

F TEC 5050

Fr 1:30 - 4:20

150, E1



NOTEBOOK

CATCH THE STAR THAT HOLDS YOUR DESTINY. THE ONE THAT FOREVER
TWINKLES WITHIN YOUR HEART. TAKE ADVANTAGE OF PRECIOUS
OPPORTUNITIES WHILE THEY STILL SPARKLE BEFORE YOU. ALWAYS
BELIEVE THAT YOUR ULTIMATE GOAL IS ATTAINABLE AS LONG AS YOU
COMMIT YOURSELF TO IT.



Machine Learning and

Artificial Intelligence



Define Learning Problem

- Task, T
Predicting whether or not this startup will receive the series A round of investment.
- Performance measure, P
Percent of "not bad+" responses received in the interview process.
- Experience, E
Examples of question, answer, feedback triplets

Task: Given $X \in \mathcal{X}$, predict $Y \in \mathcal{Y}$.

≡ Construct prediction rule $f : \mathcal{X} \rightarrow \mathcal{Y}$

Performance: Risk $R(f) \equiv \mathbb{E}_{XY} [\text{loss}(Y, f(X))]$

$$(X, Y) \sim P_{XY}$$

Experience: Training Data $\{(X_i, Y_i)\}_{i=1}^n \stackrel{i.i.d.}{\sim} P_{XY}$ (unknown)

Supervised vs Unsupervised

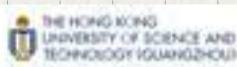
Supervised Learning – Learning with a teacher



Unsupervised Learning – Learning without a teacher



Linear Regression



A simplified model:

$$y = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

Home Owner Marital Status Annual Income

A vector version:

$$y = \mathbf{w}^T \mathbf{x} + b$$

- y is the prediction
- \mathbf{w}^T is the weight
- b is the bias
- \mathbf{w}^T and b together are the parameters

How to quantitatively handle nominal data?

A common solution: one-hot encoding

(1) Home Owner: Yes → [0, 1] No → [1, 0]

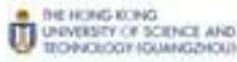
(2) Marital Status: Married → [1, 0, 0], Single → [0, 1, 0], Divorced → [0, 0, 1]

Modern typology:

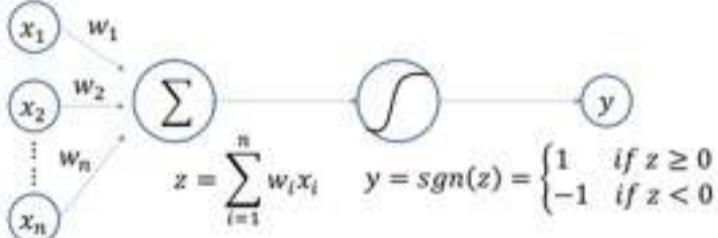
<https://dl.acm.org/doi/pdf/10.1145/2533988>

61

Linear Regression



Can be treated as a one-layer perceptron classifier:



Deep network: $y = f_0(f_1(\dots f_n(x)))$

Unlike linear regression model, $f_i(x)$ is non-linear function

62

Parameter Estimation and Loss Function



$$\text{Loss: } L = \frac{1}{N} \sum_n e_n$$

where e_n represents the prediction error of the n -th instance

Loss: how good a set of values is.

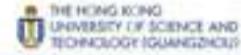
Loss is a function of parameters

$$e = |y - \hat{y}| \quad L \text{ is mean absolute error (MAE)}$$

$$e = (y - \hat{y})^2 \quad L \text{ is mean square error (MSE)}$$

If y and \hat{y} are both probability distributions → Cross-entropy

Optimization



Gradient Descent

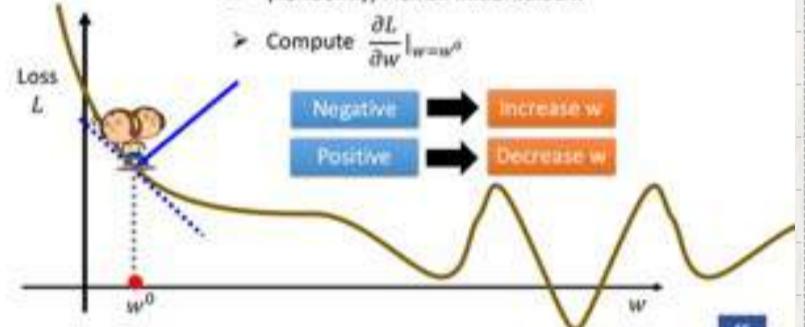
$w^* = \arg \min_w L$

> (Randomly) Pick an initial value w^0

> Compute $\frac{\partial L}{\partial w} |_{w=w^0}$

Negative → Increase w

Positive → Decrease w



Optimization



Gradient Descent

$$w^* = \arg \min_w L$$

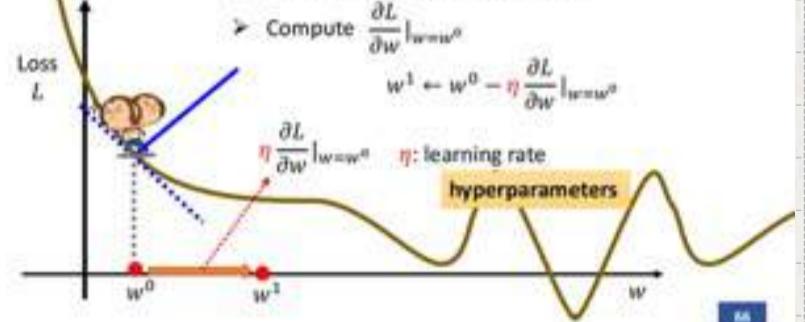
> (Randomly) Pick an initial value w^0

> Compute $\frac{\partial L}{\partial w} |_{w=w^0}$

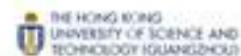
$$w^1 \leftarrow w^0 - \eta \frac{\partial L}{\partial w} |_{w=w^0}$$

η : learning rate

hyperparameters



Optimization



Gradient Descent

$$w^* = \arg \min_w L$$

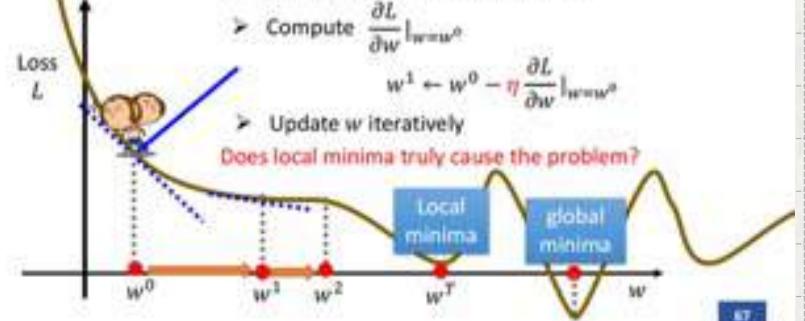
> (Randomly) Pick an initial value w^0

> Compute $\frac{\partial L}{\partial w} |_{w=w^0}$

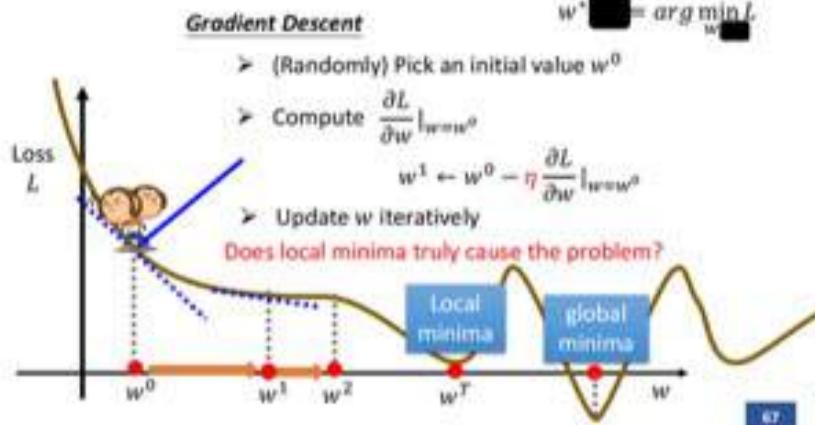
$$w^1 \leftarrow w^0 - \eta \frac{\partial L}{\partial w} |_{w=w^0}$$

> Update w iteratively

Does local minima truly cause the problem?



Optimization



Optimization

- > (Randomly) Pick initial values w^0, b^0
- > Compute

$$\frac{\partial L}{\partial w}|_{w=w^0, b=b^0}$$

$$w^1 \leftarrow w^0 - \eta \frac{\partial L}{\partial w}|_{w=w^0, b=b^0}$$

$$\frac{\partial L}{\partial b}|_{w=w^0, b=b^0}$$

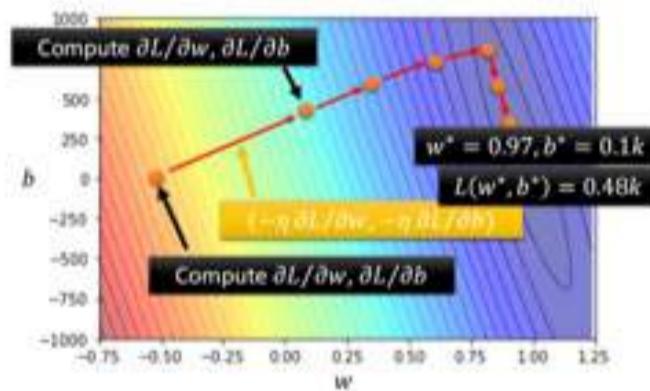
$$b^1 \leftarrow b^0 - \eta \frac{\partial L}{\partial b}|_{w=w^0, b=b^0}$$

Can be done in one line in most deep learning frameworks:

- > Update w and b iteratively

68

Optimization



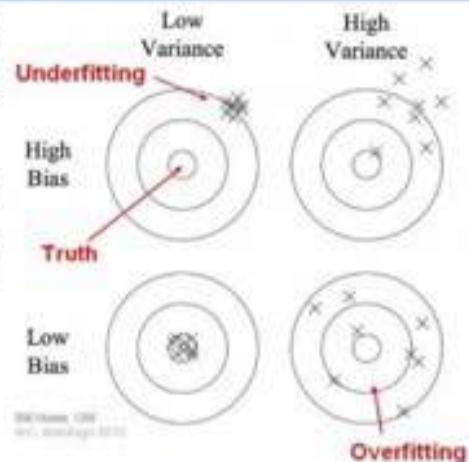
69

Model Assessment

- Generalization of a ML model refers to how well the rules/patterns/functions learned by the ML model, apply to examples not seen by the model.
- The generalization performance of a machine learning method relates to its prediction capability on independent test sets.
- Assessment of this performance is extremely important in practice, since it guides the choice of the machine learning method or model.
- Further, this gives us a measure of the quality of the ultimately chosen model.

Bias and Variance

- Bias: the difference between the predicted value and the true value.
- High bias → underfitting



or model and the training and test sets. One data point which hasn't seen before rates on test data

- Variance is tells us spread.
- High variance → overfitting

See more: <http://www.csie.ntu.edu.tw/~cjlin/paper/bias-variance.html>

Bias and Variance

If we assume that $Y = f(X) + \epsilon$ and $E[\epsilon] = 0$ and $\text{Var}(\epsilon) = \sigma^2$ then we can derive the expression for the expected prediction error of a regression fit $\hat{f}(X)$ at an input $X = x_0$ using squared error loss

$$\text{Err}(x_0) = E[(Y - \hat{f}(x_0))^2 | X = x_0]$$

For notational simplicity let $\hat{f}(x_0) = \hat{f}$, $f(x_0) = f$ and recall that $E[f] = f$ and $E[Y] = f$

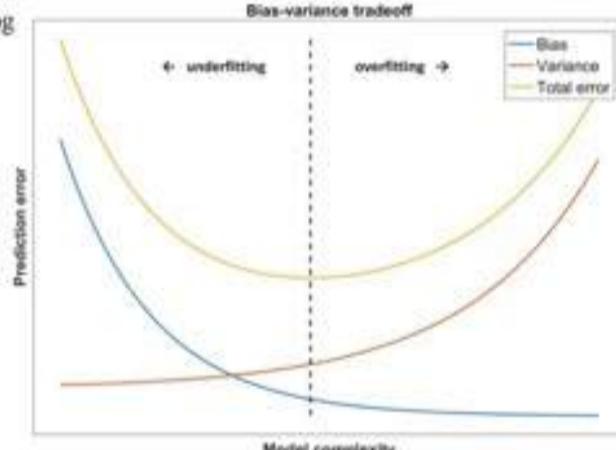
$$\begin{aligned} E[(Y - \hat{f})^2] &= E[(y - f)^2 + E[(f - \hat{f})^2] + 2E[(f - \hat{f})(y - f)] \\ &= E[(f + \epsilon - \hat{f})^2] + E[(f - \hat{f})^2] + 2E[(f - \hat{f})(y - f)] \\ &= E[\epsilon^2] + E[(f - \hat{f})^2] + 2(f^2 - \hat{f}^2 - fE[\hat{f}] + fE[\hat{f}]) \\ &= \sigma^2 + E[(f - \hat{f})^2] \end{aligned}$$

For the term $E[(f - \hat{f})^2]$ we can use a similar trick as above, adding and subtracting $E[\hat{f}]$ to get

$$\begin{aligned} E[(f - \hat{f})^2] &= E[(f + E[\hat{f}] - E[\hat{f}] - \hat{f})^2] \\ &= E[f - E[\hat{f}]]^2 + E[\hat{f} - E[\hat{f}]]^2 \\ &= |f - E[\hat{f}]|^2 + E[\hat{f} - E[\hat{f}]]^2 \\ &= \text{Bias}^2[\hat{f}] + \text{Var}[\hat{f}] \end{aligned}$$

Bias and Variance

Putting it log



73



Title

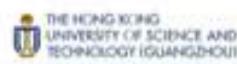
Chapter

Date

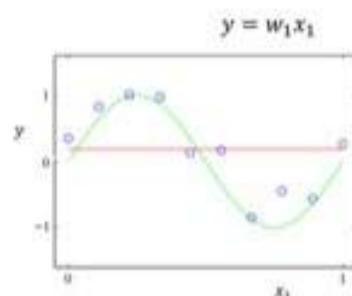
Review

Rating

Model Assessment

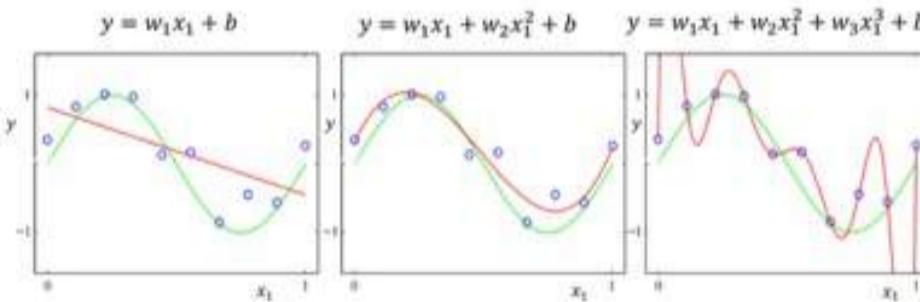


ID	Home Owner	Marital Status	Annual Income	Defaulted Borrower
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes



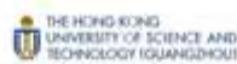
74

Model Assessment



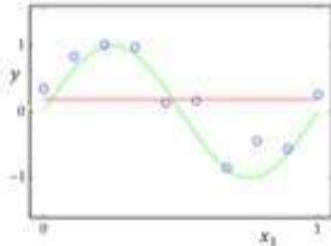
75

Model Assessment



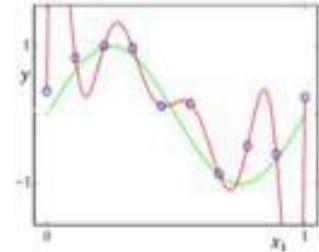
Underfitting: model is too simple — does not fit the data.

$$y = w_1 x_1$$



Overfitting: model is too complex — fits perfectly, does not generalize.

$$y = w_1 x_1 + w_2 x_1^2 + w_3 x_1^3 + b$$



76

Regularization



Regularization is a technique used to reduce errors by fitting the function appropriately on the given training set and avoiding overfitting.

Original Loss

$$\text{New Loss: } L = \frac{1}{N} \sum_n e_n + \lambda J(f)$$

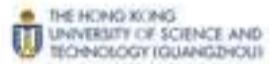
where λ is a positive hyperparameter, $J(f)$ represents the model's complexity. The term $\lambda J(f)$ is called the regularizer or penalty term.

The commonly used regularization techniques are :

1. Lasso Regularization – L1 Regularization
2. Ridge Regularization – L2 Regularization

77

Regularization



L1 Regularization

- Adds the “absolute value of magnitude” of the coefficient as a penalty term to the loss function.
- Helps achieve feature selection by penalizing the weights to approximately equal to zero if that feature does not serve any purpose in the model.

$$L = \frac{1}{N} \sum_n e_n + \lambda \sum_{i=1}^m |w_i|$$

n = number of instances
m = number of features

L2 regularization

- Adds the “squared magnitude” of the coefficient as a penalty term to the loss function.

$$L = \frac{1}{N} \sum_n e_n + \lambda \sum_{i=1}^m w_i^2$$

78

Summary – Linear Regression



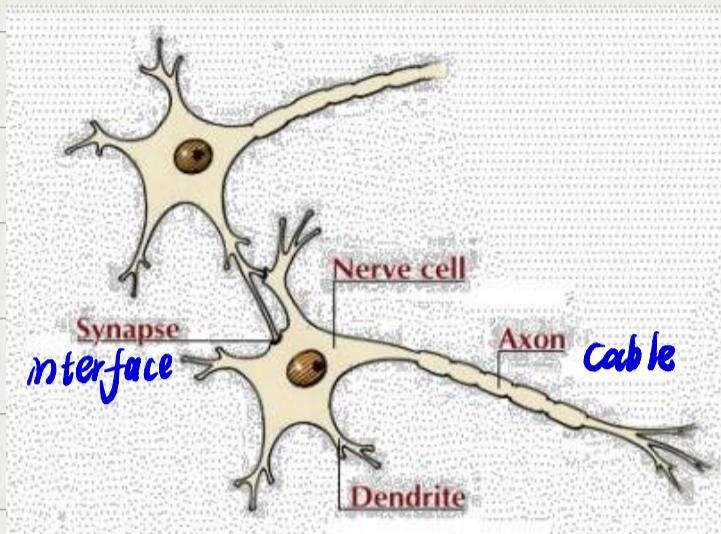
Linear regression exemplifies recurring themes of supervised learning:

- formulate an optimization problem
- choose a model and define a loss function
- improve the generalization by adding a regularizer
- direct solution (set derivatives to zero) gradient descent
- update parameters to solve the optimization problem

79

Deep Feedforward Networks

1. Perceptron 感知器



“分类”

Initialize $w=0$ and $b=0$.

repeat

if $y_i [\langle w, x_i \rangle + b] \leq 0$ then

$$w \leftarrow w + y_i x_i \text{ and } b \leftarrow b + y_i$$

end if

until all classified correctly

$$x = [x_1, \dots, x_n]$$

$$\hat{d}_i = \begin{cases} 1 & \text{if } x_i \in \hat{S} \\ -1 & \text{if } x_i \notin \hat{S} \end{cases}$$

$$\hat{f}(x) = w_0 + w_1 x_1 + w_2 x_2 = w^T x = 0$$

$$f(x) = w^T x_i = \begin{cases} 1 > 0 & \text{if } x_i \in S \\ 1 \leq 0 & \text{if } x_i \notin S \end{cases}$$

Nothing happened if classified correctly

Weight Vector is linear combination

Classifier is linear combination of inner products

$$f(x) = \sum_{i=1}^n y_i \langle x_i, x \rangle + b$$

Convergence Theorem

→ 该定理成立

If there exists some (w^*, b^*) with unit length and $y_i [\langle x_i, w^* \rangle + b^*] \geq \rho$ for all i ,

then the perceptron converges to a linear separator after a number of steps

bounded by $(b^{*2}+1)(r^2+1)\rho^{-2}$ where $\|x_i\| \leq r$

- ① Dimensionality independent
- ② Order independent (i.e., also worst case)
- ③ Scales with ‘difficulty’ of problem

Step 2: Cauchy-Schwartz for the Dot Product

$$\langle (w_{j+1}, b_{j+1}) \cdot (w^*, b^*) \rangle \leq \| (w_{j+1}, b_{j+1}) \| \| (w^*, b^*) \|$$

$$\leq \sqrt{1 + (b^*)^2} \| (w_{j+1}, b_{j+1}) \|$$

Step 3: Upper Bound on $\| (w_j, b_j) \|$

$$\begin{aligned} \| (w_{j+1}, b_{j+1}) \|^2 &= \| (w_j, b_j) + y_i (x_i, 1) \| ^2 \\ &= \| (w_j, b_j) \|^2 + 2y_i \langle (x_i, 1), (w_j, b_j) \rangle + y_i^2 \| (x_i, 1) \|^2 \\ &\leq \| (w_j, b_j) \|^2 + \| (x_i, 1) \|^2 \leq j(r^2 + 1) \end{aligned}$$

Step 4: Combination of first three steps

$$j\rho \leq \sqrt{1 + (b^*)^2} \| (w_{j+1}, b_{j+1}) \| \leq \sqrt{j(r^2 + 1)(1 + (b^*)^2)}$$

Solving for j proves the theorem.

Proof: Starting Point: $w_0=0, b_0=0$

Step 1: Bound on the increase of alignment

Denote by w_i the value of w at step i .

Analogously b_i .

Alignment: $\langle (w_i, b_i), (w^*, b^*) \rangle$

For error in observation (x_i, y_i) we get

$$\langle (w_{j+1}, b_{j+1}) * (w^*, b^*) \rangle$$

$$= \langle [(w_j, b_j) + y_i (x_i, 1)], (w^*, b^*) \rangle$$

$$= \langle (w_j, b_j), (w^*, b^*) \rangle + y_i \langle (x_i, 1) * (w^*, b^*) \rangle$$

$$\geq \langle (w_j, b_j), (w^*, b^*) \rangle + \rho \geq j\rho$$

Alignment increases with number of errors.

随机梯度

Stochastic Gradient With Hinge Loss

① Only need to store errors. This gives a compression bound for perceptron.

② This is stochastic gradient descent on hinge loss

$$l(x_i, y_i; w, b) = \max(0, 1 - y_i(\langle w, x_i \rangle + b))$$

③ Fails miserably with noisy data

2. Loss Function and Objectives

Objectives



- (Linear) Function

$$f(x) = \langle w, x \rangle + b$$

We need to define what to do with it

- Regression

Real valued y , goodness measured by $y - f(x)$

$$l(y, f(x)) = |y - f(x)| \quad \text{absolute value}$$

$$l(y, f(x)) = \frac{1}{2}(y - f(x))^2 \quad \text{least mean squares}$$

$$l(y, f(x)) = \max(0, |y - f(x)| - \epsilon) \quad \epsilon\text{-insensitive}$$

$$l(y, f(x)) = \begin{cases} |y - f(x)| - 0.5 & \text{if } |y - f(x)| > 1 \\ 0.5(y - f(x))^2 & \text{otherwise} \end{cases} \quad \text{Huber}$$

27

Objectives



- Classification

- Binary y , e.g. {apples, oranges}
- Multiple categories, e.g. {red, green, blue}
- Ordinal relationship, e.g. {A, B, C, D, Fail}

$$l(y, f(x)) = \log(1 + e^{-yf(x)}) \quad \text{logistic}$$

$$l(y, f(x)) = \max(0, 1 - yf(x)) \quad \text{soft-margin}$$

$$l(y, f(x, \cdot)) = \log \sum_{y'} e^{f(x, y')} - f(x, y) \quad \text{multiclass logistic}$$

30

多层次神经网络处理非线性可分问题，采用阶跃函数
多层次神经网络可进行任意分类

3. Stochastic Gradient Descent

随机梯度下降

$$\text{Empirical Risk: } \text{Remp}[f] := \frac{1}{m} \sum_{i=1}^m l(y_i, f(x_i))$$

Stochastic gradient descent (pick random x, y) $w_{t+1} \rightarrow w_t - \eta_t \partial f(x_i) l(y_i, f(x_i)) \partial w f(x_i)$

$$\text{Chain rule: } \partial_x [f_2 \circ f_1](x) = [\partial_x f_2 \circ f_1](x) [\partial_x f_1](x)$$

梯度下降更新公式

Convergence in Expectation

$$E_{\theta} [l(\theta)] - l^* \leq \frac{R^2 + L^2 \sum_{t=0}^{T-1} \eta_t^2}{2 \sum_{t=0}^{T-1} \eta_t} \quad \text{initial loss}$$

where $l(\theta) = E_{(x, y)} [l(y, \langle \phi(x), \theta \rangle)]$ and $l^* = \inf_{\theta \in X} l(\theta)$ and $\bar{\theta} = \frac{\sum_{t=0}^{T-1} \theta_t \eta_t}{\sum_{t=0}^{T-1} \eta_t}$

使 y 和 $\phi(x)$ 能接近

expected loss

parameter average

Proof: Show that parameters converge to minimum

$$\theta^* \in \arg \min_{\theta \in X} l(\theta) \quad \text{and set } r_t := \|\theta^* - \theta_t\|$$

Summing over inequality for t proves claim

$$r_{t+1}^2 = \|\pi_x[\theta_t - \eta_t g_t] - \theta^*\|^2$$

$$\leq \|\theta_t - \eta_t g_t - \theta^*\|^2$$

$$= r_t^2 + \eta_t^2 \|g_t\|^2 - 2\eta_t \langle \theta_t - \theta^*, g_t \rangle \quad \text{by convexity}$$

$$\text{hence } E[r_{t+1}^2 - r_t^2] \leq \eta_t^2 L^2 + 2\eta_t [L^* - E[l(\theta_t)]] \leq \eta_t^2 L^2 + 2\eta_t [L^* - E[l(\bar{\theta})]]$$

This yields randomized algorithm for minimizing objective functions

(try logarithms and pick the best/or average

median trick)



Title

梯度下降法

梯度下降法取参数值 $(w^{(0)}, b^{(0)})$
 应用迭代算法求目标函数局部极值

Chapter

Date

Review

/ / / / / /

Rating

Rates



- Guarantee $E_{\bar{\theta}}[l(\bar{\theta})] - l^* \leq \frac{R^2 + L^2 \sum_{t=0}^{T-1} \eta_t^2}{2 \sum_{t=0}^{T-1} \eta_t}$

- If we know R, L, T pick constant learning rate

$$\eta = \frac{R}{L\sqrt{T}} \text{ and hence } E_{\bar{\theta}}[l(\bar{\theta})] - l^* \leq \frac{R[1+1/T]L}{2\sqrt{T}} < \frac{LR}{\sqrt{T}}$$

- If we don't know T pick $\eta_t = O(t^{-\frac{1}{2}})$

This costs us an additional log term

$$E_{\bar{\theta}}[l(\bar{\theta})] - l^* = O\left(\frac{\log T}{\sqrt{T}}\right)$$

39

Strong Convexity



- $l_i(\theta') \geq l_i(\theta) + \langle \partial_\theta l_i(\theta), \theta' - \theta \rangle + \frac{1}{2} \lambda \|\theta - \theta'\|^2$

- Use this to bound the expected deviation

$$\begin{aligned} r_{t+1}^2 &\leq r_t^2 + \eta_t^2 \|g_t\|^2 - 2\eta_t \langle \theta_t - \theta^*, g_t \rangle \\ &\leq r_t^2 + \eta_t^2 L^2 - 2\eta_t [l_t(\theta_t) - l_t(\theta^*)] - 2\lambda\eta_t r_k^2 \end{aligned}$$

hence $E[r_{t+1}^2] \leq (1 - \lambda h_t) E[r_t^2] - 2\eta_t [E[l(\theta_t)] - l^*]$

- Exponentially decaying averaging

$$\bar{\theta} = \frac{1-\sigma}{1-\sigma^T} \sum_{t=0}^{T-1} \sigma^{T-1-t} \theta_t$$

and plugging this into the discrepancy yields

$$l(\bar{\theta}) - l^* \leq \frac{2L^2}{\lambda T} \log \left[1 + \frac{\lambda RT^{\frac{1}{2}}}{2L} \right] \text{ for } \eta = \frac{2}{\lambda T} \log \left[1 + \frac{\lambda RT^{\frac{1}{2}}}{2L} \right]$$

40

4. Backprop

反向传播算法.

Backpropagation

Chain Rules

Layers & Gradients

$$y_i = W_i x_i$$

$$\partial x_i y_i = W_i$$

从输出向输入对冲.

$$x_{i+1} = \sigma(y_i)$$

$$\partial y_i x_{i+1} = \sigma'(y_i)$$

$$\Rightarrow \partial x_i x_{i+1} = \sigma'(y_i) W_i$$

Backpropagation

$$\partial x_i l(y, y_n) = \partial x_{i+1} l(y, y_n) \partial x_i x_{i+1}$$

$$= \partial x_{i+1} l(y, y_n) \sigma'(y_i) W_i^T$$

$$\partial w_i l(y, y_n) = \partial y_i l(y, y_n) \partial w_i y_i$$

$$= \partial x_{i+1} l(y, y_n) \sigma'(y_i) x_i^T$$

Optimization

Layer Representation

$$y_i = W_i x_i$$

$$x_{i+1} = \sigma(y_i)$$

Gradient descent

$$W_i \leftarrow W_i - \eta \partial_w l(y, y_n)$$

Second order method

(use higher derivatives)

$$g_n = l'(y, y_n) W_n \quad \text{- stochastic gradient descent}$$

$$g_i = g_{i+1} \sigma'(y_i) W_i^T \quad \text{(use only one sample)}$$

$$\partial w_i l(y, y_n) = g_{i+1} \sigma'(y_i) x_i^T \quad \text{- Minibatch (small subset)}$$

5. Layers

- Full Connected

- Forward Mapping

$$y_i = W_i x_i$$

$$x_{i+1} = \sigma(y_i)$$

with subsequent nonlinearity

- Backprop gradients

$$\frac{\partial x_i}{\partial w_j} x_{i+1} = \sigma'(y_i) W_i^T$$

$$\frac{\partial w_i}{\partial x_j} x_{i+1} = \sigma'(y_i) x_i^T$$

General purpose layer

- Convolutional Layers

- Feature Locality

Relevant information only in neighborhood of pixel

$$y_{ij} = \sum_{a=-\Delta}^{\Delta} \sum_{b=-\Delta}^{\Delta} W_{ij,ab} x_{i+a, j+b}$$

- Translation Invariance

Weight invariant relative to shift in point

$$\text{of view } y_{ij} = \sum_{a=-\Delta}^{\Delta} \sum_{b=-\Delta}^{\Delta} W_{ab} x_{i+a, j+b}$$

- Momentum

局部梯度值

Helps with local minima

→ 加速梯度下降 -

- Flat (noisy) gradients

$$m_t = (1-\lambda)m_{t-1} + \lambda g_t$$

$$w_t \leftarrow w_t - \eta_t g_t - \tilde{\eta}_t m_t$$

Can lead to oscillations for large momentum

梯度步

- Nesterov's accelerated gradient

加速梯度

$$m_{t+1} = \mu m_t + \varepsilon g_t (w_t - \mu m_t)$$

$$w_{t+1} = w_t - m_{t+1}$$

- Capacity Control

Minimizing loss can lead to overfitting

Weight Decay $w_t \leftarrow w_t - \eta_t g_t$

$$w_t \leftarrow (1-\lambda) w_t - \eta_t g_t$$

- Dropout / Avoid parameter sensitivity

Distributed representation

Randomized sparsification

6. Optimization

较大的LR会使参数更

新较大，前期收敛较慢，但后期会造

Learning Rate Decay

- Constant (requires schedule for piecewise constant, tricky)
- Useful hack Constant until no more improvement on validation set
- Polynomial decay

$$\eta(t) = \frac{\alpha}{\beta+t}$$

Recall exponent of 0.5 for conventional SGD, 1 for strong convexity. Bottou picks 0.75

- Exponential decay

$$\eta(t) = \alpha e^{-\beta t}$$

not recommended since quite aggressive

AdaGrad

自适应学习率衰减

- Adaptive learning rate (preconditioner)

$$\eta_{ij}(t) = \frac{\eta_0}{\sqrt{K + \sum_t g_{ij}^2(t)}}$$

- For directions with large gradient, decrease learning rate aggressively to avoid instability
- If gradients start vanishing, learning rate decreases, too
- Local variant

$$\eta_{ij}(t) = \frac{\eta_t}{\sqrt{K + \sum_{t'=\tau}^t g_{ij}^2(t')}}$$

Duchi, Hazan, Singer, 2010

<http://www.magicbroom.info/Papers/DuchiHaSi10.pdf>



74

time window

75

7. Memory & State

Autoregressive Models AR 自回归

Time series of observations

..., $x_{t-2}, x_{t-1}, x_t, x_{t+1}, x_{t+2}, \dots$

Estimate $x_{t+1} = f(x_t, \dots, x_{t-T})$ e.g. via deep net

- Latent state models

"记忆"

$$x_{t+1} = f(x_t, \dots, x_{t-T}, z_t, \dots, z_{t-T})$$

$$z_{t+1} = g(x_{t+1}, \dots, x_{t-T}, z_t, \dots, z_{t-T})$$

RNN

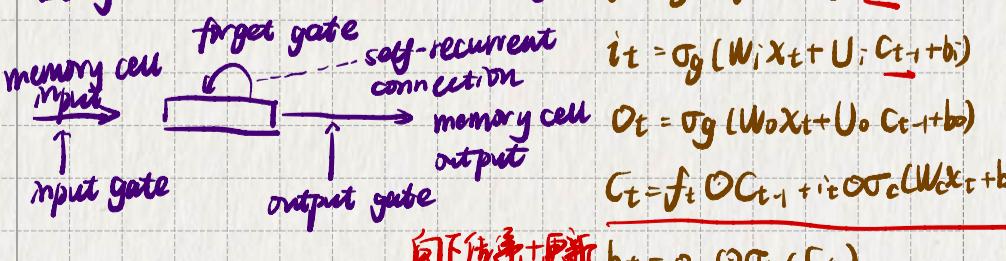
also called Recurrent Neural Network

LSTM

Long Short Term Memory

循环神经网络

$$f_t = \sigma_g(W_f x_t + U_f C_{t-1} + b_f)$$





Title

Regularization and Optimization for Deep Models

Rating: / / / / /

Framework of ML:

{ Training data $\{(x_1, y_1), \dots, (x_n, y_n)\}$ }{ Testing data $\{x^{N+1}, x^{N+2}, \dots, x^{N+m}\}$ }

Training

function with unknown - define loss function - optimization training data

Testing: Label and used in practice

Guideline

How to determine model bias?

• see right

• large

make your model complex

 $y = f(\theta, \theta^*, \theta^*(x))$
respectively

redesign and make it

flexible.

- too simple
more feature $y = b + \sum w_j x_j$ deep learning
(more neurons, layers) $y = b + \sum_i g_i \text{sigmoid}(a_i + \sum_j w_{ij} x_j)$

Loss on training data

optimization

advanced optimization



loss on training data

small

large

mismatch

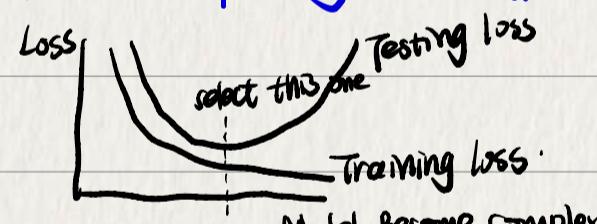
overfitting

Small loss on training data, large loss on testing data.

Create more training data
(data augmentation)

Limit model flexibility

Bias-Complexity Trade-off

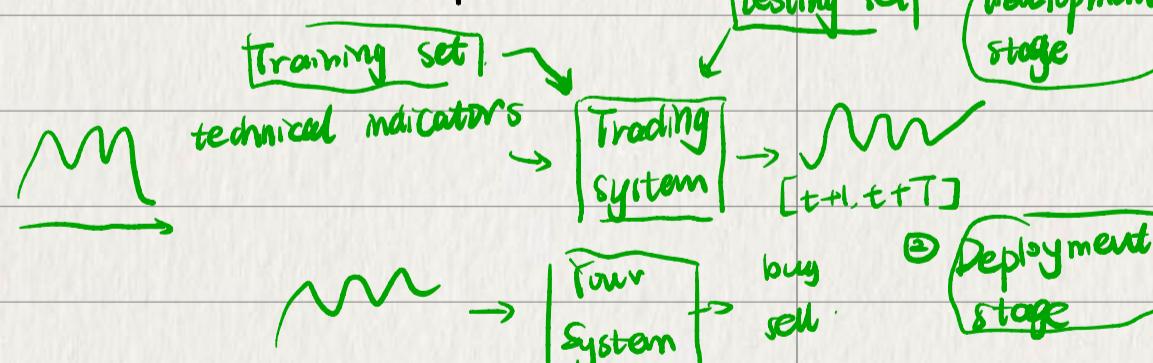
An extreme example $f(x) = \begin{cases} 0 & \exists x = x \\ \text{random} & \text{otherwise} \end{cases}$

zero training loss, large testing loss

Testing loss

Model become complex

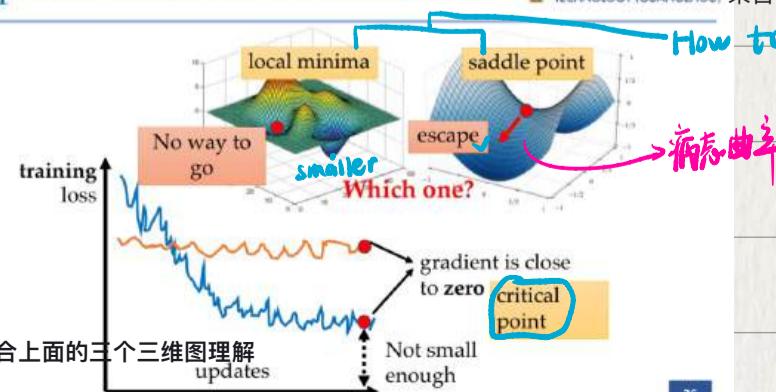
e.g.



Optimization Fails because

THE HONG KONG UNIVERSITY OF SCIENCE AND TECHNOLOGY (GUANGZHOU)

① 在附近的 loss 函数可以用泰勒展开式近似表示。下图中的纵坐标是 loss 值，横坐标是向量。一阶微分乘自变量差值是 y 轴上的绿色长度（可以用 tan 来理解），二阶微分乘两个自变量差值是 y 轴上的红色长度。

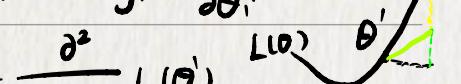
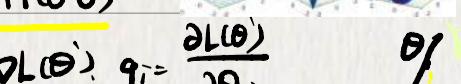


How to determine? ① Taylor Series Approximation

L(θ) around $\theta = \theta'$ can be approximated below

A critical point

$$L(\theta) \approx L(\theta') + (\theta - \theta')^T g + \frac{1}{2} (\theta - \theta')^T H(\theta) (\theta - \theta')$$

Gradient g is a vector, $g = \nabla L(\theta)$, $g_i = \frac{\partial L(\theta)}{\partial \theta_i}$ Hessian H is a matrix, $H_{ij} = \frac{\partial^2 L(\theta)}{\partial \theta_i \partial \theta_j}$ 当所有特征值都是正的, 矩阵 H 为正定矩阵时, 即红框值大于 0 时, $L(\cdot) > L(\theta')$, 那么 θ' 就是局部最小点。当特征值有正有负时, 即红框有时大于 0 有时小于 0, 则 θ' 就是鞍点。

$$\text{e.g. } y = w_1 w_2 x \xrightarrow{x=1} \begin{matrix} w_1 \\ w_2 \end{matrix} \xrightarrow{\otimes} \begin{matrix} w_1 \\ w_2 \end{matrix} \xrightarrow{\otimes} y \xleftarrow{=} \hat{y}$$

② Hessian: At critical point, $L(\theta) \approx L(\theta') + \frac{1}{2} (\theta - \theta')^T H (\theta - \theta')$

For all v , $v^T H v > 0 \rightarrow$ Around θ' : $L(\theta) > L(\theta')$

= H is positive definite = All eigenvalues are positive.

For all v , $v^T H v < 0 \rightarrow$ Around θ' : $L(\theta) < L(\theta')$

= H is negative definite = All eigenvalues are negative.

Sometimes $v^T H v > 0$, sometimes $v^T H v < 0 \rightarrow$ Saddle point

Some eigenvalues are positive, and some are negative.

e.g. Explain:

y 的 head是真实值，采用均方误差来计算 loss (L)，计算梯度即一阶微分使之为 0，求出参数 w_1 和 w_2 的值，此值对应位置就是临界点；然后将这两个值代入二阶微分求得 H 矩阵；再通过 H 矩阵去算特征值，算得两个特征值为 2 和 -2，所以该点是个鞍点。

如何确定参数更新的方向？

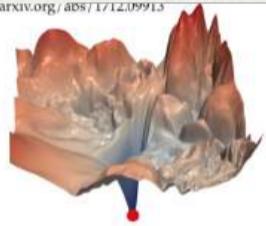
鞍点处的 H 矩阵中，特征值有正有负，这些正负能指导我们更新参数的方向。由下图中间式子推导可知，特征值小于 0，则红色框框里面的值始终小于 0，也就是说 $L(\theta) < L(\theta')$ ，这就是鞍点处 loss 能继续下降的方向。所以我们只要令 $v = u$ 即可，也就是 $-\lambda = u$, $\lambda = u$ ，此时处的 L 比处的 L 要低。

举例说明

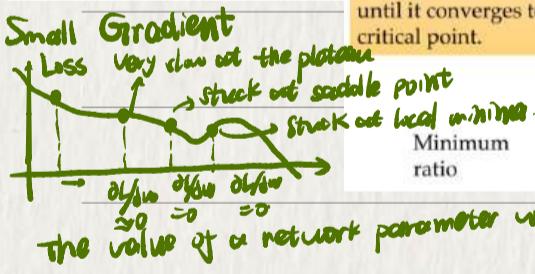
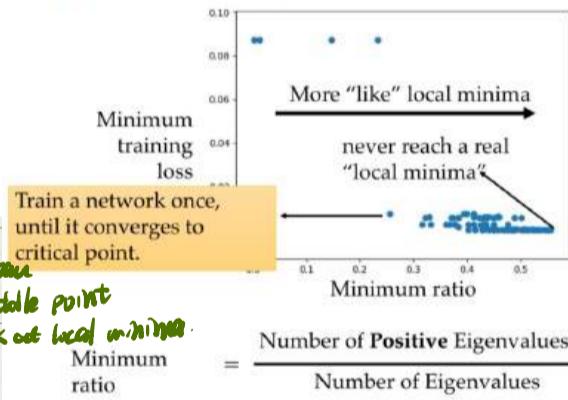
还是刚刚那个例子，求出负的特征值对应的特征向量，沿着特征向量的方向更新参数就可以减低 loss。特征向量有无穷多个，取一个即可。

上面讲的利用泰勒展开式判断临界点类型，并用 H 矩阵确定参数更新方向的方法。在实际中很少用到，因为二次微分和找特征值、特征向量，需要很大的运算量，不值得。

“方法”



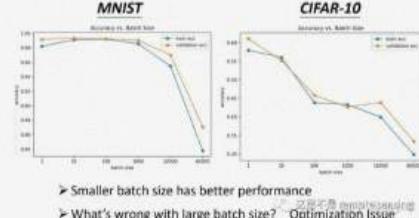
Empirical Study



大小批次的性能对比

小批次有更好的性能，由图可知同一个模型，同一个网络，training 误差随着 batch size 的增大而增大，testing 的误差也是。如果是 model bias 的问题，那么在 size 小的时候也会表现差，而不会等到 size 变大才差。所以这是 Optimization issue (优化问题) 导致大批次性能差。

Small Batch v.s. Large Batch

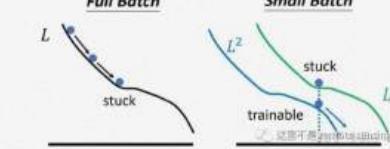


小批次训练时性能好

每次更新的时候，用的 loss 函数会有差异（如右图），因为不同的 batch 用不同的 loss function，换了个 loss 函数就更不容易卡住！！！

Small Batch v.s. Large Batch

- Smaller batch size has better performance
- “Noisy” update is better for training



总结

小批次在 1 个 epoch 中的速度很慢，耗时很长，但是小批次在优化过程性能更好，在测试时的性能也更好。而批次大小是一个超参数，需要自行设定。

$$L = (\hat{y} - w_1 w_2 x)^2 = (1 - w_1 w_2)^2$$

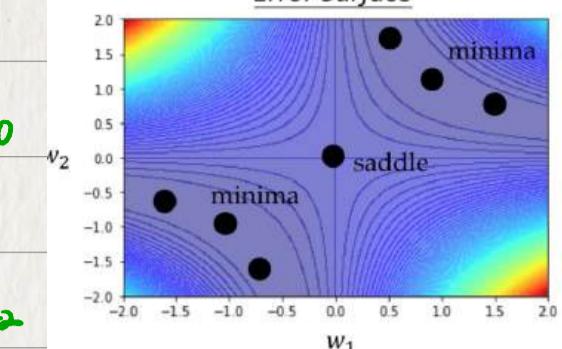
$$\textcircled{9} \frac{\partial L}{\partial w_1} = 2(1 - w_1 w_2)(-w_1) = 0$$

$$\frac{\partial L}{\partial w_2} = 2(1 - w_1 w_2)(-w_2) = 0$$

$$\textcircled{10} \frac{\partial^2 L}{\partial w_1^2} = 2(-w_2)(-w_1) = 2$$

$$\frac{\partial^2 L}{\partial w_2^2} = -2 + 4w_1 w_2 = -2$$

$$\frac{\partial^2 L}{\partial w_1 \partial w_2} = -2 + 4w_1 w_2 = -2$$



Critical point: $w_1 = 0, w_2 = 0$

$$H = \begin{bmatrix} 2 & -2 \\ -2 & 0 \end{bmatrix}, \lambda_1 = 2, \lambda_2 = -2$$

Saddle point 鞍点

Don't afraid of saddle point?

THE HONG KONG UNIVERSITY OF SCIENCE AND TECHNOLOGY (GUANGZHOU)

At critical point: $L(\theta) \approx L(\theta') + \frac{1}{2} (\theta - \theta')^T H (\theta - \theta') v^T H v$

Sometimes $v^T H v > 0$, sometimes $v^T H v < 0 \rightarrow$ Saddle point

H may tell us parameter update direction!

鞍点：损失函数值不再减小或增大
根据训练模型不会得到更好的结果

u is an eigenvector of $H \rightarrow u^T H u = u^T (\lambda u) = \lambda u^T u$

λ is the eigenvalue of u

$\lambda < 0$

$$L(\theta) \approx L(\theta') + \frac{1}{2} (\theta - \theta')^T H (\theta - \theta') \rightarrow L(\theta) < L(\theta')$$

$$\theta - \theta' = u \quad \theta = \theta' + u \quad \text{Decrease } L$$

$$\lambda_2 = -2 \quad \text{Has eigenvector } u = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Update the parameter along the direction of u

You can escape the saddle point and decrease the loss.

(this method is seldom used in practice)

Batch

对抗临界点方法②

Optimization with Batch

THE HONG KONG UNIVERSITY OF SCIENCE AND TECHNOLOGY (GUANGZHOU)

$\theta^* = \arg \min_{\theta} L$
 ➤ (Randomly) Pick initial values θ^0
 ➤ Compute gradient $g^0 = \nabla L(\theta^0) L$
 ➤ update $\theta^1 \leftarrow \theta^0 - \eta g^0$
 ➤ Compute gradient $g^1 = \nabla L^2(\theta^1) L$
 ➤ update $\theta^2 \leftarrow \theta^1 - \eta g^1$
 ➤ Compute gradient $g^2 = \nabla L^3(\theta^2) L$
 ➤ update $\theta^3 \leftarrow \theta^2 - \eta g^2$
 1 epoch = see all the batches once
 ➤ Shuffle after each epoch

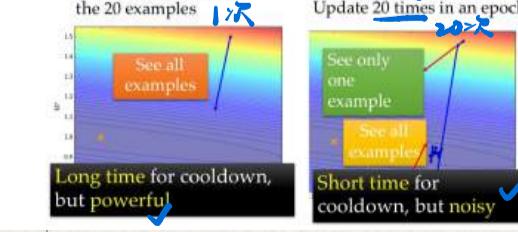
epoch batch
将一批大型资料若干次打乱，然后重新分批次
loss 和梯度，从而更新参数，每有一个 epoch 就把这一批大型资料打乱，然后重新分批次
→ 保证每个 epoch 中 batch 资料不同，避免偶然性。

Small Batch vs. Large Batch

THE HONG KONG UNIVERSITY OF SCIENCE AND TECHNOLOGY (GUANGZHOU)

大批量
Batch size = N (Full batch)
Update after seeing all the 20 examples 1次

小批次
Batch size = 1
Update for each example
Update 20 times in an epoch 20次



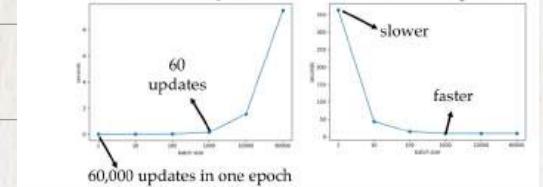
Small Batch v.s. Large Batch

	Small	Large
Speed for one update (no parallel)	Faster	Slower
Speed for one update (with parallel)	Same	Same (not too large)
Time for one epoch	Slower	Faster
Gradient	Noisy	Stable
Optimization	Better	Worse
Generalization	Better	Worse

更新速度差不多，小批次耗时长，大批量耗时短。

没有平行计算时，单次更新大批量耗时更长，有平行计算时，单次更新大批量耗时不变。
时间不长，而 1 epoch 中大批量耗时更短。

Smaller batch requires longer time for one epoch (longer time for seeing all data once)



BN: $\sigma^2 = \text{Var}[x_i], \beta = E[x_i]$, 保证整个 network 的 capacity 还原原来输入

Notes: 25 50 75 100

② How: chain rule ③ Where: 在非线性层的输出，对 $x = W_1 + b_1 + \dots + b_n$

④ Why: 克服深度神经网络训练时的梯度消失问题，但均值差和方差差分布很广。

Batch Normalization

Feature Normalization Adaptive Learning Rate

$x_i^r \leftarrow \frac{x_i - \mu}{\sigma} \sim N(0, 1)$

In general, feature normalization makes landscape gradient descent converge faster.

Training stuck ≠ Small Gradient

People believe training stuck because the parameters are around a critical point...

Different parameters need different learning rate

Considering Deep Learning

Same μ, σ . Different dims have different ranges.

Batch normalization - Testing

Computing the moving average of μ and σ of the batches during training. $\bar{\mu} \leftarrow p\bar{\mu} + (1-p)\mu^t$

Others: Layer/ Instance/ Group/ Weight/Spectrum Normalization

Momentum ↗ 对于全局最优点方法

(Vanilla) Gradient Descent

$\theta^0 = \theta^0 - g^0 - m^0 = \theta^0 - \eta g^0 \rightarrow \theta^1 = \theta^0 + m^0 - g^0 - m^1 = \theta^0 - \eta g^0 - \eta^2 g^1 \rightarrow \theta^2 = \theta^0 + m^0 - g^0 - m^1 - g^1 - m^2 = \theta^0 - \eta g^0 - \eta^2 g^1 - \eta^3 g^2$

movement = Negative of $\partial L / \partial w$ + Last Movement

Loss

① Critical points have zero gradients.

Summary

② Critical points can be either saddle points or local minima. Can be determined by the Hessian matrix.

- It is possible to escape saddle points along the direction of eigenvectors of the Hessian matrix.
- Local minima may be rare.

③ Smaller batch size and momentum help escape critical points

那 BN 到底是什么原理呢？说到底还是为了防止“梯度弥散”。关于梯度弥散，大家都知道一个简单的栗子： $0.9^{10} \approx 0.04$ 。在 BN 中，是通过将 activation 规范化为均值和方差一致的手段使得原本会减小的 activation 的 scale 变大。可以说是一种更有效的 local response normalization 方法（见

⑤ 例如，在神经网络训练时遇到收敛速度很慢，或梯度爆炸等无法训练的状况时可以尝试 BN 来解决。另外，在一般使用情况下也可以加入 BN 来加快训练速度，提高模型精度。

Training stuck ≠ Small Gradient

Formulation for one parameter: $\theta_i^{t+1} \leftarrow \theta_i^t - \eta g_i^t$

$$g_i^t = \frac{\partial L}{\partial \theta_i} \Big|_{\theta=\theta^t}, \theta_i^{t+1} \leftarrow \theta_i^t - \frac{\eta}{\sigma_i^t} g_i^t$$

parameter dependent

Root Mean Square

$$\theta_i^{t+1} \leftarrow \theta_i^t - \frac{\eta}{\sigma_i^t} g_i^t$$

$$\sigma_i^t = \sqrt{\frac{1}{t+1} \sum_{i=0}^t (g_i^t)^2}$$

Used in Adagrad

Learning rate adapts dynamically

Error surface can be very complex.

① 不需要手动设置学习率参数，由算法自动完成
② 可以为每个参数选择不同学习率

RMSProp

结合了梯度方向梯度和历史梯度信息

Learning Rate Scheduling

Learning Rate Decay:
As the training goes, we are closer to the destination, so we reduce the learning rate.

RMSprop Momentum

Algorithm 1: Adam, our proposed algorithm for stochastic optimization. See section 2 for details, and for a slightly more efficient (but less clear) order of computation, g_i^t indicates the elementwise square $g_i \odot g_i$. Good default settings for the tested machine learning problems are $\alpha = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$. All operations on vectors are element-wise. With β_1^t and β_2^t we denote β_1 and β_2 to the power t .

```

Require:  $\alpha$ : Stepsize
Require:  $\beta_1, \beta_2 \in [0, 1]$ : Exponential decay rates for the moment estimates
Require:  $f(\theta)$ : Stochastic objective function with parameters  $\theta$ 
Require:  $\theta_0$ : Initial parameter vector
 $m_0 \leftarrow 0$  (Initialize 1st moment vector)
 $v_0 \leftarrow 0$  (Initialize 2nd moment vector)
 $t \leftarrow 0$  (Initialize timestep)
while  $\theta_t$  not converged do
     $t \leftarrow t + 1$ 
     $g_t \leftarrow \nabla f_t(\theta_{t-1})$  (Get gradients w.r.t. stochastic objective at timestep  $t$ )
     $m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$  (Update biased first moment estimate)
     $v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$  (Update biased second raw moment estimate)
     $\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$  (Compute bias-corrected first moment estimate)
     $\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$  (Compute bias-corrected second raw moment estimate)
     $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$  (Update parameters)
end while
return  $\theta_t$  (Resulting parameters)

```

将 Momentum 和 RMSprop 结合

Summary of Optimization

(Vanilla) Gradient Descent: $\theta_i^{t+1} \leftarrow \theta_i^t - \eta g_i^t$

Various Improvements

- Learning rate scheduling
- Momentum: weighted sum of the previous gradients
- Root mean square of the gradients only magnitude
- Consider direction



Title

Chapter

Date

Review

/ / / / /

Rating

A row of five empty circles, each with a thin black outline, intended for a child to draw a face in.



Title

Chapter

Date

Review

/ / / / /

Rating

A row of five empty circles, each enclosed in a thin black outline, arranged side-by-side.

Summary

Summary



Title

Chapter

Date

Review

/ / / / /

Rating

○ ○ ○ ○ ○



Title

Chapter

Date

Review

/ / / / /

Rating

A row of five empty circles, each enclosed in a thin black border, intended for a student to draw a shape inside.

Summary

Summary



Title

Chapter

Date

Review

Rating



Title

Chapter

Date

Review

Rating

A horizontal row of five empty speech bubbles, each containing a single diagonal slash (/), used for recording audio responses.



Chapter 01

Keywords / 关键词

Important and difficult points / 重点难点

Date / 复习日期	Total Time / 复习时长	Efficiency / 学习效率	Note / 备注
		优秀 <input type="radio"/> 良好 <input type="radio"/> 一般 <input type="radio"/> 差 <input type="radio"/>	
		优秀 <input type="radio"/> 良好 <input type="radio"/> 一般 <input type="radio"/> 差 <input type="radio"/>	
		优秀 <input type="radio"/> 良好 <input type="radio"/> 一般 <input type="radio"/> 差 <input type="radio"/>	
		优秀 <input type="radio"/> 良好 <input type="radio"/> 一般 <input type="radio"/> 差 <input type="radio"/>	
		优秀 <input type="radio"/> 良好 <input type="radio"/> 一般 <input type="radio"/> 差 <input type="radio"/>	
		优秀 <input type="radio"/> 良好 <input type="radio"/> 一般 <input type="radio"/> 差 <input type="radio"/>	
		优秀 <input type="radio"/> 良好 <input type="radio"/> 一般 <input type="radio"/> 差 <input type="radio"/>	
		优秀 <input type="radio"/> 良好 <input type="radio"/> 一般 <input type="radio"/> 差 <input type="radio"/>	