Machine Learning Basics: Stochastic Gradient Descent

Sargur N. Srihari srihari@cedar.buffalo.edu

- 1. Learning Algorithms
- 2. Capacity, Overfitting and Underfitting
- 3. Hyperparameters and Validation Sets
- 4. Estimators, Bias and Variance
- 5. Maximum Likelihood Estimation
- 6. Bayesian Statistics
- 7. Supervised Learning Algorithms
- 8. Unsupervised Learning Algorithms
- 9. Stochastic Gradient Descent
- 10. Building a Machine Learning Algorithm
- 11. Challenges Motivating Deep Learning

Stochastic Gradient Descent (SGD)

- Nearly all deep learning is powered by SGD
 SGD extends the gradient descent algorithm
- Recall gradient descent:
 - Suppose y=f(x) where both x and y are real nos
 - Derivative is a function denoted as f'(x) or dy/dx
 - It gives the slope of f(x) at the point x
 - i.e., it specifies how to make a small change in the input to make a corresponding change in the output

 $- f(x+\varepsilon) \approx f(x) + \varepsilon f'(x)$

How Gradient Descent uses derivatives

• Criterion f(x) minimized by moving from current solution in direction of the negative of gradient



- Gradient descent proposes a new point

$$oldsymbol{x}' = oldsymbol{x}$$
 - $arepsilon
abla_{oldsymbol{x}} fig(oldsymbol{x}ig)$

– where ϵ is the learning rate, a positive scalar.

• Set to a small constant

Learning rate for deep learning

- Useful to reduce ϵ as training progresses
- Constant learning rate is default in Keras
 - Momentum and decay are set to $0\ {\rm by}\ {\rm default}$
 - keras.optimizers.SGD(lr=0.1, momentum=0.0, decay=0.0, nesterov=False)





Computational bottleneck

- A recurring problem in machine learning:
 - large training sets are necessary for good generalization
 - but large training sets are also computationally expensive
- SGD is an extension of gradient descent that offers a solution
 - Moreover it is a method of generalization beyond the training set

Cost function is sum over samples

- Criterion in machine learning is a cost function
- Cost function often decomposes as a sum of per sample loss function
 - E.g., Negative conditional log-likelihood of training data is

$$J(\boldsymbol{\theta}) = E_{\boldsymbol{x}, \boldsymbol{y} \sim \hat{p}_{\text{data}}} L(\boldsymbol{x}, \boldsymbol{y}, \boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^{m} L\left(\boldsymbol{x}^{(i)}, \boldsymbol{y}^{(i)}, \boldsymbol{\theta}\right)$$

where m is the no. of samples and

L is the per-example loss $L(x, y, \theta) = -\log p(y|x; \theta)$

In linear regression
we minimize
$$J(\theta) = \frac{1}{2} \sum_{i=1}^{m} \left\{ y^{(i)} - \theta^{T} x^{(i)} \right\}^{2}$$

Gradient is also sum over samples

• For these additive cost functions, gradient descent requires computing

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^{m} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{x}^{(i)}, y^{(i)}, \boldsymbol{\theta})$$

In linear regression

$$\nabla \ln p(y \mid X, \theta, \beta) = \beta \sum_{i=1}^{m} \left\{ y^{(i)} - \theta^{T} \boldsymbol{x}^{(i)} \right\} \boldsymbol{x}^{(i)T}$$

- Computational cost of this operation is O(m)
- As training set size grows to billions, time taken for single gradient step becomes prohibitively long

Insight of SGD

- Insight: Gradient is an expectation $\nabla_{\theta} J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \nabla_{\theta} L(\boldsymbol{x}^{(i)}, \boldsymbol{y}^{(i)}, \boldsymbol{\theta})$
 - Expectation may be approximated using small set of samples
- In each step of SGD we can sample a minibatch of examples $B = \{x^{(1)}, \dots, x^{(m')}\}$
 - drawn uniformly from the training set
 - Minibatch size m' is typically chosen to be small: 1 to a hundred
 - Crucially m' is held fixed even if sample set is in billions
 - We may fit a training set with billions of examples using updates computed on only a hundred examples

10

SGD Estimate on minibatch

• Estimate of gradient is formed as $g = \frac{1}{m'} \nabla_{\theta} \sum_{i=1}^{m'} L(\mathbf{x}^{(i)}, y^{(i)}, \theta)$

– using examples from minibatch *B*

 SGD then follows the estimated gradient downhill

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \varepsilon \boldsymbol{g}$$

– where ε is the learning rate

How good is SGD?

- In the past gradient descent was regarded as slow and unreliable
- Application of gradient descent to non-convex optimization problems was regarded as unprincipled
- SGD is not guaranteed to arrive at even a local minumum in reasonable time
- But it often finds a very low value of the cost function quickly enough

SGD and Training Set Size

- Outside of deep learning
 - SGD is the main way to train large linear models on very large data sets
- Without SGD cost per update increases with m
 - Cost per SGD update does not depend on the training set size m (it depends only on m)
 - As $m \rightarrow \infty$ model will eventually converge to its best possible test error before SGD has sampled every example in the training set
 - Asymptotic cost of training a model with SGD is O(1) as a function of m13

Deep Learning vs SVM

 Prior to advent of DL main way to learn nonlinear models was to use the kernel trick in combination with a linear model

- SVM:

$$f(\boldsymbol{x}) = \boldsymbol{w}^T \boldsymbol{x} + b = b + \sum_{i=1}^m \boldsymbol{\alpha}_i \boldsymbol{x}^T \boldsymbol{x}^{(i)}$$

- Replace *x* by a feature function \$\phi(x)\$ and the dot product with a kernel function \$k(x,x^{(i)})=\phi(x)•\phi(x^{(i)})\$
- Requires constructing an $m \times m$ matrix $G_{i,j} = k(x^{(i)}, x^{(j)})$
 - Constructing this matrix is $O(m^2)$

Growth of interest in Deep Learning

- In academia (with medium sized data sets)
 - Starting in 2006, deep learning was interesting because it performed better on data sets with thousands of examples
- In industry (with large data sets)
 - Because it provided a scalable way of training nonlinear models on large datasets