

亮点: } 可以建模  $p(\tilde{y}|x, y)$  的噪声  
加一层 softmax 就可以.

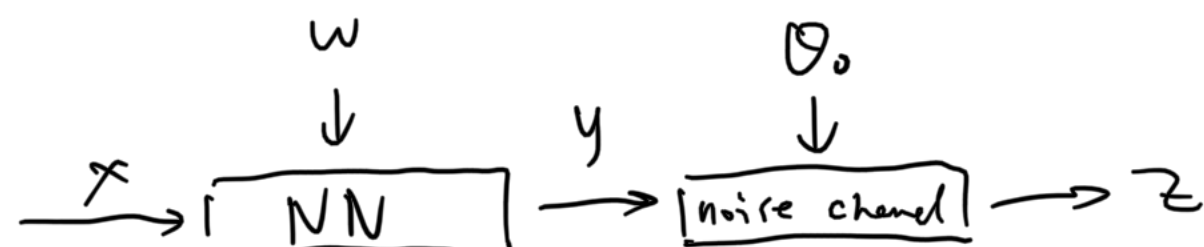
## EM 算法 - 利用 EM 算法.

Intro:

以前的算法利用  $\tilde{y}$  和  $y$  关系为 asymmetric 将真实标签作为变量.

$$P(z = \tilde{y} | x; w, \theta) = \sum_{\tilde{y}=1}^K P(z = \tilde{y} | y = \tilde{y}; \theta) P(y = \tilde{y} | x; w)$$

假设



函数不好优化,  $w$  不好求出.

用 EM 算法求解. 即

$$\arg \max_w P(y_t = \tilde{y} | x_t, z_t, w_0, \theta_0) \cdot P(y | x; w) \Rightarrow$$

$$\left\{ \begin{array}{l} \theta \text{ 有显式表达 } \theta(c_{ij}) = \frac{\sum_t c_{ti} 1_{\{z_t = j\}}}{\sum_t c_{ti}} \\ w \text{ 通过梯度下降 } \frac{\partial \log \mathcal{L}}{\partial w_i} \end{array} \right.$$

(变种) ① 不学习  $\theta$ , 设置一个验证集手工标.

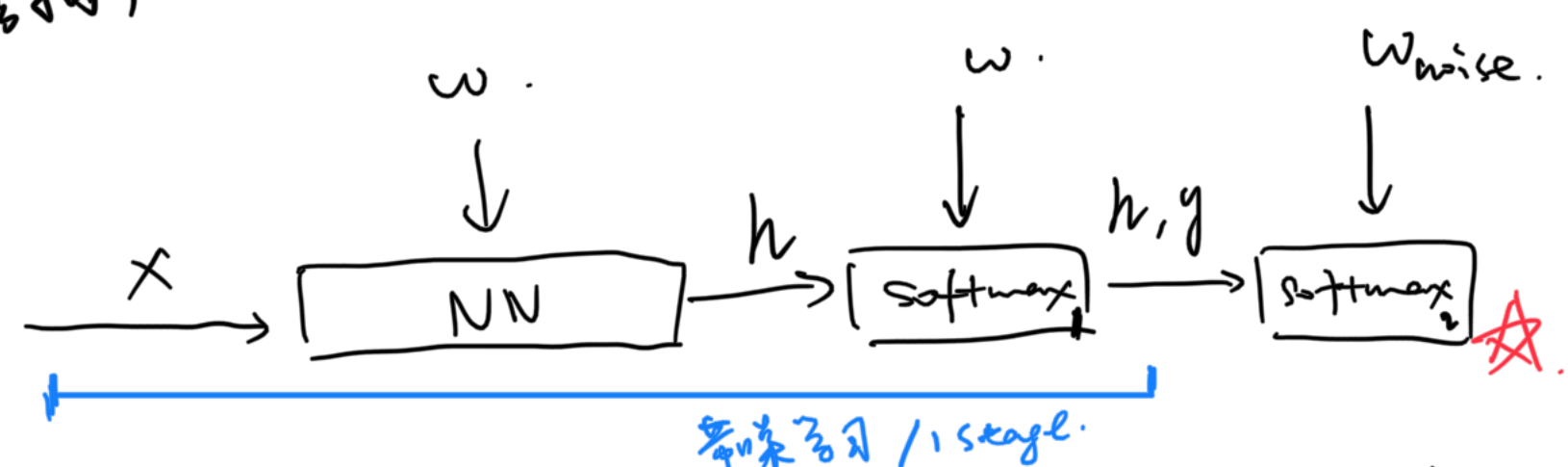
② 蒸馏学习, 先用大网络蒸馏学习. 后将前面网络参数  $w_0$  用于初始化小网络.

(缺点): EM 算法易陷入局部最优. EM 不易收敛. asymmetric 标签假设太强.

## 提出算法 - 只基于 DNN 的梯度下降

Method:

(结构)



S-model: 假设  $\tilde{y}$  是 asymmetric 的. 所以,  $w_{noise}$  是不变参数  $u_{ij}$  的.

C-model: 假设  $\tilde{y}$  是 instance-dependent 的.  $w_{noise}$  含参数  $u_{ij}$ .

$$\text{softmax} = \frac{\exp(u_{ij}h + b)}{\sum_i \exp(u_{ij}h + b)}$$

(初始化)

$w_0$  为 1 stage 的模型参数  $w_{noise}$  }  $u_{ij}$  置 0.

$b_{ij}$ . 1 stage - 混淆矩阵

(步骤)

① 首先带噪训练 DNN. 后假设  $\hat{\tilde{y}}$  为  $y$ . 初始化出 T, 即 softmax<sub>2</sub> 的偏置参数

② DNN + softmax<sub>2</sub> 再训练一次. 如果 S-model 就使  $u_{ij}$  恒为 0

(优点)

① 利用一层 softmax 显式建模  $p(\tilde{y}|y, x)$ . 并且初始化使得模型可能学到  $p(y|x)$ .

$$\arg \max_{w, w_{noise}} \underbrace{p(y|x, w)}_{\text{hidden layer + softmax}} + \underbrace{p(z|y, x; w_{noise})}_{\text{通过一层 softmax 建模}}$$

假设第  $n$  次变化小. 那么  $p(y|x, w)$  和  $p(z|y, x, w, w_{noise})$  是 statistically consistent 的

② 可扩展性好. 仅需要做: ① 统计 T ② 加一层 softmax.

## 实验

1. 比较对象.

baseline: 带噪训练的 DNN.

read-soft:  $\beta = 0.95$ .

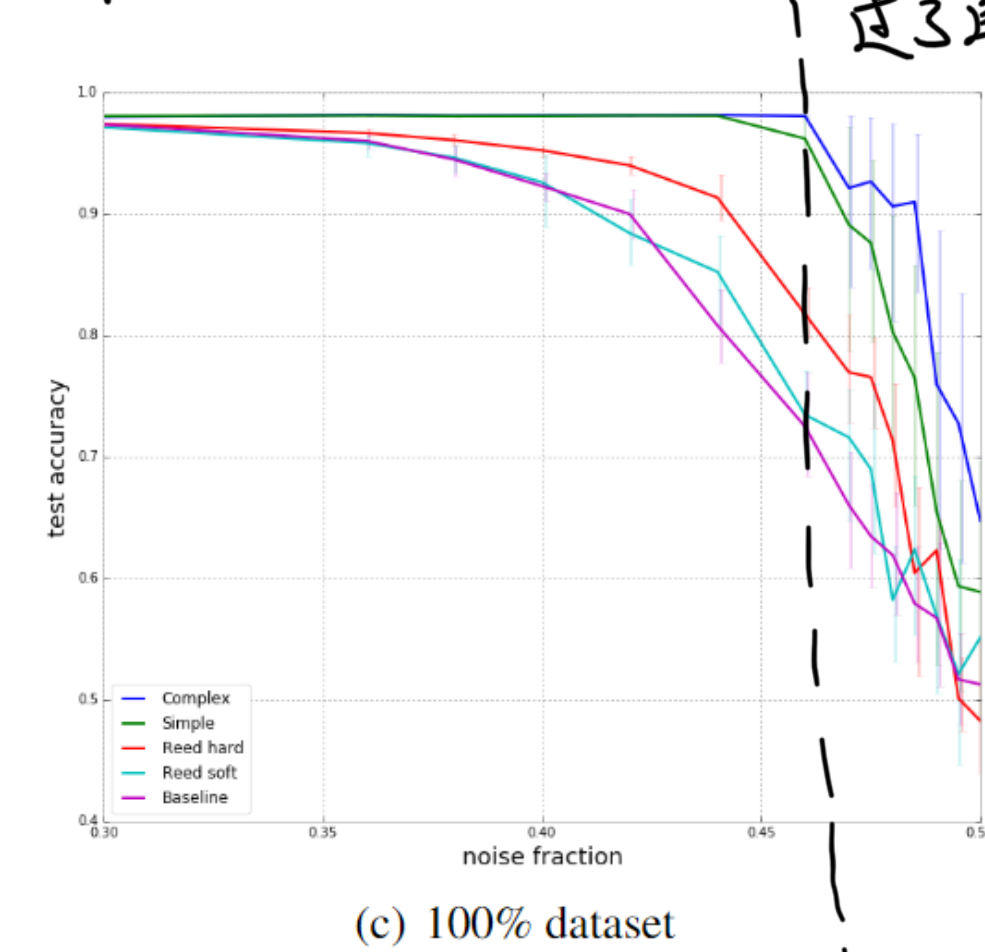
read-hard:  $\beta = 0.8$

S-model

C-model.

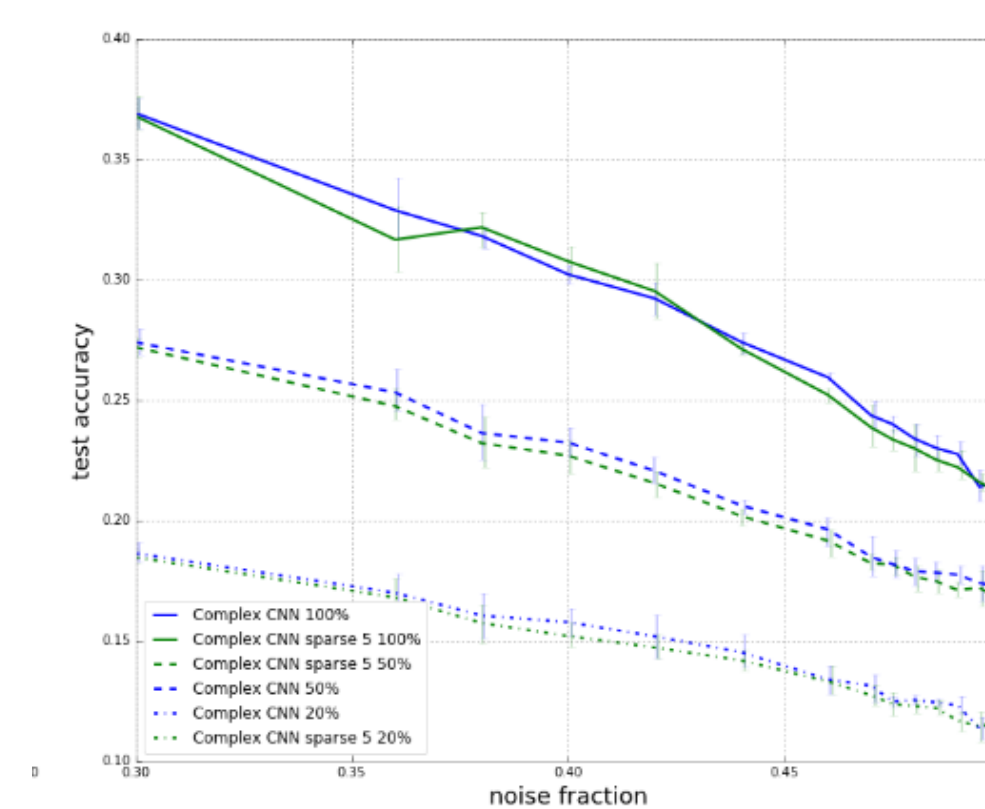
2.

MNIST.



过了这个点, 准确度下降? 值得研究.

CIFAR-100.



在类别数量大 (100) 的情况下.  
T 只记录每一类前 5 名最易转移的标签概率.  
↓  
so dense 初始化相比没有较大性能的下降.