

神经网络的训练技巧

主讲：王皓 副研究员| 硕士导师

邮箱：wanghao3@ustc.edu.cn

主页：http://staff.ustc.edu.cn/~wanghao3



本章内容

- 归一化 (normalization)
 - 数据归一化
 - 逐层归一化

- 训练目标和方式
 - 多任务、迁移学习、课程学习

- 数据增强

- 正则化
 - Early stop
 - Dropout

归一化



归一化

➤数据归一化

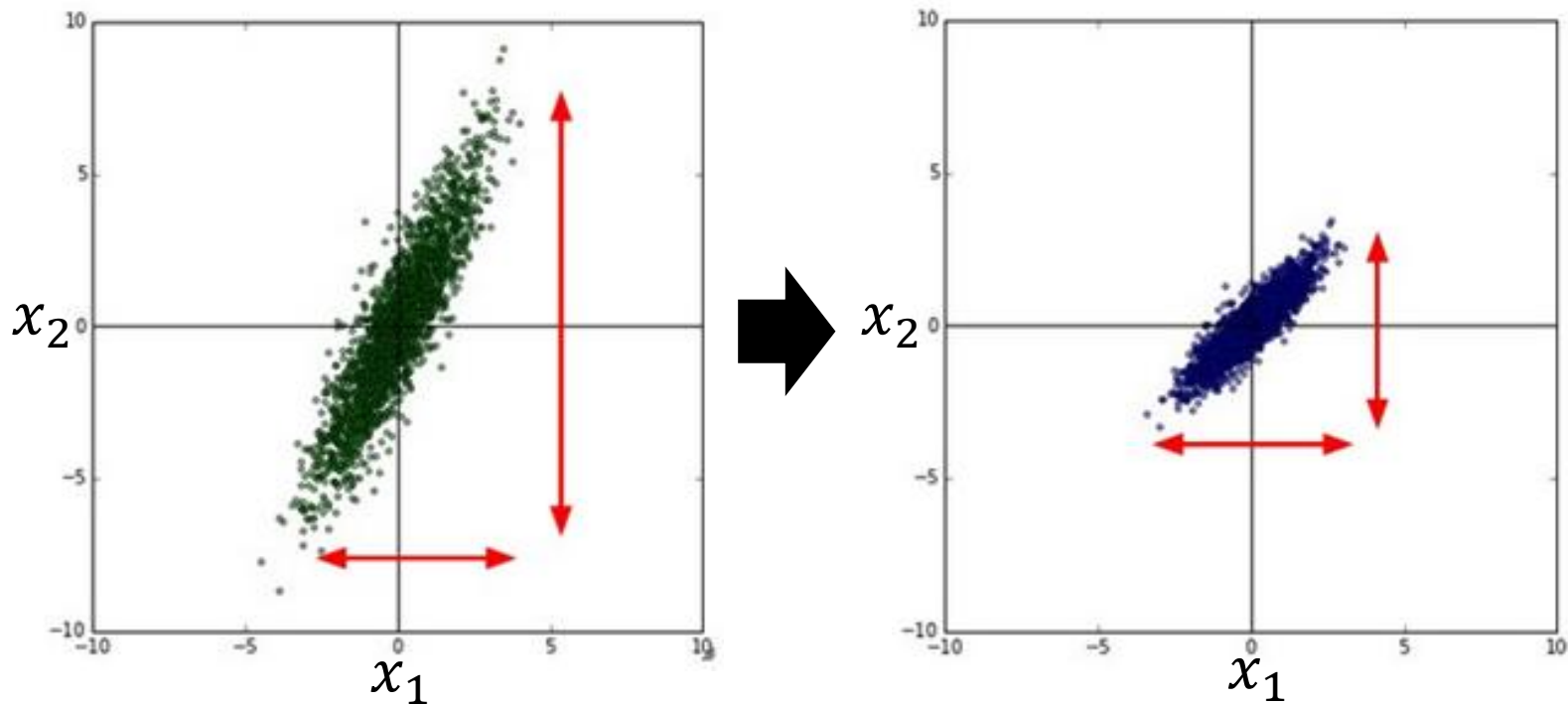
- 标准化 (Standardization) / Z值归一化 (Z-Score Normalization)
- 最大最小值归一化 (Min-Max Normalization)

➤逐层归一化

- Batch Normalization, BN
- Layer Normalization, LN
- Instance Normalization, IN
- Group Normalization, GN

数据归一化

$$\hat{y} = b + w_1x_1 + w_2x_2$$

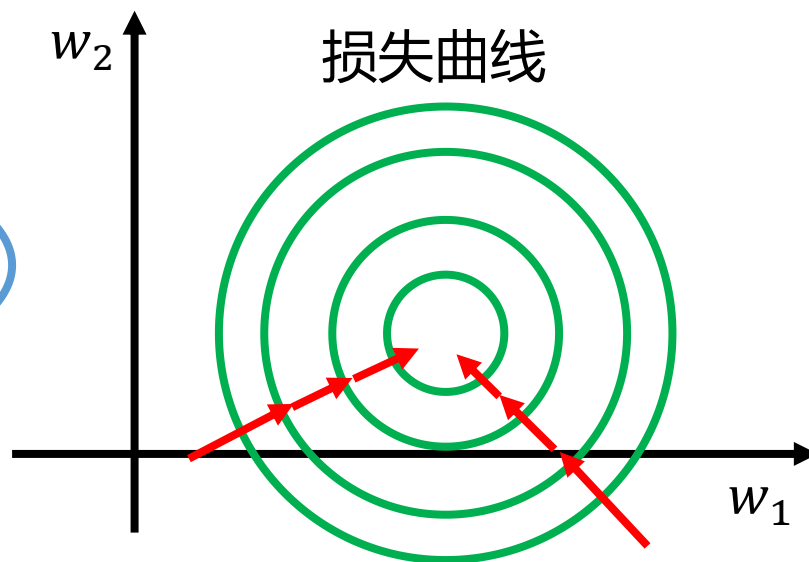
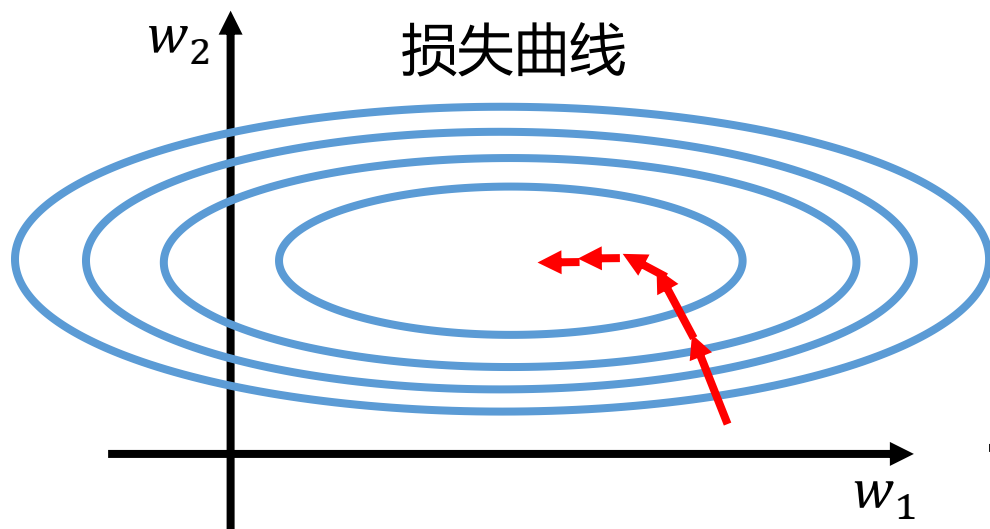
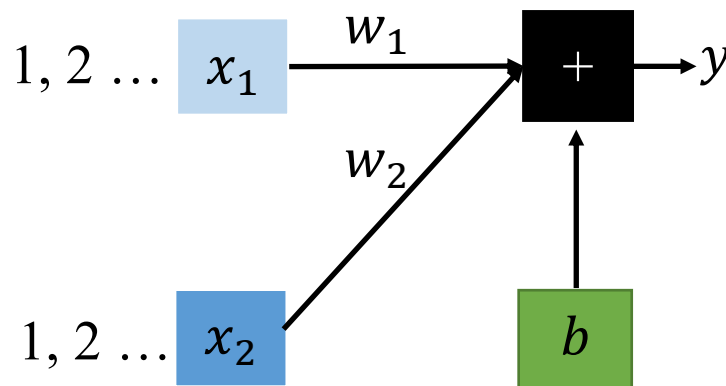
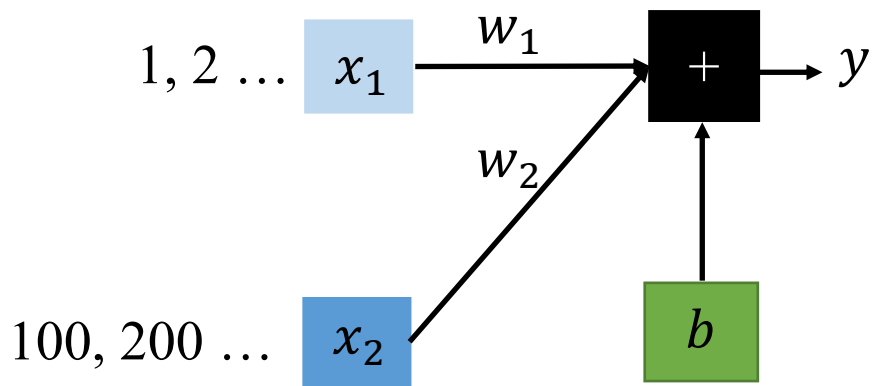


数据归一化：使得不同的特征有相同的尺度

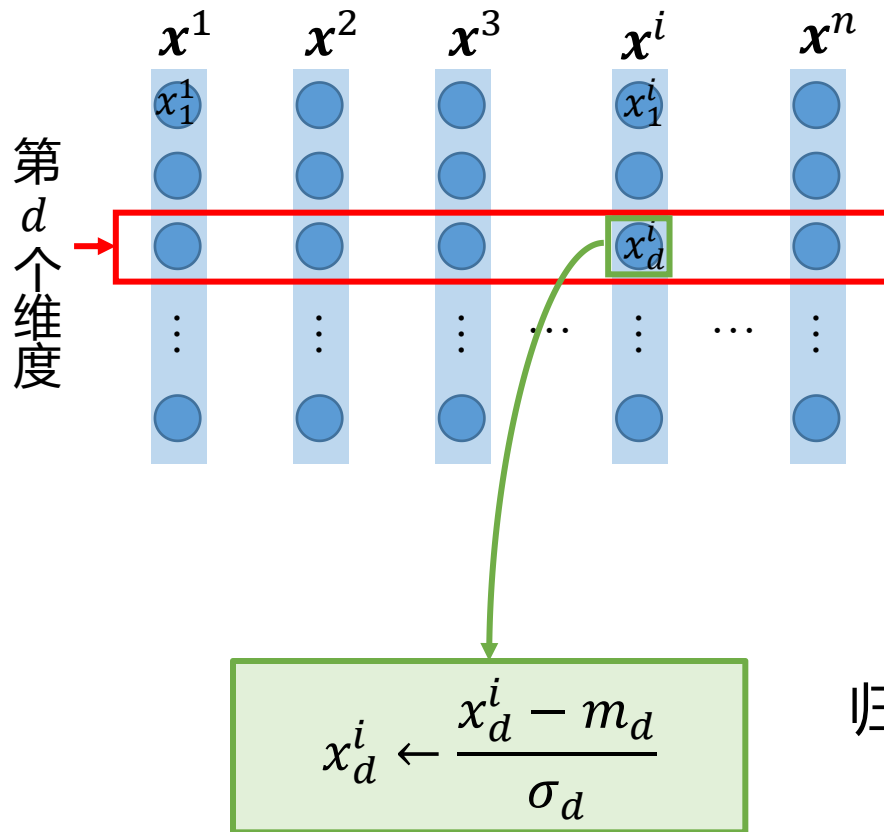
数据归一化

$$\hat{y} = b + w_1x_1 + w_2x_2$$

$$\mathcal{L} = (y - \hat{y})^2$$



数据归一化——标准化

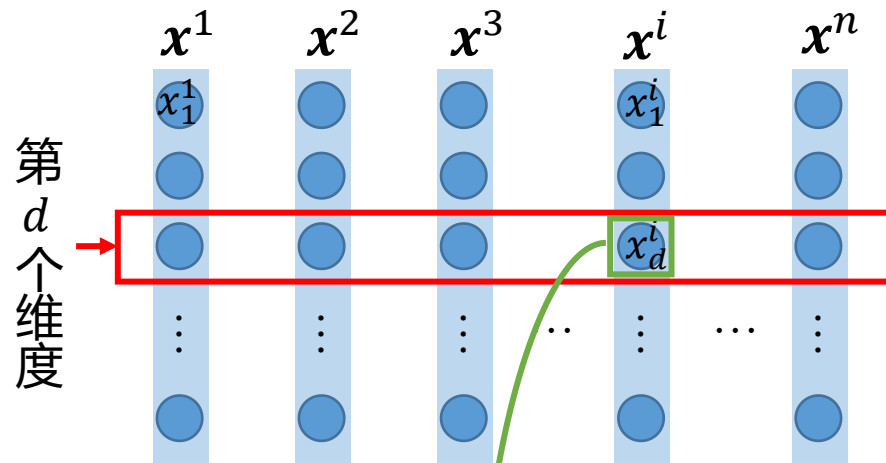


$$\text{计算均值 } m_d = \frac{1}{n} \sum_{i=1}^n x_d^i$$

$$\text{计算方差 } \sigma_d^2 = \frac{1}{n} \sum_{i=1}^n (x_d^i - m_d)^2$$

归一化后数据的所有维度均为
0均值 1方差

数据归一化——最大最小值归一化



计算最大值 \max_d

计算最小值 \min_d

$$x_d^i \leftarrow \frac{x_d^i - \min_d}{\max_d - \min_d}$$

归一化后数据的所有维度均在 $[0,1]$

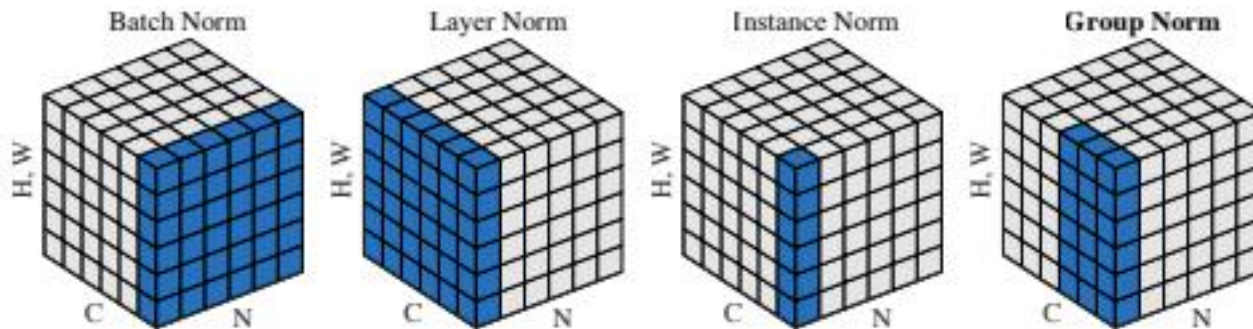
逐层归一化

➤ 更好的尺度不变性：

- 如果一个神经层的输入分布发生了改变，那么其参数需要重新学习，这种现象叫作**内部协变量偏移**
- 为了缓解上述问题，可以对每一个**神经层的输入**进行归一化操作，使其分布保持稳定

➤ 更平滑的优化地形：

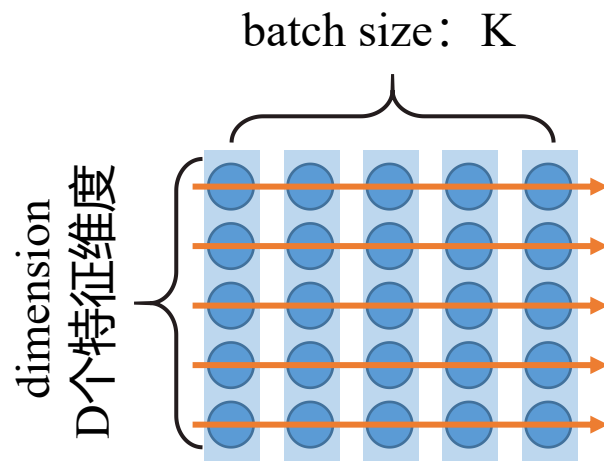
- 使得大部分神经层的输入处于不饱和区域，从而让梯度变大，**避免梯度消失**问题
- 使梯度变得更加稳定，从而允许使用**更大的学习率**，并**提高收敛速度**



Batch Normalization (BN)

神经网络中第 l 层的净输入为 \mathbf{z}^l ，神经元的输出为 \mathbf{a}^l ，

$$\mathbf{a}^{(l)} = f_l(\mathbf{z}^{(l)}) = f_l(\mathbf{W}^{(l)} \mathbf{a}^{(l-1)} + \mathbf{b}^{(l)})$$



1. 给定一个包含K个样本的小批量样本集合，独立计算该层每维度的均值和方差

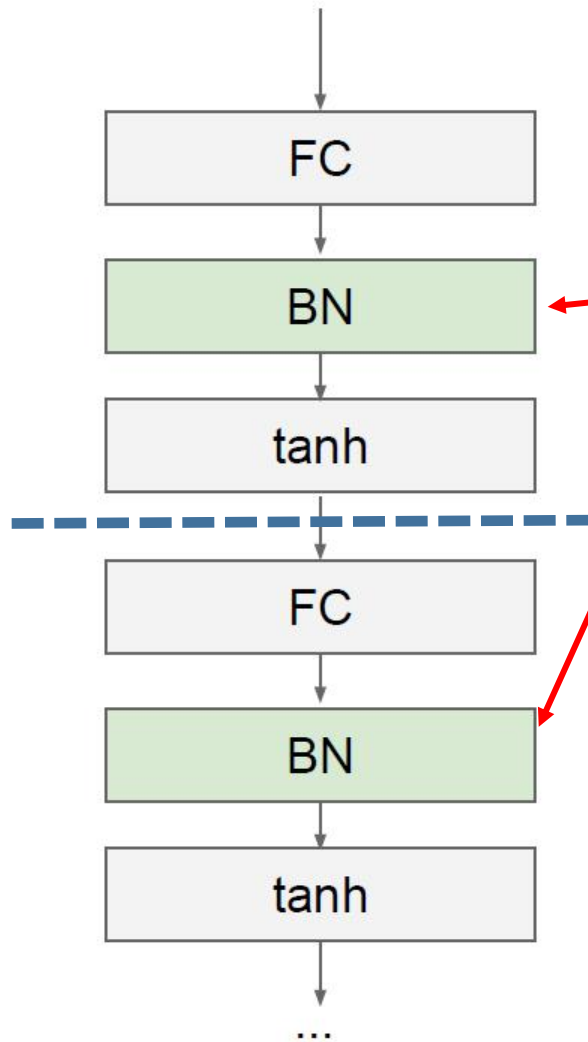
$$\mu_B = \frac{1}{K} \sum_{k=1}^K \mathbf{z}_k$$

$$\sigma_B^2 = \frac{1}{K} \sum_{k=1}^K (\mathbf{z}_k - \mu_B) \odot (\mathbf{z}_k - \mu_B)$$

2. 归一化

$$\mathbf{z}_k = \frac{\mathbf{z}_k - \mu_B}{\sigma_B}$$

Batch Normalization (BN)



BN层一般在仿射变换之后，
非线性激活函数之前

- 可以应用在深度神经网络中的任何一个**中间层**，并不需要对所有层进行归一化
- 每层的输入通常是前一层的非线性激活函数的输出
- 由于非线性激活函数，每层输入通常是非高斯的，不适合进行标准化



Batch Normalization (BN)

- 标准化单元的均值和标准差会降低该单位的表达能力
- 为了使得归一化不对网络的表示能力造成负面影响，引入 γ 和 β 两个可学习参数使得各单元的净输出有任意均值和方差

$$z_k = \frac{z_k - \mu_B}{\sigma_B} \quad \rightarrow \quad z_k = \frac{z_k - \mu_B}{\sigma_B} \odot \gamma + \beta$$

- 网络可能学习出这两个参数，即得到一个恒等映射

$$\gamma = \sigma_B \quad \beta = \mu_B$$

旧参数

- μ_B 和 σ_B 取决于低层神经网络的复杂关联

新参数

- γ 和 β 解除了与下层计算的密切耦合，直接通过梯度下降来学习

Batch Normalization (BN)

- 测试过程和训练过程有所不同：均值和方差不是在测试集的 mini-batch 上计算，而是整个数据集上的均值和方差。
- 实践中，可以通过移动平均来计算样本均值和方差

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots x_m\}$;

Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

BN的作用

- 改善梯度的传播
- 允许比较大的学习率
- 减小初始化的影响
- 一种隐形的正则化方法

- batch size不能太小，否则会导致算法性能下降
- 对于样本维度不统一的数据无法很好处理，例如RNN



Batch Normalization (BN)

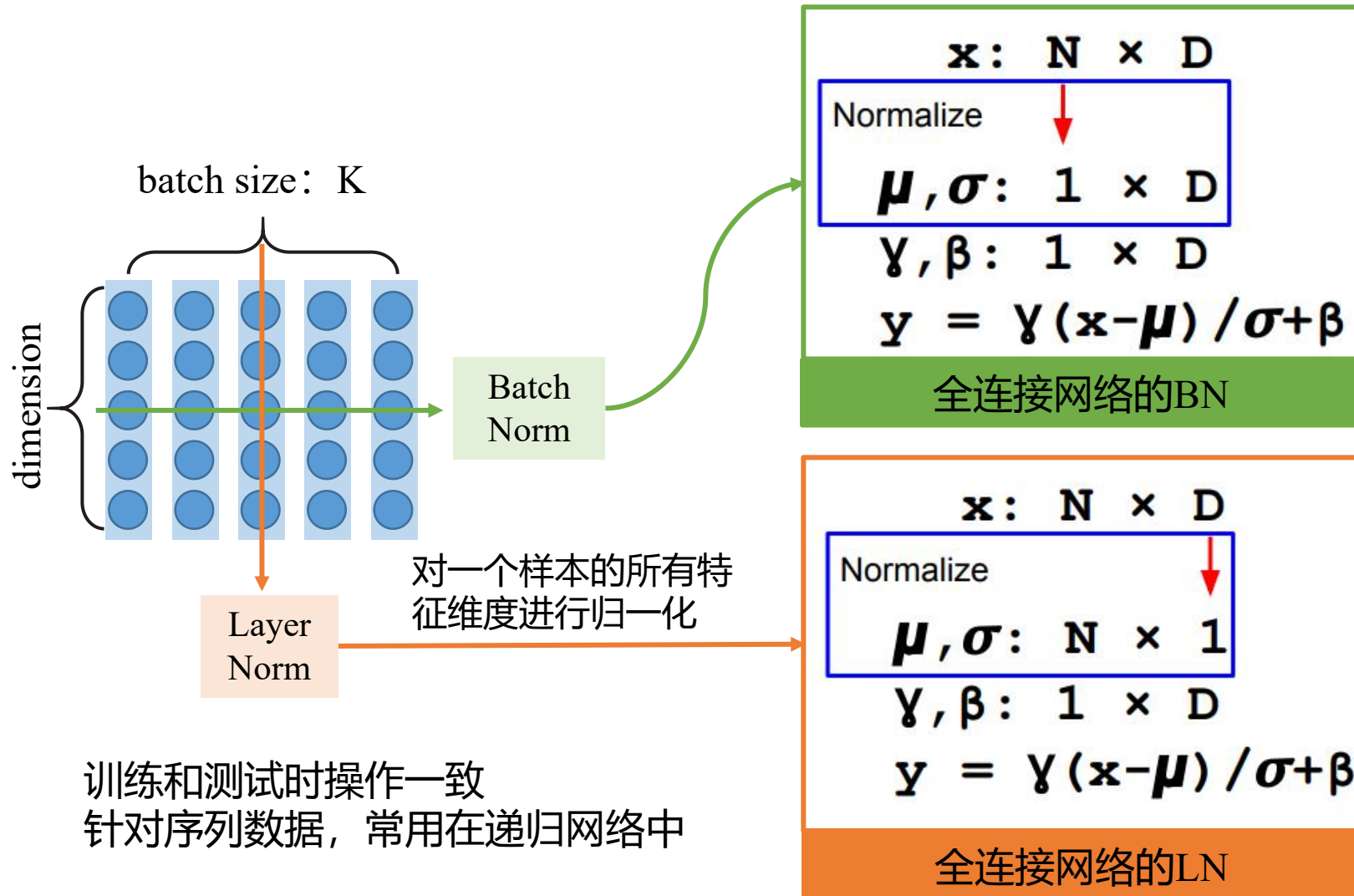
全连接网络的BN

$$\begin{aligned} \mathbf{x} &: \mathbf{N} \times \mathbf{D} \\ \text{Normalize} & \quad \downarrow \\ \boldsymbol{\mu}, \boldsymbol{\sigma} &: \mathbf{1} \times \mathbf{D} \\ \boldsymbol{\gamma}, \boldsymbol{\beta} &: \mathbf{1} \times \mathbf{D} \\ \mathbf{y} &= \boldsymbol{\gamma}(\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta} \end{aligned}$$

卷积网络的BN

$$\begin{aligned} \mathbf{x} &: \mathbf{N} \times \mathbf{C} \times \mathbf{H} \times \mathbf{W} \\ \text{Normalize} & \quad \downarrow \quad \downarrow \quad \downarrow \\ \boldsymbol{\mu}, \boldsymbol{\sigma} &: \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1} \\ \boldsymbol{\gamma}, \boldsymbol{\beta} &: \mathbf{1} \times \mathbf{C} \times \mathbf{1} \times \mathbf{1} \\ \mathbf{y} &= \boldsymbol{\gamma}(\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \boldsymbol{\beta} \end{aligned}$$

Layer Normalization (LN)



训练和测试时操作一致
针对序列数据，常用在递归网络中



Instance Normalization (IN)

➤ 常用于处理图像数据，作用于每个样本的像素或通道等维度

卷积网络的BN

$$\begin{array}{l} \mathbf{x}: N \times C \times H \times W \\ \text{Normalize} \quad \downarrow \quad \downarrow \quad \downarrow \\ \boldsymbol{\mu}, \boldsymbol{\sigma}: 1 \times C \times 1 \times 1 \\ \gamma, \beta: 1 \times C \times 1 \times 1 \\ \mathbf{y} = \gamma(\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \beta \end{array}$$

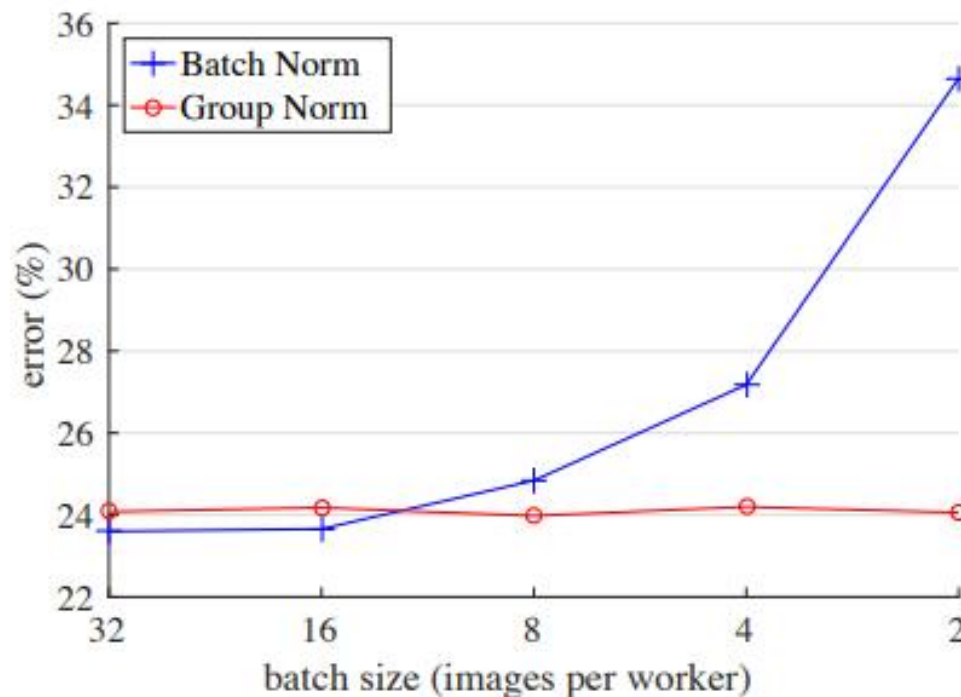
卷积网络的IN

$$\begin{array}{l} \mathbf{x}: N \times C \times H \times W \\ \text{Normalize} \quad \quad \quad \downarrow \quad \downarrow \\ \boldsymbol{\mu}, \boldsymbol{\sigma}: N \times C \times 1 \times 1 \\ \gamma, \beta: 1 \times C \times 1 \times 1 \\ \mathbf{y} = \gamma(\mathbf{x} - \boldsymbol{\mu}) / \boldsymbol{\sigma} + \beta \end{array}$$

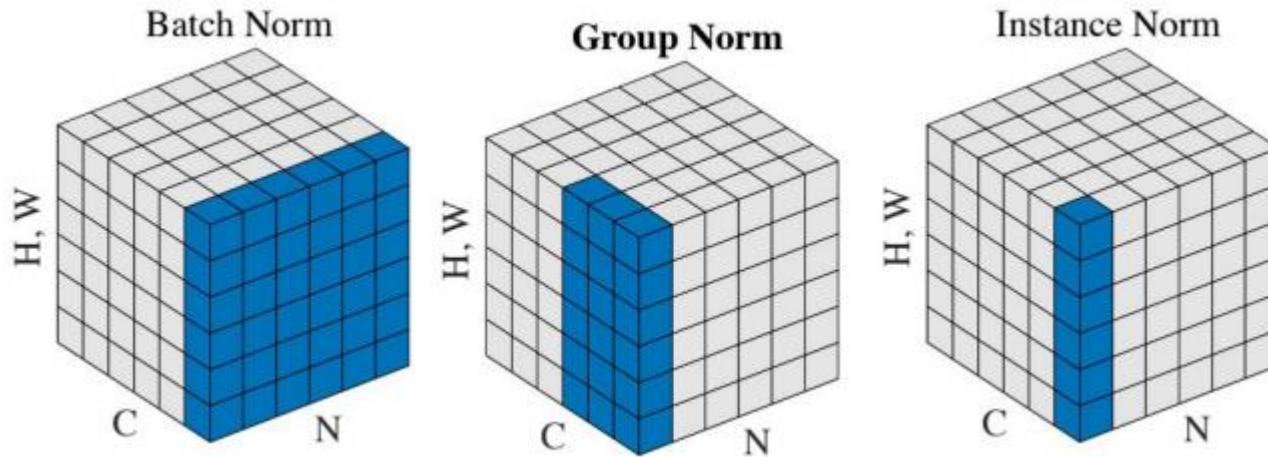


Group Normalization (GN)

- 针对BN在batch size较小时错误率较高而提出的改进算法
- GN将通道分成几组，并在每组内计算归一化的均值和方差，其准确性在很宽的批量大小范围内都很稳定
- GN归一化操作的计算不依赖batch size的大小



归一化层的比较



- 当GN将通道分成1组，此时即为LN
- 当GN将通道分成C组，此时即为IN

训练目标和方式

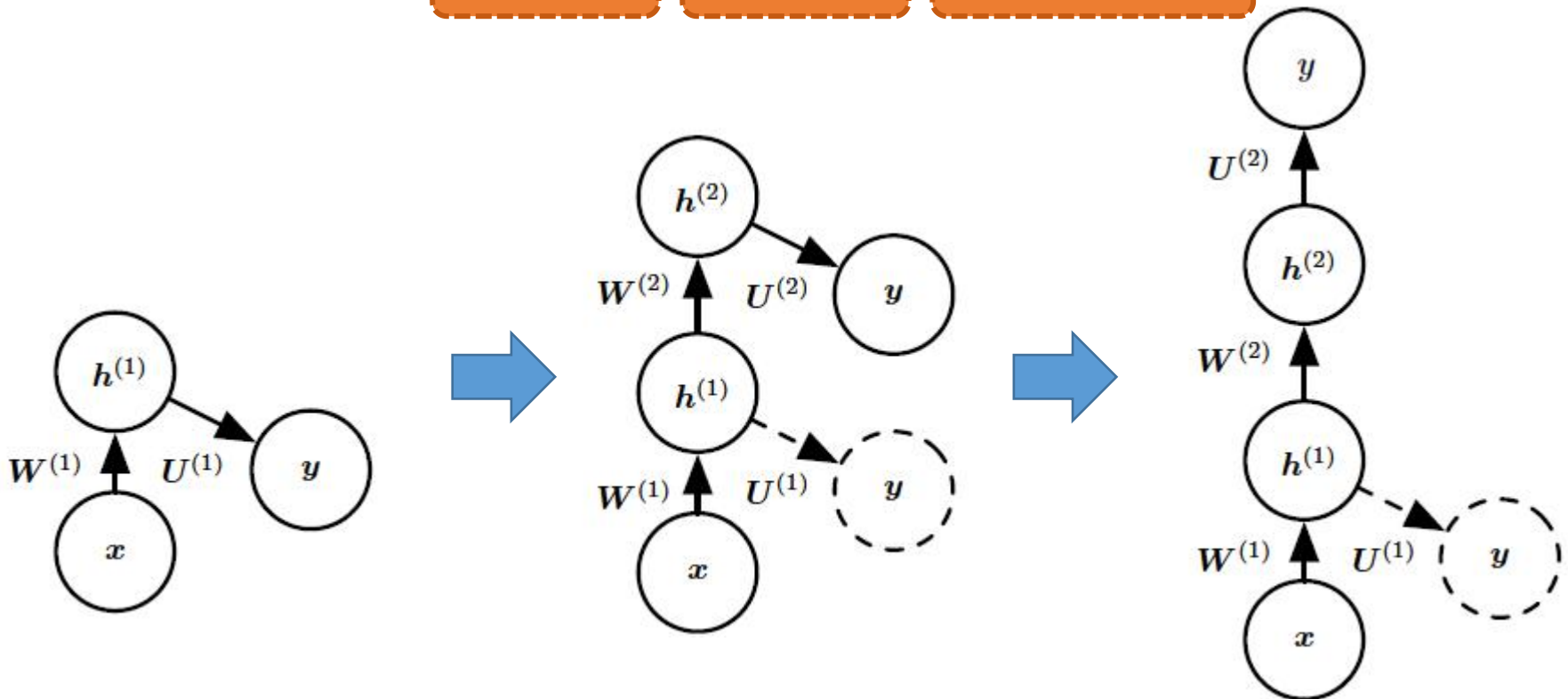
监督预训练

- 不同任务的模型往往都是从零开始来训练的，一切知识都从需要与真实数据分布保持一致的训练数据中得到。
- 但实际应用中：

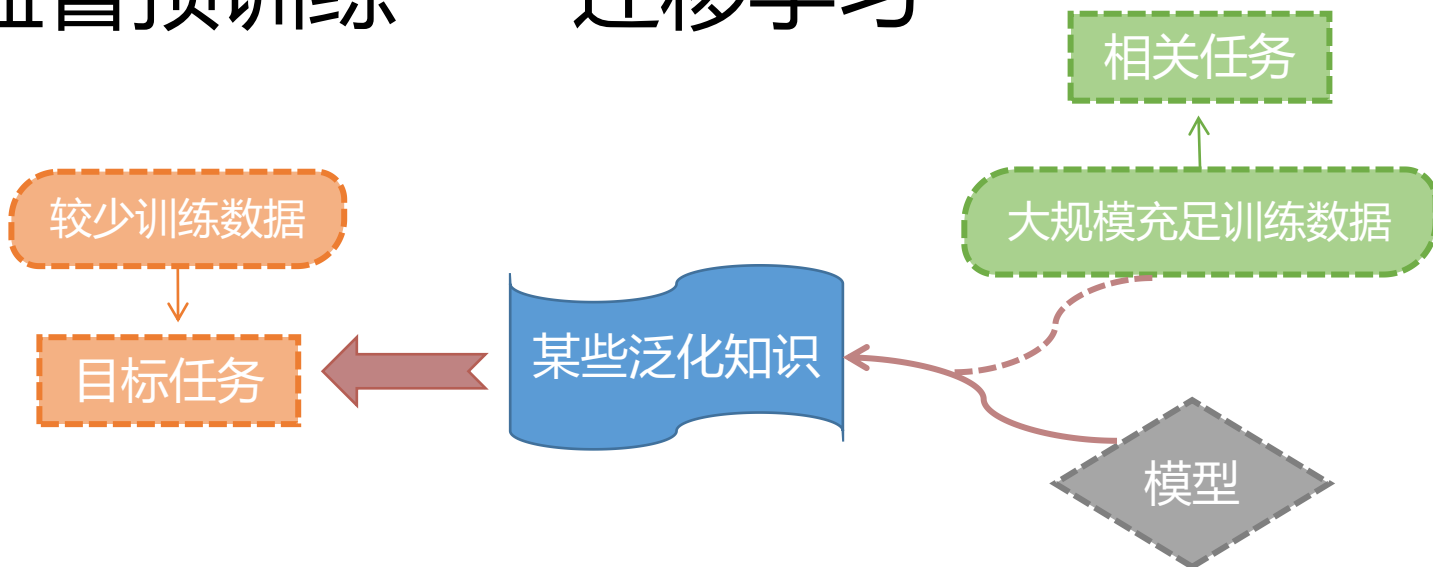
模型复杂

任务困难

训练数据过少



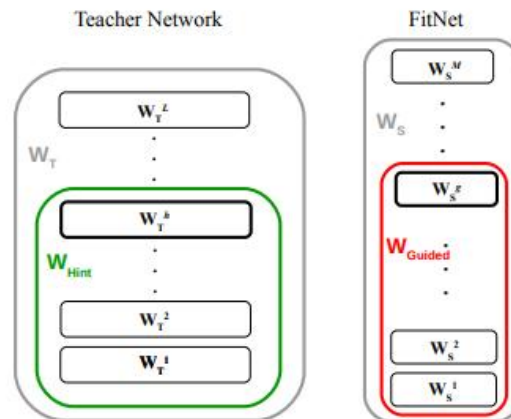
监督预训练——迁移学习



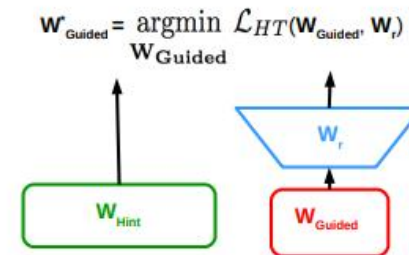
- 在一组任务上预训练了8层权重的深度卷积网络
 - 1000个ImageNet对象类的子集
- 用该网络的前 k 层初始化同样规模的网络，其余层随机初始化
 - 联合训练以执行不同的任务（1000个ImageNet对象类的另一个子集）
 - 训练样本少于第一个任务

监督预训练——FitNets

- 教师网络：深度浅（5层）且足够宽的网络，较容易训练
- 学生网络：深度深（11-19层）且窄的网络，很难用SGD训练
- 学生网络预测原任务的输出而且预测教师网络中间层的值
 - 学生网络中间层去回归教师网络的中间层



(a) Teacher and Student Networks

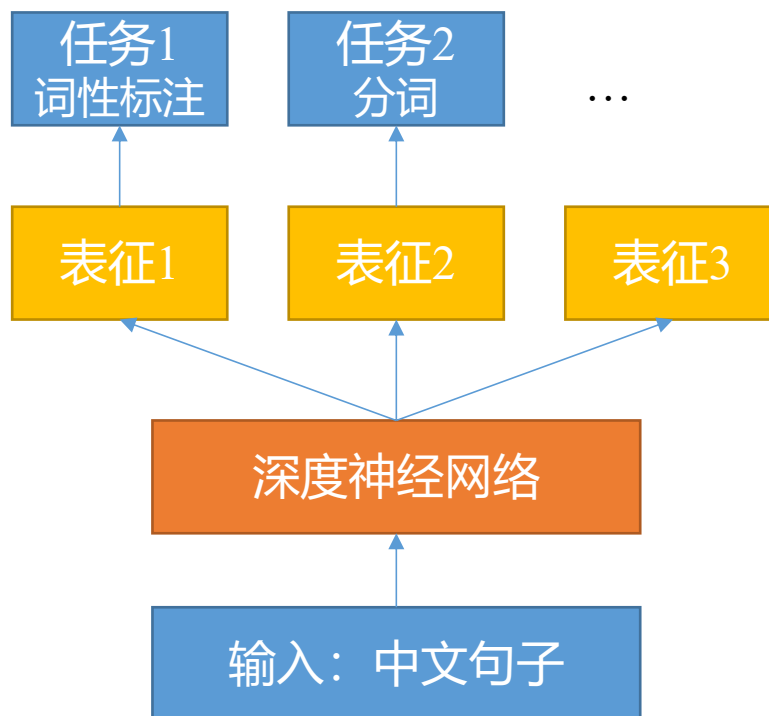


(b) Hints Training

$$\mathcal{L}_{HT}(\mathbf{W}_{\text{Guided}}, \mathbf{W}_{\text{r}}) = \frac{1}{2} \|u_h(\mathbf{x}; \mathbf{W}_{\text{Hint}}) - r(v_g(\mathbf{x}; \mathbf{W}_{\text{Guided}}); \mathbf{W}_{\text{r}})\|^2$$

多任务学习

➤通过合并几个任务中的样例来提高泛化的一种方式





课程学习

- 基于规划学习过程的想法，首先学习简单的概念，然后逐步学习依赖于这些简化概念的复杂概念
 - 在大规模的神经语言模型任务上使用课程学习，可以获得更好的结果
- 课程学习被证实为与人类教学方式一致
 - 教师刚开始会展示更容易、更典型的示例，然后帮助学习者在不太显然的情况下提炼决策面
 - 在人类教学上，基于课程学习的策略比基于样本均匀采样的策略更有效

课程学习

- 首先学习简单的、普适性的知识结构，然后逐渐增加难度，过渡到学习更复杂、更专业化的知识。



Samples of "Dog" to learn earlier.



Samples of "Dog" to learn later.

课程学习

- 首先学习简单的、普适性的知识结构，然后逐渐增加难度，过渡到学习更复杂、更专业化的知识。



“bus” samples to learn earlier



“bus” samples to learn later



Age



数据增强

数据集增强——水平翻转

- 需要保证水平翻转不会改变类别
 - 反例：OCR中的字符b和d，6和9



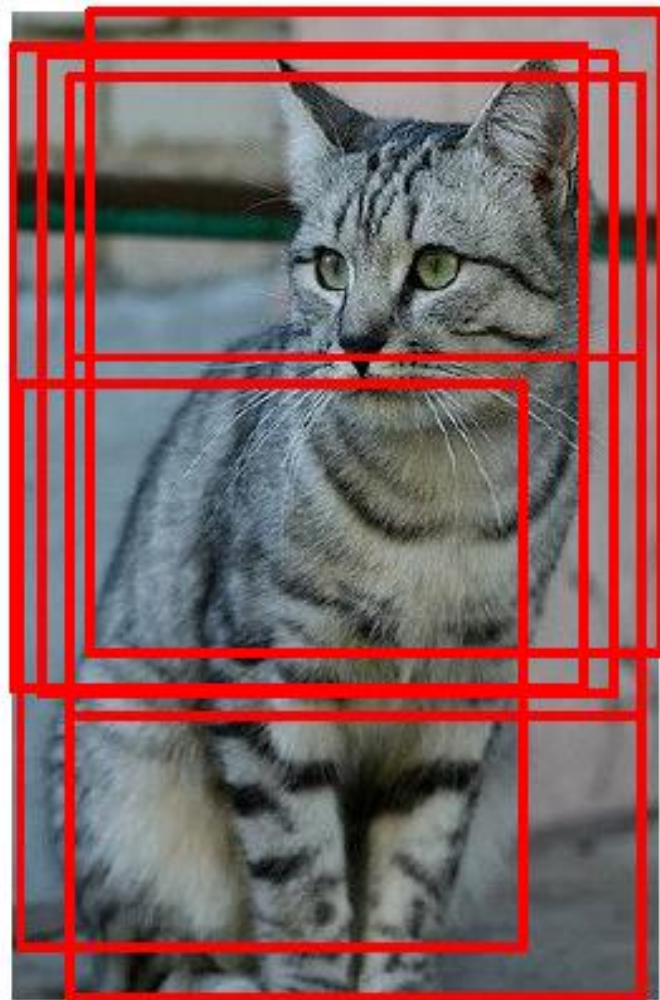
数据集增强——随机裁剪和缩放

➤训练

- 从 $[256, 480]$ 的范围内随机取 L
- 按比例缩放，其中短边长度为 L
- 从缩放的图上采样 224×224 的图

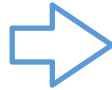
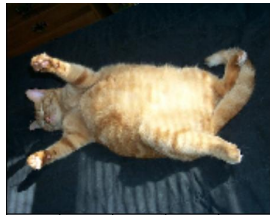
➤测试

- 五个尺度上的缩放 $\{224, 256, 384, 480, 640\}$
- 每个大小的图片，用10个 224×224 的裁剪
 - 2 (水平翻转) \times (4 角 + 1 中心)



数据集增强——变色

➤调整色调，饱和度和亮度（例如 $[0.5, 1.5]$ ）



亮度

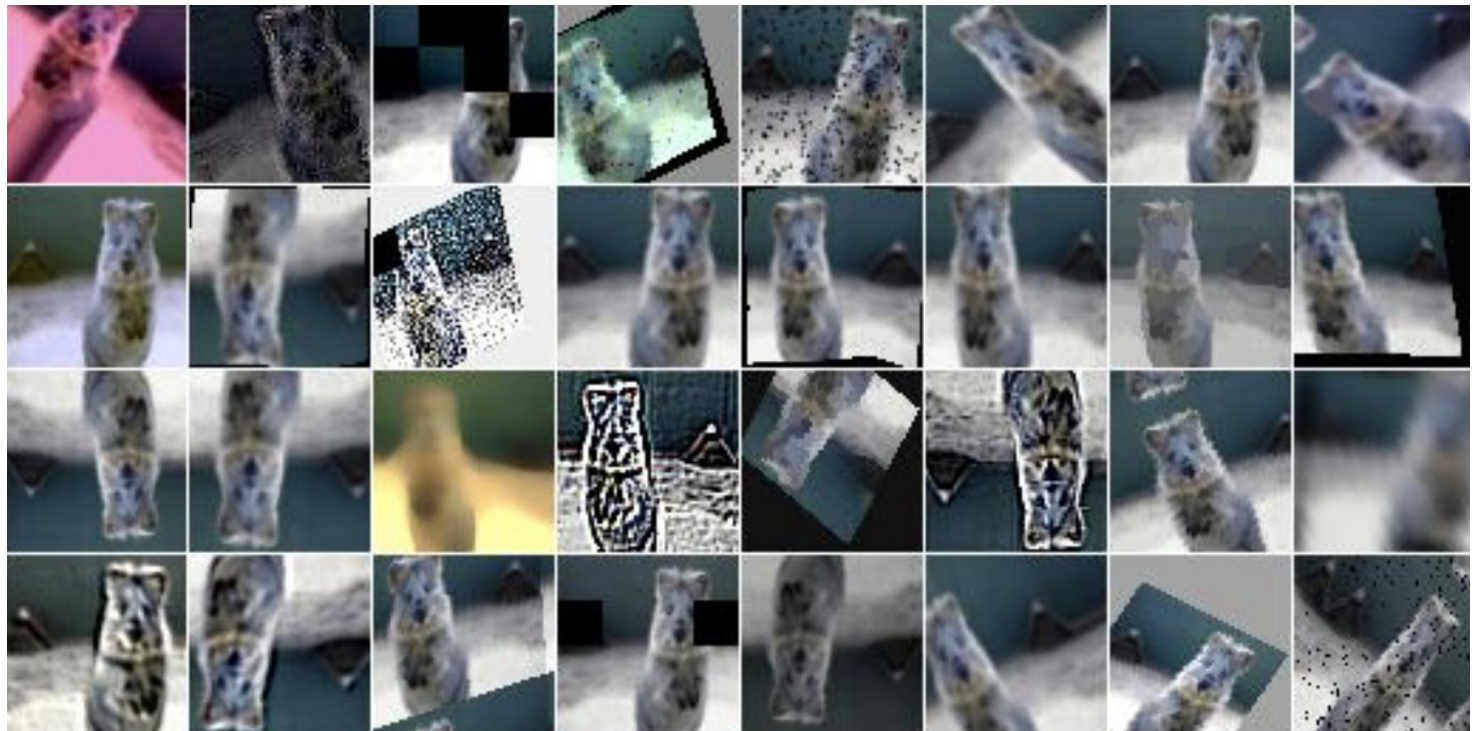


色调



数据集增强——更多方法

➤ 平移、旋转、伸缩、裁剪、扭曲的任意组合



数据集增强——注入噪声

➤ 向神经网络的输入层注入噪声也是一种数据集增强的方法

原图



加入 $\mathcal{N}(0, 0.1^2)$ 的噪声



也可以向隐藏单元施加噪声



数据集增强——标签平滑

- 在输出标签中添加噪声来避免模型过拟合
- 大多数数据集的 y 标签都有一定错误，错误的 y 不利于最大化 $\log p(y|x)$
- 一个样本 x 的标签一般用onehot向量表示

$$y = [0, \dots, 0, 1, 0, \dots, 0]$$




- 引入一个噪声对标签进行平滑，即假设样本以 ϵ 的概率为其它类。
平滑后的标签为

$$y = \left[\frac{\epsilon}{k-1}, \dots, \frac{\epsilon}{k-1}, 1 - \epsilon, \frac{\epsilon}{k-1}, \dots, \frac{\epsilon}{k-1} \right]$$

软目标
(Soft Target)

对抗学习

- 对抗样本: x 与 x' 非常相似, 人类难以察觉差异, 但神经网络会作出不同的预测

	$+ .007 \times$		$=$	
x		$\text{sign}(\nabla_x J(\theta, x, y))$		$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$
$y =$ 熊猫		线虫		长臂猿
w/ 57.7%		w/ 8.2%		w/ 99.3 %
confidence		confidence		confidence

对抗学习

在对抗扰动的训练集样本上训练网络, 以减少原有独立同分布的测试集的错误率

正则化



正则化

- 正则化 (Regularization) 是一类通过**限制模型复杂度**，从而**避免过拟合**，**提高泛化能力**的方法，比如引入约束、增加先验、提前停止等
- L1&L2正则化：通过约束参数的**L1 和L2范数**来减小模型在训练数据集上的过拟合现象

$$\mathcal{L}'(\boldsymbol{\theta}) = \underbrace{\mathcal{L}(\boldsymbol{\theta})}_{\text{损失函数}} + \frac{1}{2} \lambda \underbrace{\|\boldsymbol{\theta}\|_2^2}_{\text{正则化项}}$$

损失函数，比如平方损失，交叉熵损失等等

模型参数： $\boldsymbol{\theta} = \{w_1, w_2, \dots\}$

L1 正则:

$$\|\boldsymbol{\theta}\|_1 = |w_1| + |w_2| + \dots$$

L2 正则:

$$\|\boldsymbol{\theta}\|_2^2 = (w_1)^2 + (w_2)^2 + \dots$$

一般不考虑bias

L1正则化

➤ L1正则化:

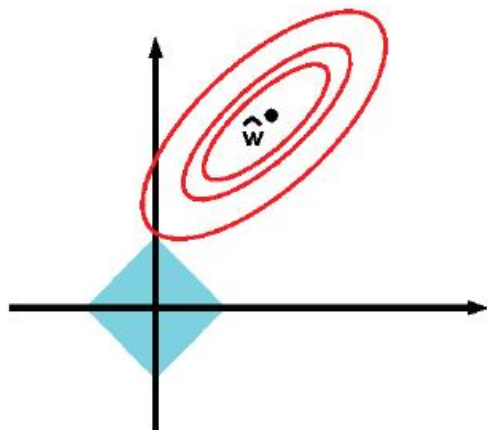
$$\mathcal{L}'(\boldsymbol{\theta}) = \mathcal{L}(\boldsymbol{\theta}) + \frac{1}{2}\lambda\|\boldsymbol{\theta}\|_1 \quad \|\boldsymbol{\theta}\|_1 = |w_1| + |w_2| + \dots$$

➤ 梯度: $\nabla\mathcal{L}' = \nabla\mathcal{L} + \lambda \text{sign}(\boldsymbol{\theta})$

➤ 梯度下降更新:

$$\boldsymbol{\theta}^t \leftarrow \boldsymbol{\theta}^{t-1} - \eta\nabla\mathcal{L}' = \boldsymbol{\theta}^{t-1} - \eta\nabla\mathcal{L} - \eta\lambda \text{sign}(\boldsymbol{\theta}^{t-1})$$

$$= \underbrace{\boldsymbol{\theta}^{t-1} - \eta\lambda \text{sign}(\boldsymbol{\theta}^{t-1})}_{\text{L1 regularization effect}} - \eta\nabla\mathcal{L}$$



- 如果 $w^{t-1} > 0$, w^{t-1} 变小
- 如果 $w^{t-1} < 0$, w^{t-1} 变大

➡ $w^{t-1} \rightarrow 0$

具有稀疏的特性

L2正则化

➤ L2正则化:

$$\mathcal{L}'(\boldsymbol{\theta}) = \mathcal{L}(\boldsymbol{\theta}) + \frac{1}{2} \lambda \|\boldsymbol{\theta}\|_2^2 \quad \|\boldsymbol{\theta}\|_2^2 = (w_1)^2 + (w_2)^2 + \dots$$

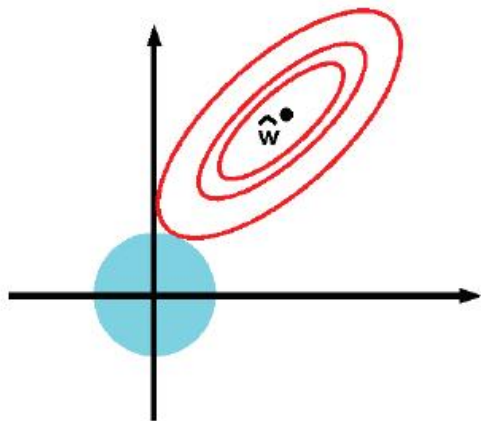
➤ 梯度: $\nabla \mathcal{L}' = \nabla \mathcal{L} + \lambda \boldsymbol{\theta}$

➤ 梯度下降更新:

$$\boldsymbol{\theta}^t \leftarrow \boldsymbol{\theta}^{t-1} - \eta \nabla \mathcal{L}' = \boldsymbol{\theta}^{t-1} - \eta \nabla \mathcal{L} - \eta \lambda \boldsymbol{\theta}^{t-1}$$

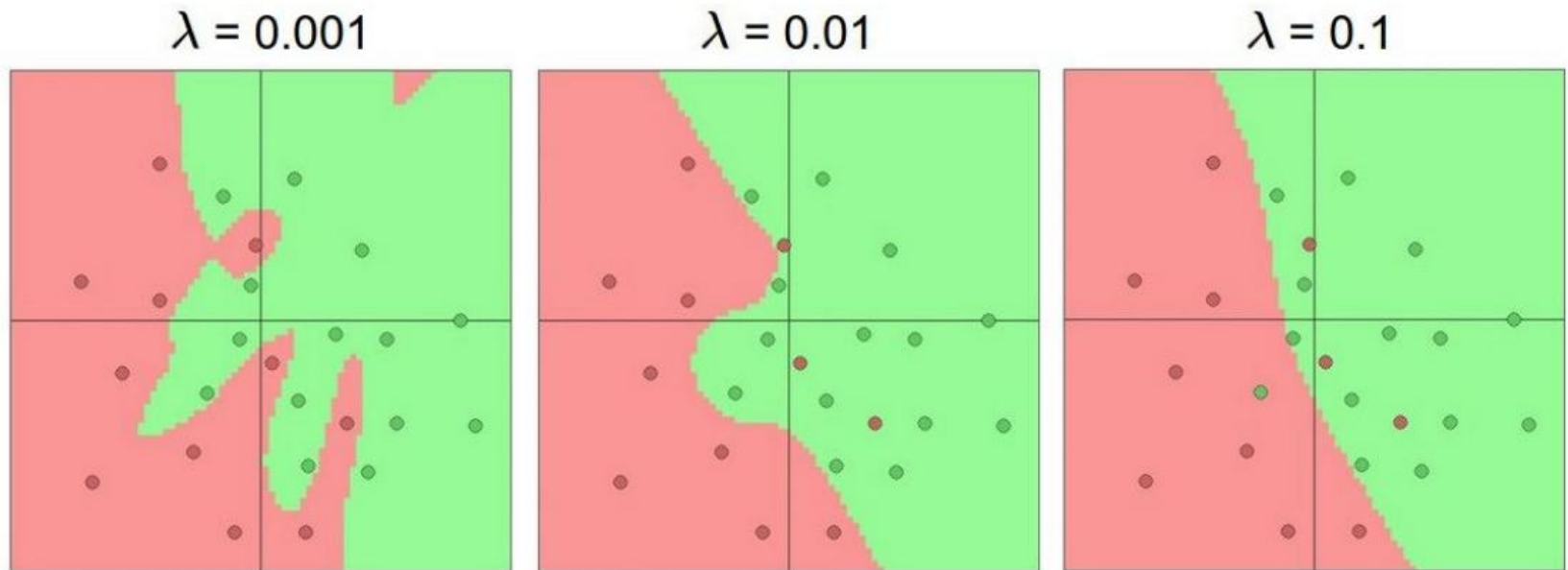
$$= \underbrace{(1 - \eta \lambda) \boldsymbol{\theta}^{t-1}}_{\downarrow} - \eta \nabla \mathcal{L} \quad \eta < 1, \lambda < 1$$

越来越小, 但由于 $-\eta \nabla \mathcal{L}$ 项, 使得参数不会变为0



权重衰减 (weight decay) 可用L2正则化来实现

神经网络示例



<http://playground.tensorflow.org/>



添加噪声

- 假设网络权重添加随机扰动 $\epsilon_w \sim \mathcal{N}(0, \eta I)$, 扰动模型为: $\hat{y}_{\epsilon_w}(x)$
- 将平方损失作为目标函数

$$\tilde{J}_W = \mathbb{E}_{p(x,y, \epsilon_w)} [(\hat{y}_{\epsilon_w}(x) - y)^2]$$

$$= \mathbb{E}_{p(x,y, \epsilon_w)} [\hat{y}_{\epsilon_w}^2(x) - 2y\hat{y}_{\epsilon_w}(x) + y^2]$$

$$\approx \mathbb{E}_{p(x,y, \epsilon_w)} [y_w^2(x) + 2y_w(x)\nabla y_w(x)^\top \epsilon_w + \epsilon_w^\top \nabla y_w(x)\nabla y_w(x)^\top \epsilon_w]$$

$$+ \mathbb{E}_{p(x,y, \epsilon_w)} [-2yy_w(x) - 2y\nabla y_w(x)^\top \epsilon_w + y^2]$$

$$= \mathbb{E}_{p(x,y)} [(y_w(x) - y)^2] + \mathbb{E}_{p(x,y)} [\text{tr}(\nabla y_w(x)\nabla y_w(x)^\top \mathbb{E}_{p(\epsilon_w)}[\epsilon_w \epsilon_w^\top])]$$

$$= J_W + \eta \mathbb{E}_{p(x,y)} [\|\nabla y_w(x)\|^2]$$

一阶泰勒展开近似

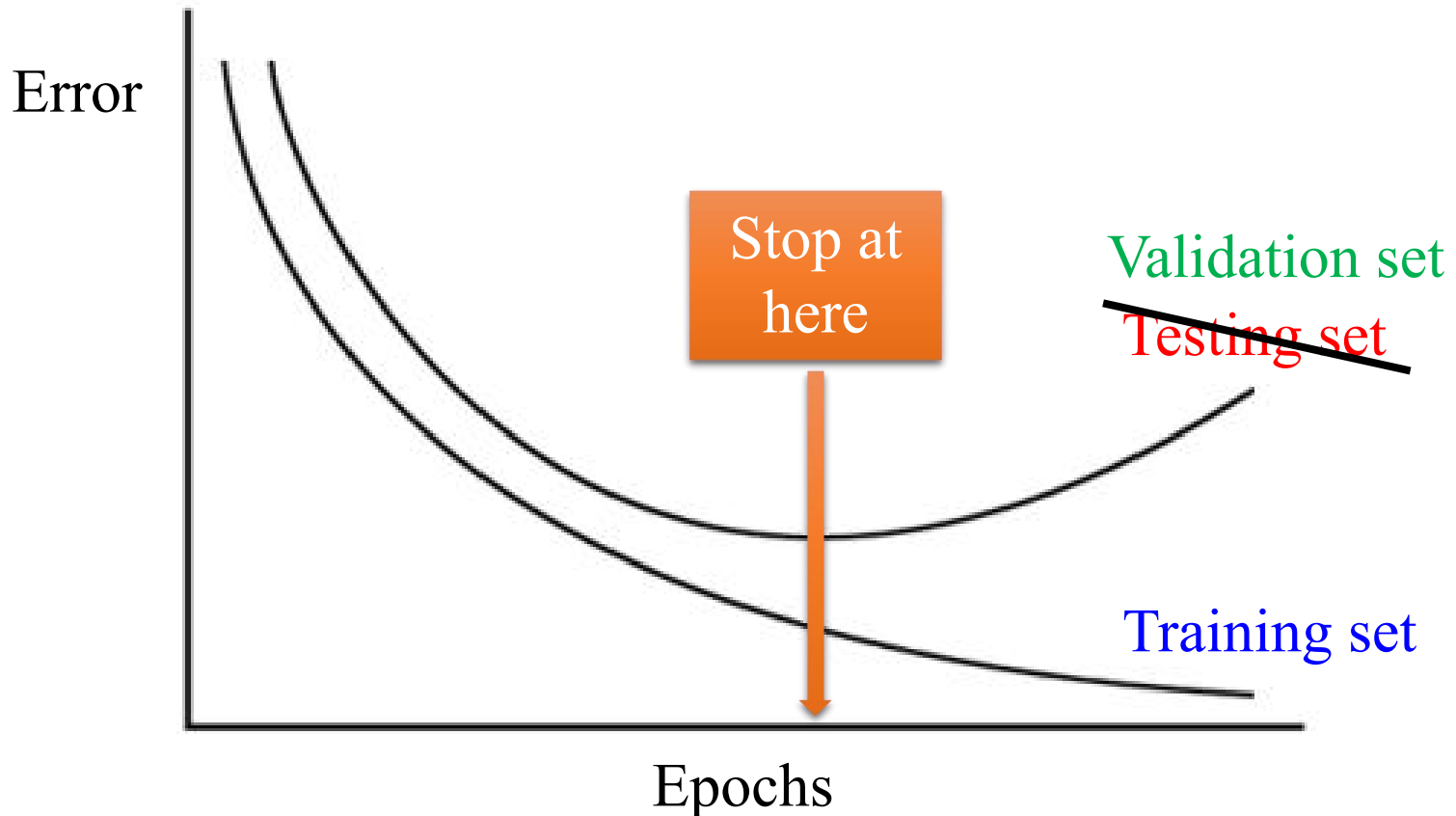
$$\hat{y}_{\epsilon_w}(x) \approx y_w(x) + \nabla_w y_w(x)^\top \epsilon_w$$

$$\mathbb{E}_{p(\epsilon_w)}[\epsilon_w \epsilon_w^\top] = \eta I$$

鼓励参数进入 权重小扰动对输出相对影响较小的参数空间区域

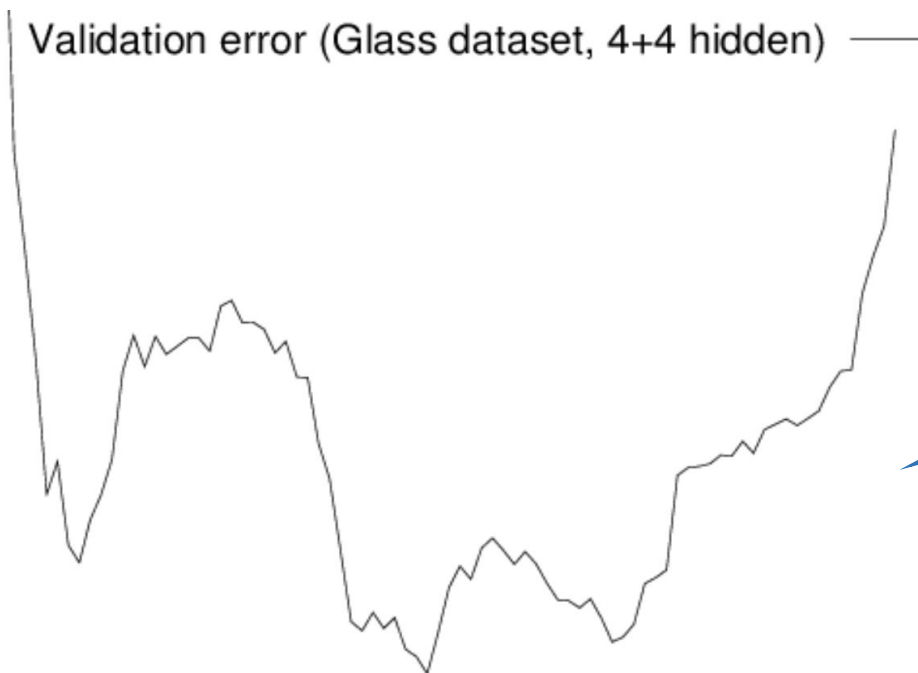
提前终止

- 验证集上**准确率下降**（**损失上升**）的时候停止训练
- 训练很长时间，保存在**验证集**上最优的模型



提前终止

- 早停标准：什么时候停？寻找更好的训练时间和泛化错误之间的权衡。



模型在验证集上的表现可能短暂的变差之后继续变好

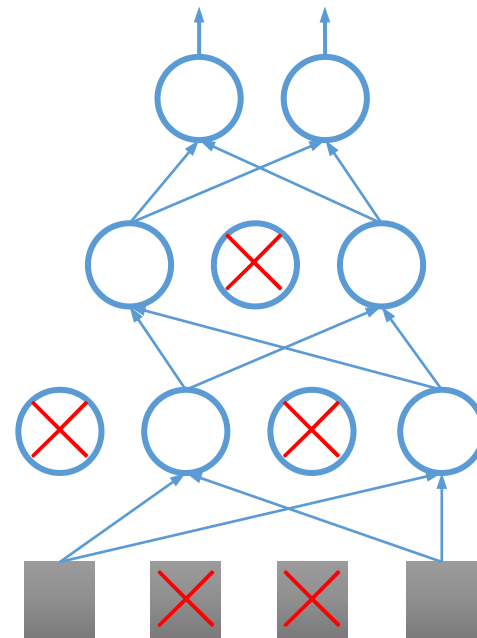
真实的验证集误差变化曲线

https://blog.csdn.net/weixin_40519315

Dropout

➤在训练时，取得一个batch后，以概率 p 设置一些神经元为0

- 网络结构发生了变化
- 用新网络在batch上计算梯度
- 在新网络上进行参数更新
- 每次取新batch时，都需要重新随机对神经元置为0





Dropout

- 测试时，**没有dropout**，即所有的神经元都处于激活状态
- **缩放激活函数的输出**，使得每个神经元测试时输出等于训练时的期望输出

权重比例推断规则

如果训练时dropout 的概率为 p ，那么测试时所有的权重要**乘以 $1 - p$**



Dropout—训练和测试时实现

```
p = 0.5 # probability of keeping a unit active. higher = less dropout
```

该概率为保留概率

```
def train_step(X):
```

```
    """ X contains the data """
```

```
    # forward pass for example 3-layer neural network
```

```
    H1 = np.maximum(0, np.dot(W1, X) + b1)
```

```
    U1 = np.random.rand(*H1.shape) < p # first dropout mask
```

```
    H1 *= U1 # drop!
```

```
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
```

```
    U2 = np.random.rand(*H2.shape) < p # second dropout mask
```

```
    H2 *= U2 # drop!
```

```
    out = np.dot(W3, H2) + b3
```

```
def predict(X):
```

```
    # ensembled forward pass
```

```
    H1 = np.maximum(0, np.dot(W1, X) + b1) * p # NOTE: scale the activations
```

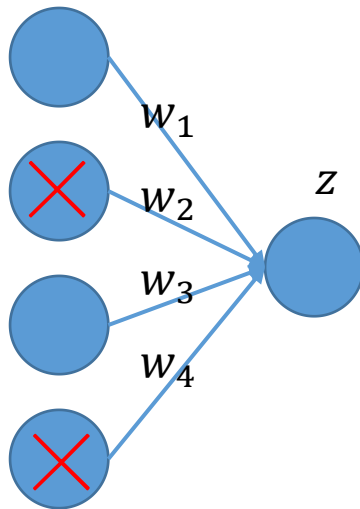
```
    H2 = np.maximum(0, np.dot(W2, H1) + b2) * p # NOTE: scale the activations
```

```
    out = np.dot(W3, H2) + b3
```

Dropout—直观解释

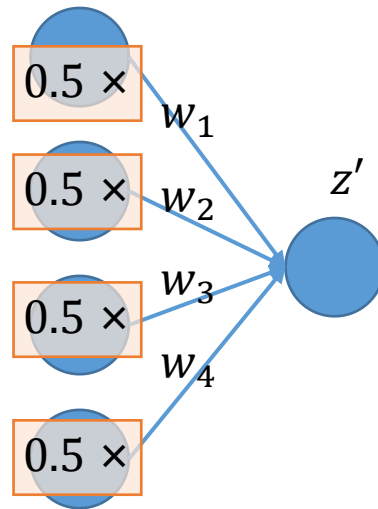
训练时dropout

dropout 概率为0.5



测试时dropout

No dropout



× 直接用训练后的权重

→ $z' \approx 2z$

✓ 权重乘以 $1 - p$

→ $z' \approx z$

Dropout动机

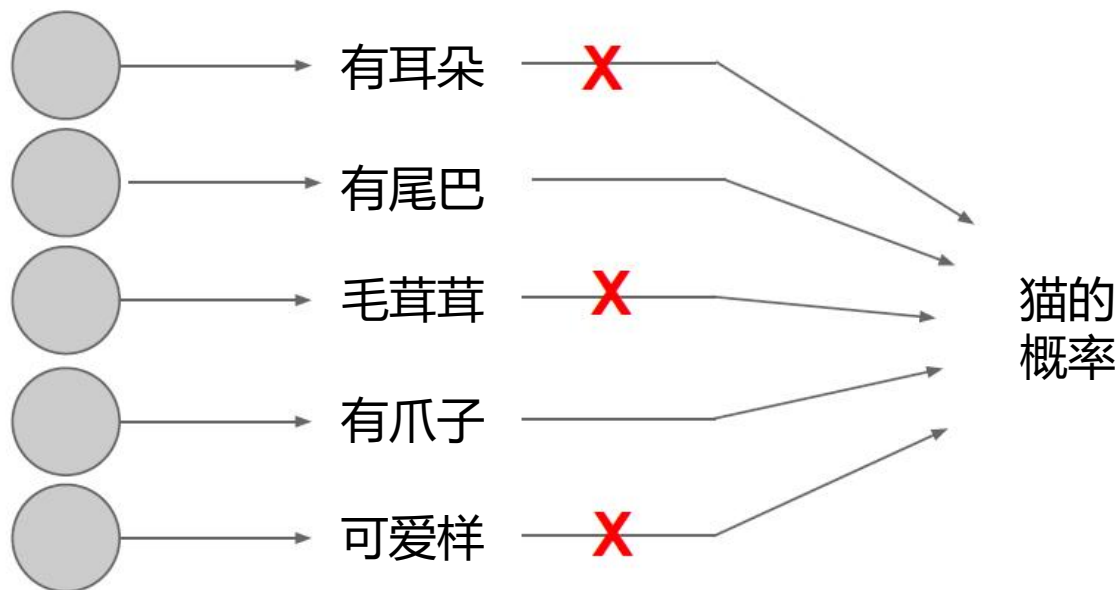
1

强制网络有冗余特征表示

2

防止特征的co-adaptation

co-adaptation: feature detectors只有在一些其它特定的feature detectors存在时才能发挥作用的情况





Dropout动机

➤ 人脸识别情形下的动机:

1

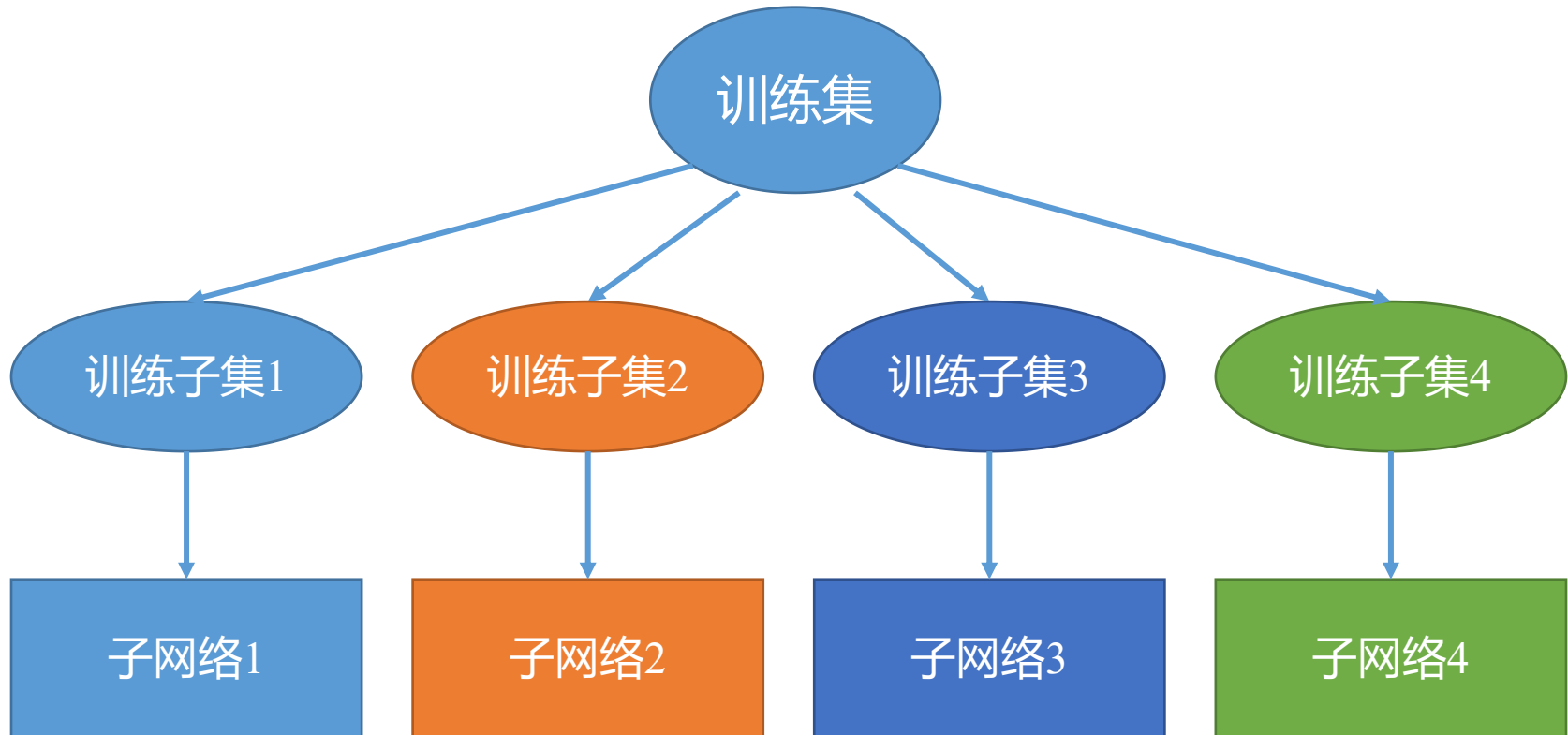
模型学得通过鼻检测脸的隐藏单元 h_i , 丢失 h_i 对应于擦除图像中有鼻子的信息

2

模型必须学习另一种 h_i , 要么是鼻子存在的冗余编码, 要么是像嘴这样的脸部的另一特征

Dropout 是一种模型集成 (Bagging)

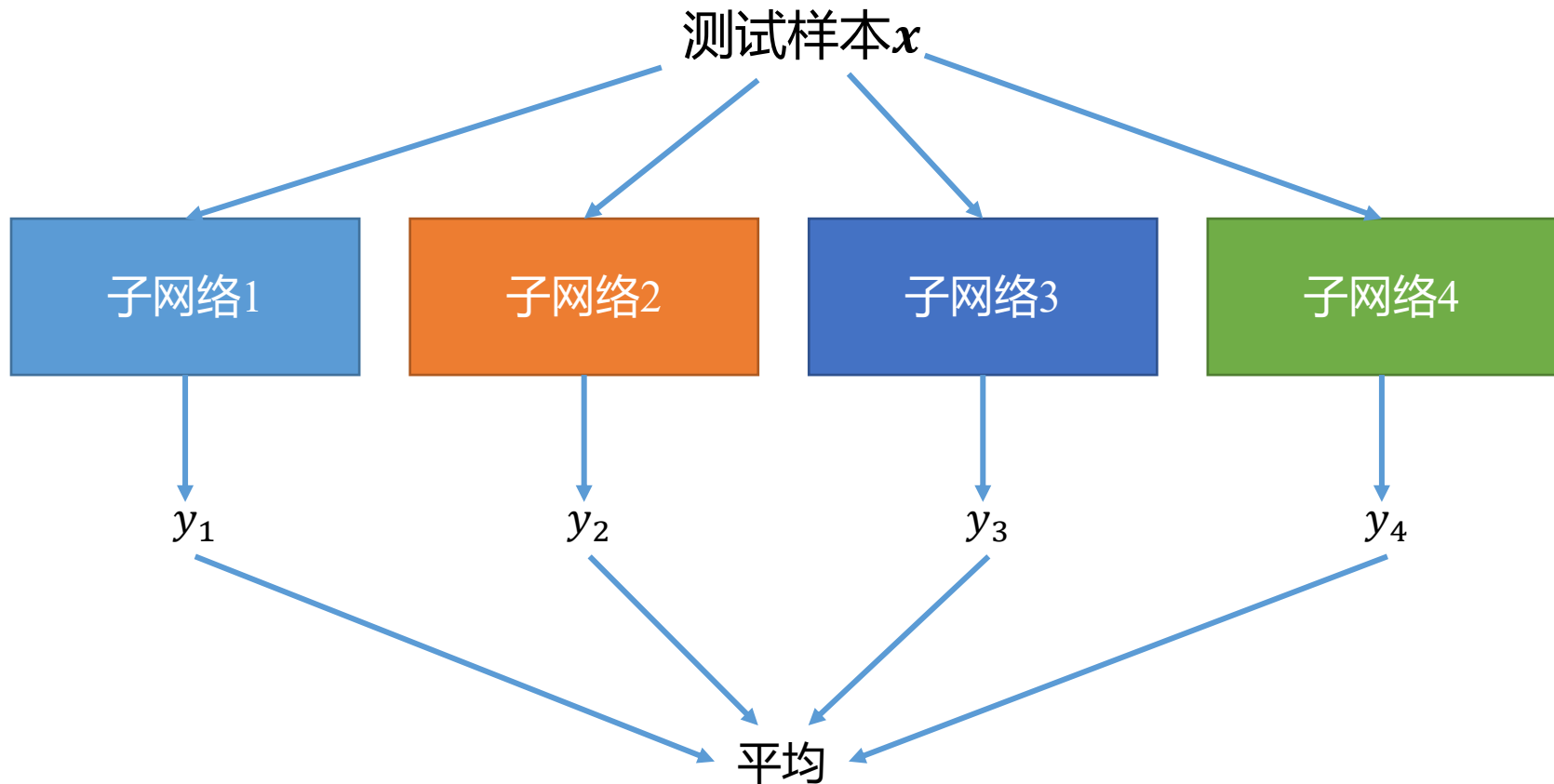
➤ Bagging:



训练一系列不同结构的神经网络

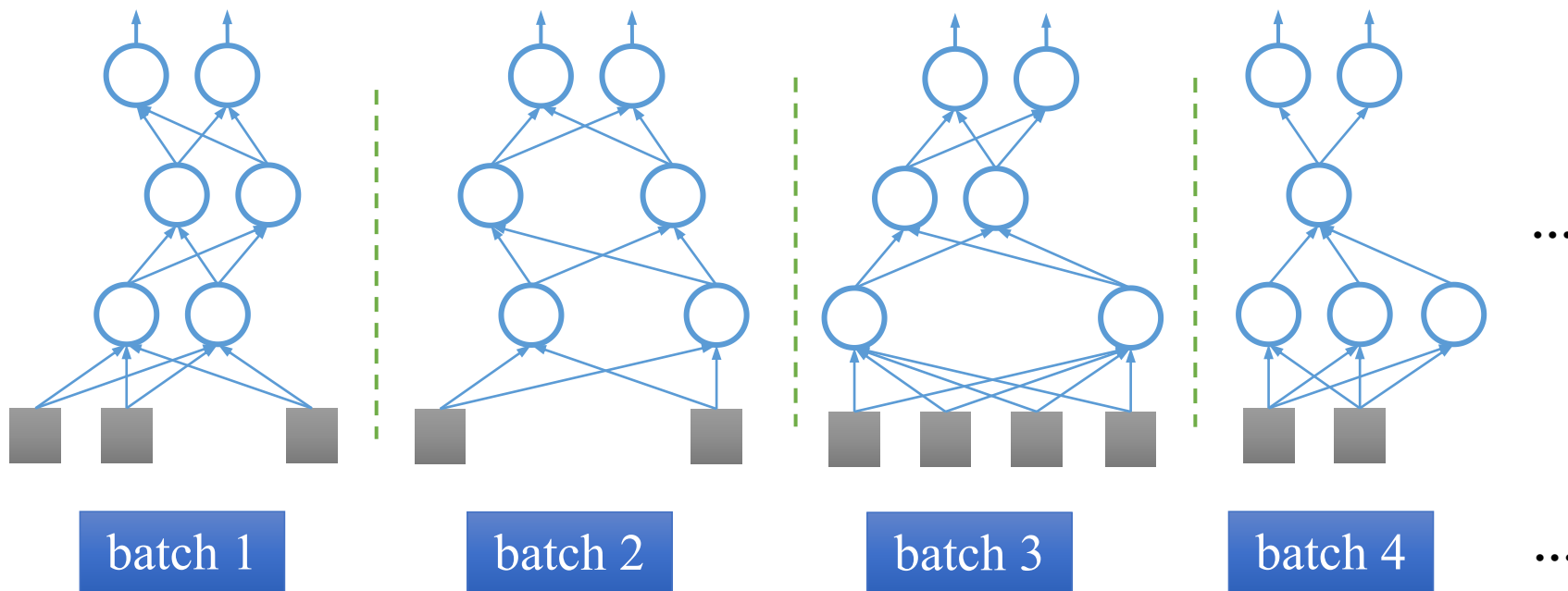
Dropout 是一种模型集成 (Bagging)

➤ Bagging:



Dropout 是一种模型集成 (Bagging)

- **训练时**：每做一次Dropout，相当于从原始的网络中**采样**得到一个子网络



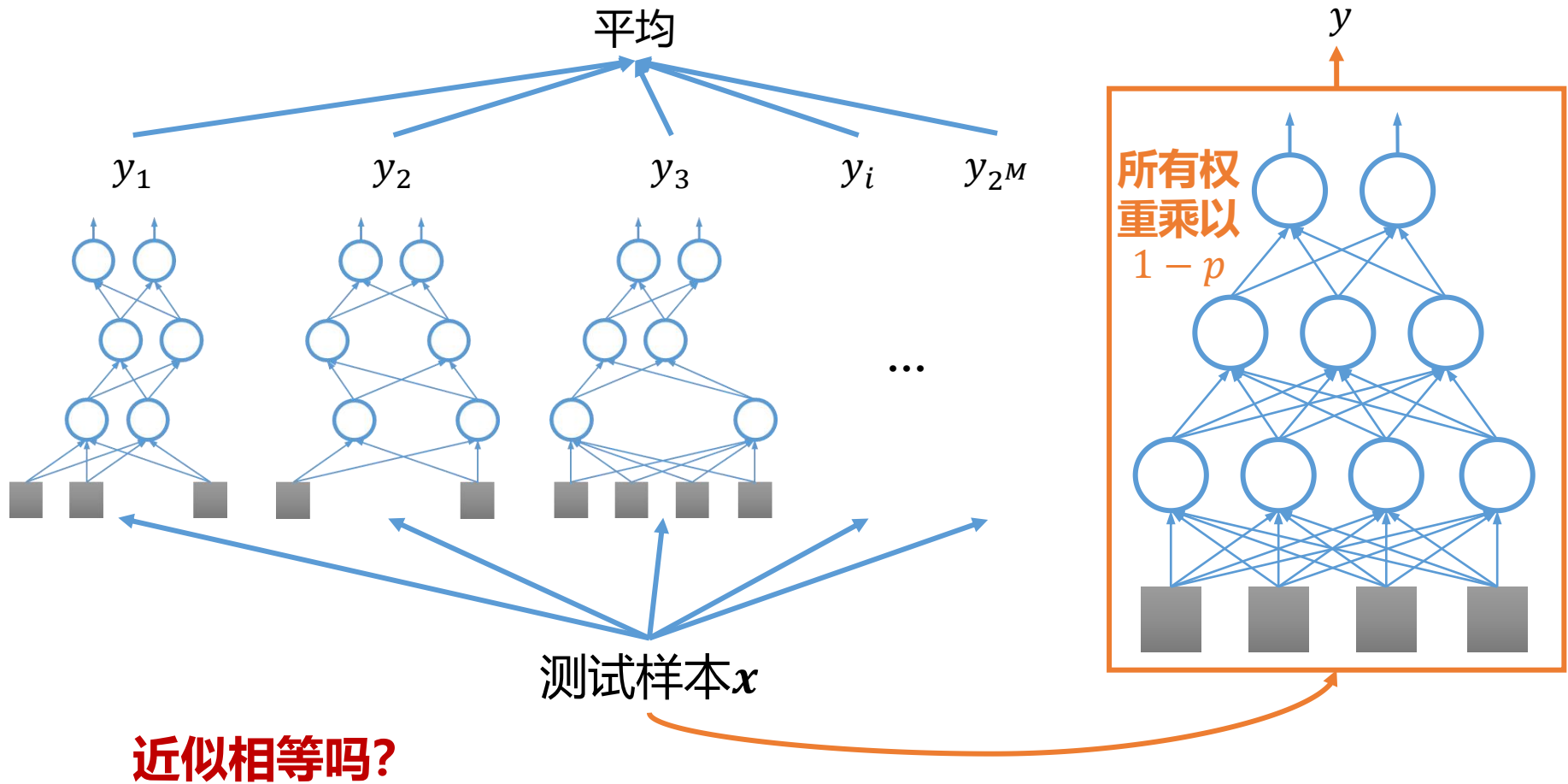
- 若神经网络有 M 个神经元，则共有 2^M 个不同结构的网络
- 一个batch数据训练一个结构的网络，只**更新一次参数**
- 不同结构的网络**共享部分参数**



Dropout 是一种模型集成 (Bagging)

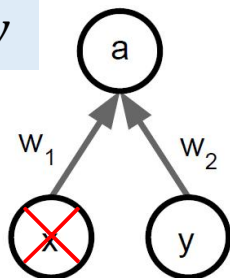
Dropout	Bagging
所有模型共享参数	所有模型相互独立
每一步训练一小部分的网络	每个模型在相应数据集上训练到收敛

Dropout 是一种模型集成 (Bagging)

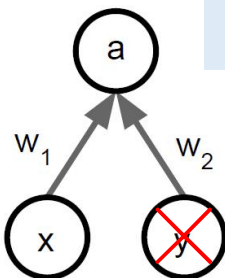


Dropout 是一种模型集成 (Bagging)

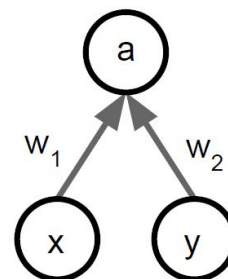
$$a = 0x + w_2y$$



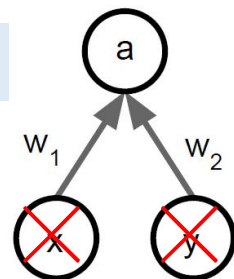
$$a = w_1x + 0y$$



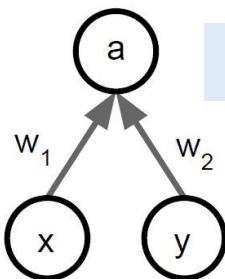
$$a = w_1x + w_2y$$



$$a = 0x + 0y$$



$$a = w_1x + w_2y$$



训练时:
$$\begin{aligned}\mathbb{E}[a] &= \frac{1}{4}(w_1x + w_2y) + \frac{1}{4}(w_1x + 0y) + \frac{1}{4}(0x + 0y) + \frac{1}{4}(0x + w_2y) \\ &= \frac{1}{2}(w_1x + w_2y)\end{aligned}$$

测试时:
$$\mathbb{E}[a] = w_1x + w_2y$$
 所有权重都乘以 $1/2$



Dropout 是一种模型集成 (Bagging)

➤ 权重比例推断规则对许多无非线性隐藏单元(model族)是精确的

➤ 考虑 $P(y = y|\mathbf{v}, \boldsymbol{\mu}) = \text{softmax}(\mathbf{W}^T(\mathbf{v} \odot \boldsymbol{\mu}) + \mathbf{b})_y$

$$\tilde{P}_{\text{ensemble}}(y = y|\mathbf{v}) = \sqrt[2^n]{\prod_{\boldsymbol{\mu} \in \{0,1\}^n} P(y = y|\mathbf{v}; \boldsymbol{\mu})}$$

使用集成成员预测分布的几何平均而不是算术平均

$$= \sqrt[2^n]{\prod_{\boldsymbol{\mu} \in \{0,1\}^n} \frac{\exp(\mathbf{W}_{y,:}^T(\mathbf{v} \odot \boldsymbol{\mu}) + b_y)}{\sum_{y'} \exp(\mathbf{W}_{y',:}^T(\mathbf{v} \odot \boldsymbol{\mu}) + b_{y'})}}$$

$$\propto \sqrt[2^n]{\prod_{\boldsymbol{\mu} \in \{0,1\}^n} \exp(\mathbf{W}_{y,:}^T(\mathbf{v} \odot \boldsymbol{\mu}) + b_y)}$$

$$= \exp\left(\frac{1}{2^n} \sum_{\boldsymbol{\mu} \in \{0,1\}^n} (\mathbf{W}_{y,:}^T(\mathbf{v} \odot \boldsymbol{\mu}) + b_y)\right)$$

$$= \exp\left(\frac{1}{2} \mathbf{W}_{y,:}^T \mathbf{v} + b_y\right)$$



Dropout 是一种模型集成 (Bagging)

➤ 权重比例推断规则对许多无非线性隐藏单元(model族)是精确的

- softmax函数回归分类
- 条件正态输出的回归网络
- 隐层不含非线性的深度网络



Dropout是一种正则化

➤考虑基于指数分布族的广义线性模型，囊括线性回归、对率回归

$$P_{\boldsymbol{\beta}}(y|\mathbf{x}) = h(y) \exp \{y \mathbf{x}^{\top} \boldsymbol{\beta} - A(\mathbf{x}^{\top} \boldsymbol{\beta})\}$$

$$A(\mathbf{x}^{\top} \boldsymbol{\beta}) = \log \int h(y) \exp \{y \mathbf{x}^{\top} \boldsymbol{\beta}\} dy \quad \text{为对数配分函数，用于归一化}$$

➤指数分布族的重要性质

$$A'(\mathbf{x}^{\top} \boldsymbol{\beta}) = \mathbb{E}_{P_{\boldsymbol{\beta}}(y|\mathbf{x})}[y]$$

$$A''(\mathbf{x}^{\top} \boldsymbol{\beta}) = \mathbb{V}_{P_{\boldsymbol{\beta}}(y|\mathbf{x})}[y]$$

Dropout是一种正则化

➤对于样本 (\mathbf{x}, y) , 其损失定义为 $\ell_{\mathbf{x}, y}(\boldsymbol{\beta}) = -\log P_{\boldsymbol{\beta}}(y|\mathbf{x})$

➤考虑Dropout 噪声 $\tilde{\mathbf{x}} = \mathbf{x} \odot \boldsymbol{\xi}$

$$\mathbb{E}[\tilde{\mathbf{x}}] = \mathbf{x}$$

$$\mathbb{V}[\tilde{\mathbf{x}}] = \frac{\delta}{1-\delta} \text{diag}(\mathbf{x}\mathbf{x}^T) \quad \leftarrow \quad \mathbb{E}[\tilde{x}_j] = x_j \quad \mathbb{V}[\tilde{x}_j] = \frac{\delta}{1-\delta} x_j^2$$

ξ_j	0	$1/(1-\delta)$
概率p	δ	$1-\delta$

➤含噪声的损失函数定义为

$$\begin{aligned}
 \mathcal{L}(\boldsymbol{\beta}) &= \sum_{i=1}^n \mathbb{E}_{\boldsymbol{\xi}} [\ell_{\tilde{\mathbf{x}}_i, y_i}(\boldsymbol{\beta})] \\
 &= - \sum_{i=1}^n (y_i \mathbf{x}_i^T \boldsymbol{\beta} - \mathbb{E}_{\boldsymbol{\xi}} [A(\tilde{\mathbf{x}}_i^T \boldsymbol{\beta})]) \\
 &= \sum_{i=1}^n \left[- (y_i \mathbf{x}_i^T \boldsymbol{\beta} - A(\mathbf{x}_i^T \boldsymbol{\beta})) \right] + \left(\mathbb{E}_{\boldsymbol{\xi}} [A(\tilde{\mathbf{x}}_i^T \boldsymbol{\beta})] - A(\mathbf{x}_i^T \boldsymbol{\beta}) \right) \\
 &= \sum_{i=1}^n \ell_{\mathbf{x}_i, y_i}(\boldsymbol{\beta}) + R(\boldsymbol{\beta})
 \end{aligned}$$

正则项



Dropout是一种正则化

➤ 正则项
$$R(\boldsymbol{\beta}) = \sum_{i=1}^n \mathbb{E}_{\xi}[A(\tilde{\mathbf{x}}_i^{\top} \boldsymbol{\beta})] - A(\mathbf{x}_i^{\top} \boldsymbol{\beta})$$

➤ 考虑 $A(\tilde{\mathbf{x}}_i^{\top} \boldsymbol{\beta})$ 在 $\mathbf{x}_i^{\top} \boldsymbol{\beta}$ 附近的二阶泰勒展开

$$\begin{aligned} \mathbb{E}_{\xi}[A(\tilde{\mathbf{x}}_i^{\top} \boldsymbol{\beta})] &\approx A(\mathbf{x}_i^{\top} \boldsymbol{\beta}) + A'(\mathbf{x}_i^{\top} \boldsymbol{\beta}) \mathbb{E}_{\xi}[(\tilde{\mathbf{x}}_i^{\top} \boldsymbol{\beta} - \mathbf{x}_i^{\top} \boldsymbol{\beta})] + 1/2 A''(\mathbf{x}_i^{\top} \boldsymbol{\beta}) \mathbb{E}_{\xi}[(\tilde{\mathbf{x}}_i^{\top} \boldsymbol{\beta} - \mathbf{x}_i^{\top} \boldsymbol{\beta})^2] \\ &= A(\mathbf{x}_i^{\top} \boldsymbol{\beta}) + \frac{1}{2} A''(\mathbf{x}_i^{\top} \boldsymbol{\beta}) \boldsymbol{\beta}^{\top} \mathbb{V}_{\xi}(\tilde{\mathbf{x}}_i) \boldsymbol{\beta} \end{aligned}$$

➡
$$R(\boldsymbol{\beta}) = \frac{1}{2} \frac{\delta}{1 - \delta} \boldsymbol{\beta}^{\top} \sum_{i=1}^n A''(\mathbf{x}_i^{\top} \boldsymbol{\beta}) \text{diag}(\mathbf{x}_i \mathbf{x}_i^{\top}) \boldsymbol{\beta}$$

$$= \frac{1}{2} \frac{\delta}{1 - \delta} \boldsymbol{\beta}^{\top} \text{diag}(\mathbf{X}^{\top} V(\boldsymbol{\beta}) \mathbf{X}) \boldsymbol{\beta}$$

记 $V(\boldsymbol{\beta})$ 为 $n \times n$ 对角阵,
对角线上元素为 $A''(\mathbf{x}_i^{\top} \boldsymbol{\beta})$



Dropout是一种正则化

➤对率回归, 记 $P_{\beta}(y = 1|x_i) = p_i$,

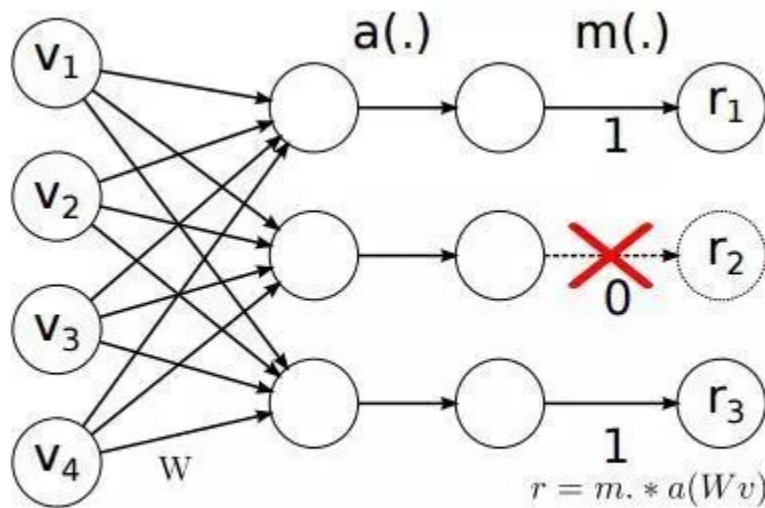
$$\text{则 } A''(\mathbf{x}_i^{\top} \boldsymbol{\beta}) = \mathbb{V}_{P_{\beta}(y|x_i)}[y] = p_i(1 - p_i)$$

$$\begin{aligned} R(\boldsymbol{\beta}) &= \frac{1}{2} \frac{\delta}{1 - \delta} \boldsymbol{\beta}^{\top} \text{diag}(\mathbf{X}^{\top} V(\boldsymbol{\beta}) \mathbf{X}) \boldsymbol{\beta} \\ &= \frac{1}{2} \frac{\delta}{1 - \delta} \sum_{i,j} x_{ij}^2 \beta_j^2 p_i(1 - p_i) \end{aligned}$$

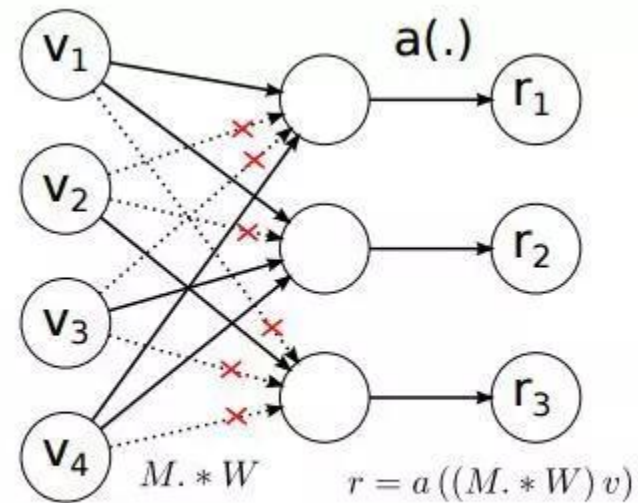
Dropout偏好置信度高的预测和小的 β

DropConnect

- DropConnect: 将节点中的每个与其相连的输入权值以 $1-p$ 的概率变成0
- Dropout: 随机的将隐层节点的输出变成0



DropOut Network



DropConnect Network