

Neural Network Basics

0100_kNN

- What is Image Classification
- Nearest Neighbor Classifier
- Validation Set
 - Hyperparam tuning
 - Cross validation

0200_Linear Classification

- Overview
- 线性分类器的直观理解
- Loss
 - SVM Loss (hinge loss)
 - Regularization
 - Softmax (cross entropy loss)

为了计算中的梯度, 通常写成向量形式

0300_Optimization

- Strategy
 - 1. Random Search
 - 2. Random Local Search
 - 3. Gradient
- Compute the gradient
 - Numerical
 - 用于梯度检查
 - centered difference formula
 - Analytical
- Gradient descent
 - Vanilla version (单样本更新)
 - Mini-batch

0400_Backpropagation

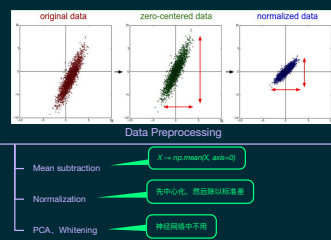
用相似设置比较两种梯度

$$\frac{df(x)}{dx} = \frac{f(x+h) - f(x-h)}{2h}$$

0500_NN1

- Intro to Neural Network
- Activation functions
 - Sigmoid
 - Tanh
 - ReLU
 - Leaky ReLU
- Neural Network architectures

0600_NN2



Weight Initialization

```
w = np.random.randn(n) * sqrt(2./n)
b = 0
```

Batch Normalization

Input: Values of x over a mini-batch: $B = \{x_{1..m}\}$;
Parameters to be learned: γ, β
Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_B \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$
$$\sigma_B^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2 \quad // \text{ mini-batch variance}$$
$$\hat{x}_i \leftarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \quad // \text{ normalize}$$
$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.

Regularization

- L2 regularization: $\frac{1}{2} \|w\|^2$
- L1 regularization: $\lambda \|w\|_1$
- Dropout

```
def train_step(b):
    # forward pass for example 3-layer neural network
    H1 = np.maximum(0, np.dot(W1, X) + b1) #np.dot, np.maximum
    U1 = (np.random.rand(*H1.shape) < p) / p # first dropout mask.
    H1 *= U1 # dropout
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    U2 = (np.random.rand(*H2.shape) < p) / p # second dropout mask.
    H2 *= U2 # dropout
    out = np.dot(W3, H2) + b3

def predict(x):
    # ensemble forward pass
    H1 = np.maximum(0, np.dot(W1, X) + b1) # no scaling necessary
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    out = np.dot(W3, H2) + b3
```

各种Loss

0700_NN3

- Parameter updates
 - # Vanilla update
 - # Momentum update
 - # Nesterov Momentum
- Annealing learning rate
- Per-parameter adaptive learning rate
 - Adagrad
 - RMSprop
 - Adam
- Master and worker
- Model ensemble

动量参数更新的方法是保持动量参数使用训练时的学习率

```
# Vanilla update
x += - learning_rate * dx

# Momentum update
v = dv + learning_rate * dx # integrate velocities
x += v # integrate position

# Nesterov Momentum
x_ahead = x + v * h
# evaluate dx_ahead (the gradient at x_ahead instead of x)
p = - dv + learning_rate * dx_ahead
x += v
```

Annealing learning rate

Per-parameter adaptive learning rate

```
# Adaptive the gradient dx and parameter vector x
cache = dx**2

# Adagrad
x += - learning_rate * dx / (np.sqrt(cache) + epsilon)

# RMSprop
cache = decay_cache * cache + (1 - decay_rate) * dx**2
x += - learning_rate * dx / (np.sqrt(cache) + epsilon)

# Adam
m = np.zeros_like(x)
v = np.zeros_like(x)
m_hat = m / (1 - beta1**t)
v_hat = v / (1 - beta2**t)
x += - learning_rate * m_hat / (np.sqrt(v_hat**2 + epsilon))
```

单维梯度的累加, 可能导致学习以更快的速度

增加了一个衰减系数, 会随你一些之前累加的梯度

与RMSprop相比, 把dx乘了一个更平滑的项
= 对dx看做是随机的, 做的一阶和二阶估计

