- Problem to Solve
  - ① Internal Covariate Shift
  - ② pathological curvature of first-order gradient descent

→ makes sense to apply normalization + saturating nonlinearity

→ data-dependent initialization

→ layer norm in conv is against convolution property

- Previous Approach

  → Weighting / domain adaptation for covariate shift

  → Whitening activation

* Batch normalization

* Advantages / Acceleration / Criticism of BN
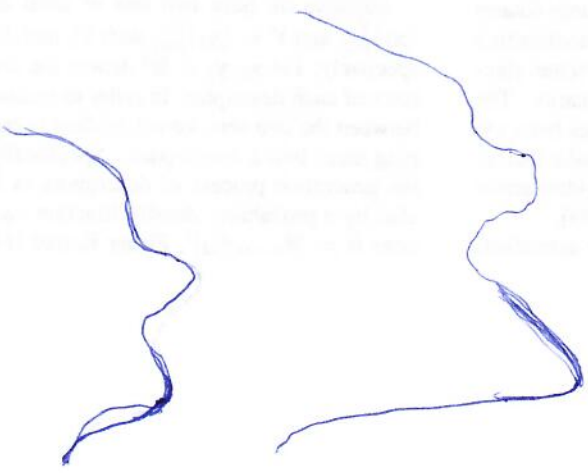
* Layer normalization

* Weight normalization

* comparisons regarding ⟨ invariance / purposes

# (✶) internal <u>covariate</u> shift
##            ‖
##      dependent variable

loss $\ell = F_2(F_1(u, \theta_1), \theta_2)$

$\qquad = F_2(\underbrace{x}_{x}, \theta_2)$ : sub network with sub input

$g$ = saturating (regime of) the nonlinearity;

- if $F_2$ contains

if input of $g \uparrow$ ; $g'|_v \to 0$    eg. $g(v) = \frac{1}{1+\exp(-v)}$

thus $F_1$ trains slowly
and $F_1(u, \theta_1)$ moves to saturated regimes.

the goal of batch normalization
: make the distribution of $x$
 (input of saturating nonlinearity)
 fixed of time & batch

---

# (✶✶) naive whitening activation

if normalization parameters are outside gradient descent, input $\hat{x} = Wu + b$

eg. $\hat{x} = x - E[x]$ where $x = Wu + b$

$b \to b + \Delta b$ do not affect $\hat{x}$

$x - E[x] = Wu + b - E[Wu + b]$
$\qquad\qquad = Wu + b + \Delta b - E[Wu + b + \Delta b]$

$b$ will explode without reducing loss function

To solve ✶ and ✶✶
- fixed distribution over time.
- differentiable
- preserve normalization param to network

result (advantages)
• use of saturating nonlinearity
• increase of learning rate
• model regularization due to sampling
• More resilient parameter scales to <u>initialization</u>
• conjecture: condition # ~1

⊛ BN  (④ more tricks in 4.2.1)

① mini-batch statistics to estimate
   mean and variance (decorrelated features)

$$\mu_B = \frac{1}{m} \sum_{i=1}^{m} x_i \qquad \sigma_B^2 = \sum_{i=1}^{m} (x_i - \mu_B)^2$$

$$\hat{x_i} = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \varepsilon}} \qquad y_i = \gamma \hat{x_i} + \beta$$

parameter to be learned $B = \{x_1, ..., x_m\}$

from each batch

② training

   - take $y_i$ in ① as inputs
   - train to optimize $\{\gamma^{(k)}, \beta^{(k)}\}_k$ — per feature map (eg? in convolution)
     together with other parameters

③ inference using unbiased estimators

   $$E[x] = E_B[\mu_B] \qquad Var[x] = \frac{m}{m-1} E_B[\sigma_B^2]$$

   $$y = \frac{\gamma}{\sqrt{Var[x] + \varepsilon}} \cdot x + \left( \beta - \frac{\gamma E[x]}{\sqrt{Var[x]+\varepsilon}} \right)$$
   ↳ input to the network

⊛ for convolution
   → immediately before nonlinearity
   → jointly normalize all activations
     in a minibatch over all locations
     of convolution layer with
     different $\gamma^{(k)}, \beta^{(k)}$ pairs — per feature map
   : for convolution property!

⊛ Criticism of BN in LN perspectives

⊕ ~~late~~ mini-batch statistics
→ are only estimates

② mini-batch size in constrained

③ different parameters for each activation

⊛Φ dependency within mini-batch

⊛ layer normalization

→ variable length of RNN ⑦

→ normalization statistics
over [all the hidden units] in [the same layer]
[per sample]

$$\mu^\ell = \frac{1}{H} \sum_{i=1}^{H} a_i^\ell \qquad \sigma^\ell = \sqrt{\frac{1}{H} \sum_{i=1}^{H} (a_i^\ell - \mu^\ell)^2}$$

→ in CNN : batch norm outperforms

→ RNN & Online
↳ "scaling"
⊛ gained parameters ←——— [incoming weight]

⊛ RNN

$$a^t = W_{hh} h^{t-1} + W_{xh} \hat{x}^t$$

$$h^t = f\left[\frac{g}{\sigma^t} \odot (a^t - \mu^t) + b\right]$$
↑
element-wise
multiplication

$$\mu^t = \frac{1}{H} \sum_{i=1}^{H} a_i^t$$

$$\sigma^t = \sqrt{\frac{1}{H} \sum_{i=1}^{H} (a_i^t - \mu^t)^2}$$

learning rate regularization

(robust
(resilient)) to input parameter scaling

**①** Pathological curvature
of the objective at optimum.

( ∘ the condition # of the Hessian matrix
( at optimum is low; "unstable gradient descent" )

→ curvature ~ parameterization

**②** whitening gradient (natural gradient)

(the cost)

- left multiply (Fisher info matrix)⁻¹

→ approximation & overhead

---

⊛ weight normalization

**①.** $\vec{w} = \dfrac{g}{\|\vec{v}\|}\vec{v}$   weight vector

$\vec{w} \rightarrow \dfrac{\vec{v}}{\|\vec{v}\|}, g$

$(g \rightarrow e^s)$

$y = \phi(\vec{w}\cdot\vec{x}+b)$

∴ $\nabla_g L = \dfrac{\nabla_{\vec{w}}L \cdot \vec{v}}{\|\vec{v}\|}$

$\nabla_{\vec{v}}L = \dfrac{g}{\|\vec{v}\|}\nabla_{\vec{w}}L - \dfrac{g\nabla_g L}{\|\vec{v}\|^2}\vec{v}$

$= \underset{\uparrow \text{ scale}}{\dfrac{g}{\|\vec{v}\|}} \underset{\uparrow \text{ projection}}{M_{\vec{w}}} \nabla_{\vec{w}}L$   where $M_{\vec{w}} = I - \dfrac{\vec{w}\vec{w}'}{\|\vec{w}\|^2}$

¡ compute orthogonal increment to the current $\vec{w}$

→ self stabilizing is not compatible with Adam; momentum optimizers.

¡ $Cov(\nabla_{\vec{v}}L) = ... (\dfrac{g^2}{\|\vec{v}\|^2}) M_{\vec{w}} Cov(\nabla_{\vec{w}}L)M_{\vec{w}}$  ~ ①

¡ stabilizing noise.

⊛ Weight normalization

2. Data-dependent Initialization
   (∵ missing scaling of features)

   $$y = \phi\left(\frac{g}{\|v\|} \cdot v \cdot x + b\right) \quad \text{then}$$

   initialize $g \leftarrow \frac{1}{\sigma[t]}$ , $b \leftarrow \frac{-\mu[t]}{\sigma[t]}$

   where $\sigma[t], \mu[t]$ : batch-statistics

   → not applicable for RNN

3. Mean-only Batch normalization
   $$\hat{z} = t - \mu[t] + b \quad \text{where } t = \vec{w} \cdot \vec{x}$$
   $$y = \phi(\hat{z})$$

   $\mu[t]$ : $\underset{\downarrow}{\text{running avg}}$ of mini batch
   $\underset{\downarrow}{\text{test time}}$

⊕ advantages
   faster , robust to noise

Weight matrix
BN, WN, LN — invariant over scaling
LN — invariant over centering.

Weight vector (feature)
BN, WN — invariant over scaling
BN, WN — invariant over scaling

Dataset ~~set~~
BN, LN — invariant over scaling
BN — invariant over centering

Single training
LN — invariant over scaling

---

⊛ Riemannian metric (curvature)

under KL

$$ds^2 = D_{KL}[P(y|\tilde{x};\theta) \| P(y|\tilde{x};\theta+\delta)]$$
$$\approx \frac{1}{2}\delta^T F(\theta)\delta$$

where $F(\theta) = \mathop{\mathbb{E}}_{\substack{\tilde{x}\sim P(x) \\ y\sim P(y|\tilde{x})}} \left[ \frac{\partial \log P(y|x;\theta)}{\partial \theta} \frac{\partial \log P(y|x;\theta)}{\partial \theta} \right]$

in GLM $\left( \log P(y|x;\omega,b) = \frac{(a+b)y - \eta(a+b)}{\phi} + C(y,\phi) \right.$ ← log partition
$\quad \mathbb{E}[y|x] = \hat{f}(a+b) \quad Var[y|x] = \phi \hat{f}'(a+b)$
$\quad a = \omega x$

$$\Rightarrow F(\theta) = \mathop{\mathbb{E}}_{\tilde{x}\sim P(\tilde{x})} \left[ \frac{Cov[y|\tilde{x}]}{\phi^2} \right] \otimes \begin{bmatrix} \tilde{x}\tilde{x}^T & \tilde{x}^T \\ \tilde{x}^T & 1 \end{bmatrix}$$

in normalized GLM $\quad$ g: parameter scales

$$[F_{ij}] = \mathop{\mathbb{E}}_{x\sim P(x)} \left[ \frac{Cov[y_i,y_j|\tilde{x}]}{\phi^2} \begin{bmatrix} \frac{g_ig_j}{\sigma_i\sigma_j}x_ix_j^T & \frac{g_i}{\sigma_i}x_i \\ \frac{g_j}{\sigma_j}x_j^T & 1 \end{bmatrix} \right.$
$\chi_i = \frac{x_i - \mu_i}{\sigma_i} \quad \frac{\partial \mu_i}{\partial \omega_k}$
$\chi_i = \frac{x_i - \mu_i}{\sigma_i} \quad \frac{g_i - \mu_i}{\sigma_i} \frac{\partial \sigma_i}{\partial \omega_k}$

$\begin{bmatrix} \frac{g_i g_j (a_i-\mu_i)}{\sigma_i\sigma_j} x_k \frac{g_k}{\sigma_k} & x_i \frac{g_i g_j(a_j-\mu_j)}{\sigma_i\sigma_j} \\ \frac{1}{\sigma_j} & \frac{a_j-\mu_j}{\sigma_j} \\ \frac{\chi_i^T g(a_i-\mu_i)}{\sigma_i\sigma_j} & \frac{a_i-\mu_i}{\sigma_i} \frac{(a_i-\mu_i)(a_j-\mu_j)}{\sigma_i\sigma_j} \end{bmatrix}$

↳ as ω↑ : output fixed : hard to change ω : lr ↓
↳ robust of input parameter scale due to $\frac{g_i}{\sigma_i}$ scale