolutional layers.
 - 2. Pooling layers.
 - 3. Fully-connected layers.



http://blog.csdn.net/strint/article/details/44163869

http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf

GoogLeNet

- GoogLeNet is very deep competitive-winning system.
- Training these systems is very expensive:
 - Weeks on clusters of computers with expensive GPUs.
- What if you aren't Google?
 - Some groups release their features/network code.
 - Quickly learn to identify new objects by:
 - Use features learned from millions of images and thousands classes

MetaMind

Add Files

Label 1 (0)

Multiple Files Single Zip

Label 1

Add Files

Label 2 (0)

🖌 Label 2

UPI OAD I

Add Clas

+

• Or 'fine-tune' the entire network with your data.



	-		<u> </u>	
		Ī		
	Sof	tmaxActiva	tion	
	_	t		
		EC		
		AveragePoo 7x7+1(V)	ol .	
		+		
)anthCanc	. +	
		$\angle \Delta$		
	Conv Con 1x1+1(S) 3x3+1	v (L(S) 5x	Conv 5+1(S)	Conv lxl+l(S)
	<u></u>		1	+
	Con	v (Conv	MaxPool
C	1x1+1	l(S) 1x	1+1(S)	3x3+1(S)
5.		<u>\</u> /		
5	C C	DepthConce	at	
	_		•	
	Conv Con		Conv	Conv
I ∨ I <i>I</i> ≓	1x1+1(S) 3x3+1	L(S) 5x	5+1(S)	1x1+1(S)
			Ť	Ť
in for	Con	V (C)		MaxPool
		1(5) IX	4	3X3+1(5)
		3x3+	2(S)	
		+	1	

Today: Alternatives to Squared Loss

• We've been using squared error as our default 'loss':

$$f(\hat{y}_i) = \sum_{i=1}^{n} (\gamma_i - \hat{\gamma}_i)^2$$

- A lot of models we've discussed fit in this framework:
- Linear regression: PCA: $\hat{\chi}_{i} = w^{T} \chi_{i}$ PCA: Collaborative fitering: Deep neural networks: $\hat{\chi}_{i} = w^{T} \chi_{i}$ $\hat{\chi}_{ij} = w_{j}^{T} Z_{i}$ $\hat{\chi}_{un} = b_{u} + b_{m} + w_{m} Z_{u}$ $\hat{\chi}_{i} = w^{T} h(W_{ij}h(W_{$
 - Square error is differentiable and sometimes has closed-form.
 - But, usually squared error is not the right 'loss'.
 - Today we discuss alternatives.

- We previously discussed a few alternatives:
 - L1-error and Huber loss are more robust to outliers.



- We previously discussed a few alternatives:
 - L1-error and Huber loss are more robust to outliers.
 - Non-convex losses can be even more robust.



- We previously discussed a few alternatives:
 - L1-error and Huber loss are more robust to outliers.
 - Non-convex losses can be even more robust.
 - Maximum loss gives better performance in worst case.

 $\max |y_i - \hat{y}_i|$ Х $\sum_{j=1}^{n} \left(y_i - \frac{\gamma}{y_i} \right)^{\mathsf{x}}$ > use this to make neural network focus on worst case.

- We previously discussed a few alternatives:
 - L1-error and Huber loss are more robust to outliers.
 - Non-convex losses can be even more robust.
 - Maximum loss gives better performance in worst case.
 - Hinge and logistic losses are better for binary data.



- We previously discussed a few alternatives:
 - L1-error and Huber loss are more robust to outliers.
 - Non-convex losses can be even more robust.
 - Maximum loss gives better performance in worst case.
 - Hinge and logistic losses are better for binary data.
- What about other types of discrete labels?
 - Multi-label: {'cat', 'dog', 'human'}.
 - Categorical: {'Edmonton', 'Paris', 'Philadelphia', 'Vancouver'}.
 - Ordinal: {1 star, 2 stars, 3 stars, 4 stars, 5 stars}.
 - Counts: 602 'likes'.
 - Ranking: Difficulty(A3) > Difficulty(A4) > Difficulty (A2) > DifficultyA(1).

Probabilistic Models and Loss Functions

- We can use probabilistic models to derive loss functions.
- Main idea: ${\color{black}\bullet}$
 - $\rho(y_i | \hat{\gamma}_i)$ Define probability of each possible label: 1.
 - 2. Define loss as negative logarithm of the probability. $-|_{og} p(x_i | \hat{y}_i)$
- Why??? •
 - We want predictions that maximize the probability of the label y_i.
 - Taking logarithm doesn't change location of maximum.
 - Maximizing logarithm is equivalent to minimizing negative of logarithm.

$$\underset{\hat{y}_{i}}{\operatorname{argmax}} p(y_{i} | \hat{y}_{i}) \iff \operatorname{argmax}} \log(p(y_{i} | \hat{y}_{i})) \iff \operatorname{argmin}_{\hat{y}_{i}} - \log(p(y_{i} | \hat{y}_{i}))$$

Sigmoid Probabilities and Logistic Loss

• Example of going from probabilities to loss function (binary y_i):

Define $p(y_i = +1|\hat{y}_i) = \frac{1}{1 + exp(-\hat{y}_i)}$ Signoid function Ensures that $p(y_i = +1|\hat{y}_i) + p(y_i = -1|\hat{y}_i) = \frac{1}{1 + exp(\hat{y}_i)}$ Can write both cases as $p(y_i|\hat{y}_i) = \frac{1}{1 + exp(-\hat{y}_i)}$

Sigmoid Probabilities and Logistic Loss

• Example of going from probabilities to loss function (binary y_i):

Define
$$p(y_i = +1|\hat{y}_i) = \frac{1}{1 + exp(-\hat{y}_i)}$$

Sigmoid function
Ensures that $p(y_i = +1|\hat{y}_i) = (0,1)$
Take logarithm of probability:
 $\log(p(y_i|\hat{y}_i)) = \log(\frac{1}{1 + exp(-y_i\hat{y}_i)}) = \log(1) - \log(1 + exp(-y_i\hat{y}_i)))$
Loss is negative logarithm of probability: $-\log(p(y_i|\hat{y}_i)) = \log(1 + exp(-y_i\hat{y}_i))$

Why Logarithm?

• We want loss function to be additive across examples:

$$\sum_{i=1}^{n} g(y_i, \hat{y}_i) \text{ for some function } g'$$

• If training examples are IID, probability is multiplicative:

$$P(Y_{1},Y_{2},Y_{3},\cdots,Y_{n},\hat{y}_{1},\hat{y}_{2},\hat{y}_{3},\cdots,\hat{y}_{n}) = p(y_{1},\hat{y}_{1})p(y_{2},\hat{y}_{2},\hat{y}_{3})p(y_{3},\hat{y}_{3},\cdots,p(y_{n},\hat{y}_{n}))$$

• Logarithm of the probability is additive across examples:

$$\log(p(y_1, y_2, \dots, y_n), \hat{y_1}, \hat{y_2}, \dots, \hat{y_n}) = \log(p(y_1, |\hat{y_1}|)) + \log(p(y_2|\hat{y_2}|)) + \dots + \log(p(y_n|\hat{y_n}))$$

Gaussian Probabilities and Squared Loss

Consider a continuous label y:

If
$$y_i$$
 is a normal/Gaussian distribution with mean \hat{y}_i , then
 $p(y_i | \hat{y}_i) \propto exp(-\frac{(y_i - \hat{y}_i)^2}{2\sigma^2})$
 $\prod_{multiplied} by factors not depending on y_i
Take negative of logarithm:
 $-\log(p(y_i | \hat{y}_i)) = -\log(exp(-\frac{(y_i - \hat{y}_i)^2}{2\sigma^2})) - \frac{constant}{constant}$
 $= \frac{1}{2\sigma^2}(y_i - \hat{y}_i)^2 - constant.$
Squared error when $\sigma^2 = 1$ and we ignore constant$

Multi-Label Data and Independent Logistic Losses

- Consider the case of multi-label data like {'cat', 'dog', 'human'}.
 Image can have none, some, or all of these labels. E.g., y_i = {1,-1,1}.
- We can treat this case in a similar way to latent-factor models:



- Squared error has 'bad' errors: penalized for being too right.
- We could use logistic loss on each category: $\operatorname{Megnin}_{i=1} \stackrel{n}{\leq} \stackrel{z}{\leq} \log(1 + exp(-y_i w_j^T x_i)))$

Categorical Data and One vs. All

- Categorical data: we have multiple labels but only one is correct.
 - Classifying images as taken in {'Alberta', 'Paris', 'Philadelphia', Vancouver'}.
 In this case, we could have y_i = 3.
- We can use same model as previous slide:



- To make a single prediction, take the biggest w_i^Tx_i.
 - 'One vs. all' classifier.



Categorical Data and Softmax Loss

- Disadvantage of 'one vs. all':
 - Logistic loss focuses on $\{0,1\}$ decisions, not making $w_i^T x_i$ large for correct y_{ij} .
 - We want a loss that makes $w_i^T x_i$ big for correct y_{ij} , small for others.
- 'Softmax' or multinomial logistic regression model:

- Exponential magnifies large values: most probability is on max.
- Negative of logarithm of probability gives loss function.
 - Logistic loss is special case where $w_1=0$.

Unbalanced Data and Extreme-Value Loss

- Consider binary case where:
 - One class overwhelms the other class ('unbalanced' data).
 - Really important to find the minority class (e.g., minority class is tumor).



Unbalanced Data and Extreme-Value Loss

• Extreme-value distribution:

С

$$p(y_{i} = +1 | \hat{y}_{i}) = 1 - exp(-exp(\hat{y}_{i})) \quad [+1 \text{ is majority class}] \quad asymmetric asym$$

Unbalanced Data and Extreme-Value Loss

• Extreme-value distribution:

0.8

0.6

0.4

0.2

-0.2

-1

-0.8

$$p(y_{i} = +1 | \hat{y}_{i}) = 1 - exp(-exp(\hat{y}_{i})) \quad [+1 \text{ is majority class}] \qquad \text{asymmetric}$$

$$To make it a probability_{5} \quad p(y_{i} = -1 | \hat{y}_{i}) = exp(-exp(\hat{y}_{i})) \qquad \text{Exterme Value Regression} \quad (error = 0.13)$$

$$Logistic (blue have & bigger weight) \quad (error = 0.15) \qquad \text{Exterme Value Regression} \quad (error = 0.13)$$

Ordinal Data and Proportional Odds

- Ordinal data: categorical data where the order matters:
 Rating hotels as {'1 star', '2 stars', '3 stars', '4 stars', '5 stars'}.
- Could do softmax, but softmax ignores the order.
- 'Proportional odds' or 'ordinal logistic regression':



Count Data and Poisson Loss

- Count data: predict the number of times something happens.
 The number of 'likes' that Facebook will get.
- Softmax/ordinal require finite number of categories.
- We probably don't want separate parameter for '654' and '655'.
- Poisson regression: use probability from Poisson count distribution.

Back to discussion of CNNs...

Interpreting CNNs

Hon does prediction change if we hide part of the image?



Inceptionism

- Start with random noise, move image to amplify class label.
 - (And enforce pairwise statistics on image)





http://googleresearch.blogspot.ca/2015/06/inceptionism-going-deeper-into-neural.html

Inceptionism

• Start with an image, amplify features deep in network.







Leaves











Birds & Insects

Starting from random hoise:



CNNs for Artistic Style



http://arxiv.org/pdf/1508.06576v2.pdf

Is this magic?

- For speech recognition and object detection:
 - No other methods have ever given the current level of performance.
 - But, we also don't know how to scale up other universal approximators.
- No baseline methods in many deep learning papers/articles/blogs.
 Would simpler methods obtain similar performance?
- Despite the 'high-level' abstraction of deep models, easily fooled:
 - But progress on fixing 'blind spots'.



CNNs for Rating Selfies

Our training data

Bad selfies



Good selfies



https://karpathy.github.io/2015/10/25/selfie

CNNs for Rating Selfies

- Be female - Face in 1/3 of image - Cut off forehead
- Show long hair
- Oversaturate face
- Use Filter

Do: N

- Add border



Don't: - use low lighting - make head too big - take group shots 1

CNNs for Rating Selfies

"Best" crop:



score 44.5





score 62.8



score 53.1



score 52.0

score 67.3

score 56.3

Summary

- Squared loss is rarely the right measure.
- -log(probability) lets us to define loss from any probability.
- Softmax is natural loss for categorical data.
- Ordinal logistic is natural loss for ordinal data.
- Exotic losses like extreme-value or Poisson for certain situations.

• Next time: finding all the cat videos on YouTube.