# LANGEVIN DYNAMICS

$$d\vec{\theta} = -\nabla_\theta U(\theta)\, dt + \sqrt{2}\, dW_t \quad (LD)$$

targets

$$p(\theta) = e^{-U(\theta)}/Z$$

Time-Evolution Obeys Fokker-Planck Eq:

$$\partial_t p(\theta) = \partial_{\theta_i}\left((\partial_{\theta_i} U)\, p(\theta)\right) + \partial_{\theta_i}\partial_{\theta_i}\left(p(\theta)\right)$$

$$\boxed{\text{Conservation of Probability}}$$

1-D: $\quad \partial_t p = \dfrac{\partial}{\partial \theta}\left(\dfrac{\partial U}{\partial \theta}\, p(\theta)\right) + \dfrac{\partial^2 p}{\partial \theta^2}$

If $\partial_t p = 0$ check $p \underset{=}{=} e^{-U}/Z$
is a solution

$$\overset{\partial p/\partial \theta}{\frac{\partial}{\partial \theta}\left(+\frac{\partial U}{\partial \theta} e^{-U}/Z\right) + \frac{\partial}{\partial \theta}\left(-U e^{-U}/Z\right) = 0}$$

$\checkmark$

# SGLD

$$\Delta \theta_t = \frac{\varepsilon_t}{2}\left(\nabla \log p(\theta_t) + \frac{N}{n}\sum_{i=1}^{n}\nabla \log p(x_i|\theta_t)\right) + \mathcal{N}(0, \varepsilon_t)$$

$$\widehat{\nabla} U(\theta) \approx \nabla U(\theta_t) + g(\theta_t)$$

Noisy Gradient          subsampling noise

$$V\left(\begin{array}{c}\text{Gradient}\\\text{Noise}\end{array}\right) = V\left(\frac{\varepsilon_t}{2} g(\theta_t)\right) = \frac{\varepsilon_t^2}{2} V(g(\theta_t))$$

$$V\left(\begin{array}{c}\text{Real}\\\text{Noise}\end{array}\right) = V(\mathcal{N}(0, \varepsilon_t)) = \varepsilon_t$$

if $\varepsilon_t \to 0$, SGLD $\approx$ LD

If $\sum_t \varepsilon_t^2 < \infty$, $\sum_t \varepsilon_t \to 0$

RM conditions

SGD converges to local mode.

As $\varepsilon_t \to 0$ MH Rejection Probability $\to 0$.

# SGLD for BNN

$$=$$

## SGD for NN

$$+$$

## Noise

Prior on Weights : $p(w)$

Probabilistic Output :

     i.e. Softmax $\xleftarrow{\text{Bregman Divergence}}$ Cross-Entropy Loss

     Gaussian $\longleftrightarrow$ $L^2$ loss

Do BACKPROP on $W + \mathcal{N}(0, \varepsilon_t)$

     to sample

# BDN

Predictive Distribution

$$\overset{\circ}{q}(y|x) = \frac{1}{S} \sum_{S=1}^{S} p(y|x, \theta^s)$$

$S$ is $\overset{MC}{\text{samples}}$ of "teachers"

$$|\theta^s| \approx 10^6 - 10^8 \text{ params}$$

$$q(y|x) \rightsquigarrow \text{"teachers"} \rightarrow \text{Bayesian Predictive Ensemble}$$

$$S(y|x,w) \rightsquigarrow \text{"student"} \rightarrow \text{Deep Net}$$

~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~

SGD trains $q(\theta)$ online while simultaneously minimizing

$$\text{mini} \quad KL\left(q(y|x) \| S(y|x,w)\right)$$

Using SGD

$$L(w|x) = KL(p(y|x, D_N) \,||\, S(y|x, w))$$

$$= -\mathbb{E}_{p(y|x, D_N)} \log S(y|x, w)$$

$$= -\int \left[ \int p(y|x, \theta) p(\theta|D_N) \, d\theta \right] \log S(y|x, w) \, dy$$

$$\doteq -\int p(\theta|D_N) \int p(y|x, \theta) \log S(y|x, w) \, dy \, d\theta$$

$$= -\int p(\theta|D_N) \left[ \mathbb{E}_{p(y|x, \theta)} \log S(y|x, w) \right] d\theta \quad ? \quad ?$$

$x$ is the input data to
## STUDENT

Monte-Carlo to integrate out $x$
is high-dimensions, near "training data" $\boxed{D}$

$$\hat{L}(w) \approx \frac{1}{|D'|} \sum_{x' \in D'} L(w|x')$$

$$\hat{L}(w) \approx -\frac{1}{|\theta||D'|} \sum_{\theta^* \in \Theta} \sum_{x' \in D'} \mathbb{E}_{p(y|x, \theta^*)} \log S(y|x, w)$$

# Distilled/Online SGLD $\begin{cases} \theta & \text{Teacher} \\ W & \text{Student} \end{cases}$

for $t = 1:T$ do

Update $\theta$: (SGLD Step)

$$\theta_{t+1} = \theta_t + \frac{\varepsilon_t}{2}\left(\nabla_\theta \log p(\theta) + \frac{N}{n} \sum_{i \in [n]} \nabla_\theta \log p(y_i | x_i, \theta)\right) + \mathcal{N}(0, \varepsilon_t)$$

Update $W$: (Student Step)

Sample $D'$ from student generator

$$W_{t+1} = W_t - \rho_t \left(\frac{1}{|D'|} \sum_{x' \in D'} \nabla_W \hat{L}(w, \theta_{t+1} | x')\right)$$

$$+ \underbrace{\gamma W_t}_{L^2 \text{ Reg.}}$$

Recall for softmax output

$$\hat{L}(W, \theta_{t+1} | x) = -\sum_{k=1}^{K} \underbrace{p(y=k | x, \theta^*)}_{\text{TEACHER}} \log \underbrace{S(y=k | x, W)}_{\text{STUDENT}}$$

USE SGD (BACKPROP)
to tRAIN