# TTIC 31230, Fundamentals of Deep Learning

David McAllester, Autumn 2023

## Continuous Time Models of SGD

Gradient Flow

The Diffusion SDE

The Langevin SDE

General SDEs

The SGD SDE

# Gradient Flow

Gradient flow is a non-stochastic (deterministic) model of stochastic gradient descent (SGD).

Gradient flow is defined by the total gradient differential equation

$$\frac{d\Phi}{dt} = -g(\Phi) \qquad g(\Phi) = \nabla_\Phi \, E_{(x,y)\sim\text{Train}} \, \mathcal{L}(\Phi, x, y)$$

We let $\Phi(t)$ be the solution to this differential equation satisfying $\Phi(0) = \Phi_{\text{init}}$.

# Gradient Flow

$$\frac{d\Phi}{dt} = -g(\Phi)$$

For small values of $\Delta t$ this differential equation can be approximated by

$$\Delta \Phi = -g(\Phi)\Delta t$$

# Time as the Sum of the Learning Rates

Consider the update.

$$d\Phi = -g(\Phi)dt \qquad \text{Gadient Flow}$$

$$\Phi_{t+\Delta t} \approx \Phi_t - g\Delta t$$

$$\Phi_{t+\eta} \approx \Phi_t - \eta_t \, \hat{g} \qquad \text{SGD}$$

We will show that as $\eta_t \to 0$ we have that SGD converges to Gradient Flow.

# Gradient Flow and SGD

Consider a sequence of model parameters $\Phi_1, \ldots, \Phi_N$ produced by SGD with

$$\Phi_{i+1} = \Phi_i - \eta \hat{g}_i$$

and where $\hat{g}_i$ is the gradient of the $i$th randomly selected training point.

Take $\eta \to 0$ and $N \to \infty$ using $N = t/\eta$. We will show that in this limit for SGD we have that $\Phi_N$ converges to $\Phi(t)$ as defined by gradient flow.

# Gradient Flow and SGD

For $\Phi_{i+1} = \Phi_i - \eta \hat{g}_i$ we divide $\Phi_1, \ldots, \Phi_N$ into $\sqrt{N}$ blocks.

$$\left(\Phi_1, \ldots, \Phi_{\sqrt{N}}\right) \left(\Phi_{\sqrt{N}+1}, \ldots, \Phi_{2\sqrt{N}}\right) \cdots \left(\Phi_{T-\sqrt{N}+1}, \ldots, \Phi_N\right)$$

For $\eta \to 0$ and $N = t/\eta$ we have $\eta\sqrt{N} \to 0$ which implies

$$\Phi_{\sqrt{N}} \sim \Phi_0 - \eta\sqrt{N}g$$

where $g$ is the average (non-stochastic) gradient.

Since the gradients within each block become non-stochastic, we are back to gradient flow.

# Stochastic Differential Equations (SDEs)

SDEs are important in deep learning for understanding SGD, incorporating Bayesian priors into SGD, and in modern diffusion models.

We will start with the simplest SDE: Brownian motion. Brownian motion is the SDE used in diffusion models.

# Diffusion Models: Brownian Motion

Consider a discrete-time process $z(0), z(1), z(2), z(3), \ldots$ with $z(n) \in \mathbb{R}^d$ defined by

$$z(n+1) = z(n) + \sigma\epsilon, \quad \epsilon \sim \mathcal{N}(0, I)$$

We can sample from $z(n)$ using

$$z(n) = z(0) + \sigma\epsilon\sqrt{n}, \quad \epsilon \sim \mathcal{N}(0, I)$$

# Brownian Motion

Fix a numerical time step $\Delta t$ and consider a discrete-time process $z(0),\ z(\Delta t),\ z(2\Delta t),\ z(3\Delta t)\ \ldots$

$$z(t + \Delta t) \;=\; z(t) + \sigma\epsilon\sqrt{\Delta t}, \quad \epsilon \sim \mathcal{N}(0, I)$$

$$z(t + n\Delta t) \;=\; z(t) + \sigma\epsilon\sqrt{n\Delta t}, \quad \epsilon \sim \mathcal{N}(0, I)$$

We now take the limit of this numerical simulation as $\Delta t \to 0$. In this limit we can sample directly from $z(t)$ using

$$z(t + \Delta t) = z(t) + \sigma\epsilon\sqrt{\Delta t} \quad \epsilon \sim \mathcal{N}(0, I)$$

# Sampling from a Bayesian Posterior: Langevin Dynamics

$$\text{Train} = (x_1, y_1), \ldots, (x_n, y_n)$$

The parameters $\Phi$ determine $P_\Phi(y|x)$.

$$p(\Phi|\text{Train}) = \frac{p(\Phi)p(\text{Train}|\Phi)}{p(\text{Train})}$$

$$= \frac{p(\Phi)p(x_1, \ldots, x_n)P_\Phi(y_1, \ldots, y_n|x_1, \ldots, x_n)}{p(x_1, \ldots, x_n)P(y_1, \ldots y_n|x_1, \ldots x_n)}$$

$$= \frac{p(\Phi)P_\Phi(y_1, \ldots, y_n|x_1, \ldots, x_n)}{P(y_1, \ldots y_n|x_1, \ldots x_n)}$$

# A Bayesian Interpretation of Langevin Dynamics

$$\text{Train} = (x_1, y_1), \dots, (x_n, y_n)$$

$$p(\Phi|\text{Train}) = \frac{p(\Phi)P_\Phi(y_1, \dots, y_n|x_1, \dots, x_n)}{P(y_1, \dots y_n|x_1, \dots x_n)}$$

The denominator does does not depend on $\Phi$ which implies

$$p(\Phi|\text{Train}) \propto p(\Phi) \prod_i P_\Phi(y_i|x_i)$$

11

# A Bayesian Interpretation of Langevin Dynamics

$$p(\Phi|\text{Train}) \propto p(\Phi) \prod_i P_\Phi(y_i|x_i)$$

$$\ln p(\Phi|\text{Train}) = \sum_i \ln P_\Phi(y_i|x_i) + \ln p(\Phi) + C$$

$$\text{Define } \mathcal{L}(\Phi) = \frac{1}{n} \sum_i -\ln P_\Phi(y_i|x_i) - \frac{1}{n} \ln p(\Phi)$$

$$= E_{(x,y)\sim\text{Train}} \left[-\ln P_\Phi(y|x)\right] - \frac{1}{n} \ln p(\Phi)$$

This Gives $\quad {\color{red} p(\Phi|\text{Train}) = \frac{1}{Z} \exp\left(-n\mathcal{L}(\Phi)\right)}$

# Sampling from a Bayesian Posterior:

# Langevin Dynamics

Consider gradient flow.

$$\frac{d\Phi(t)}{dt} = -g(\Phi)$$

$$g(\Phi) = \nabla_\Phi \, \mathcal{L}(\Phi)$$

$$\mathcal{L}(\Phi) = E_{(x,y)\sim\mathrm{Pop}} \, \mathcal{L}(\Phi, x, y)$$

# The Langevin SDE

In the Langevin SDE we add Gaussian noise to gradient flow.

$$\Phi(t + \Delta t) = \Phi(t) - g\Delta t + \sigma \epsilon \sqrt{\Delta t}, \quad \epsilon \sim \mathcal{N}(0, I)$$

We will show that the stationary distribution of Langevin Dynamics models a Bayesian posterior probability distrbution on the model parameters where $\sigma$ acts as a temperature parameter.

# The Langevin SDE

$$\Phi(t + \Delta t) = \Phi(t) - g(\Phi)\Delta t + \sigma \epsilon \sqrt{\Delta t}, \quad \epsilon \sim \mathcal{N}(0, I)$$

Let $p(\Phi)$ be a probability density on the parameter space $\Phi$. The density $p(\Phi)$ defines a gradient flow and a diffusion flow.

$$\text{gradient flow} = -p(\Phi)g(\Phi)$$

$$\text{diffusion flow} = -\frac{1}{2}\sigma^2 \nabla_\Phi p(\Phi)$$

The expression for the diffusion flow follows from the Fokker-Planck equation. A derivation of the diffusion flow expression from first princple is given in the appendix.

# The Langevin SDE

$$\text{gradient flow} = -p(\Phi)g(\Phi)$$

$$\text{diffusion flow} = -\frac{1}{2}\,\sigma^2\,\nabla_\Phi(p(\Phi))$$

For the stationary distribution these two flows cancel each other out. In one dimension we have

$$\frac{1}{2}\sigma^2\nabla_\Phi\,p = -p\nabla_\Phi\mathcal{L}$$

# The Langevin Stationary Distribution

$$\frac{1}{2}\sigma^2 \nabla_\Phi\, p = -p\nabla_\Phi\mathcal{L}$$

$$\frac{1}{2}\sigma^2 \frac{\nabla_\Phi\, p}{p} = -\nabla_\Phi\mathcal{L}$$

$$\frac{1}{2}\sigma^2 \left(\nabla_\Phi \ln p\right) = \nabla_\Phi(-\mathcal{L})$$

$$\frac{1}{2}\sigma^2 \ln p = -\mathcal{L} + C$$

$$\textcolor{red}{p(\Phi) = \frac{1}{Z}\exp\left(\frac{-2\mathcal{L}(\Phi)}{\sigma^2}\right)}$$

# A Bayesian Interpretation of Langevin Dynamics

$$p(\Phi|\text{Train}) = \frac{1}{Z}e^{-n\mathcal{L}(\Phi)}$$

$$p_{\text{Langevin}}(\Phi) = \frac{1}{Z}\exp\left(\frac{-2\mathcal{L}(\Phi)}{\sigma^2}\right)$$

Setting $\sigma^2 = \frac{1}{2n}$ gives

$$p_{\text{Langevin}}(\Phi) = p(\Phi|\text{Train})$$

# A General SDE

$$x(t + \Delta t) \; = \; x(t) + \mu(x,t)\Delta t + \sigma(x,t)\epsilon\sqrt{\Delta t}, \quad \epsilon \sim \mathcal{N}(0, I) \quad (1)$$

Here $\sigma(x,t)$ is a matrix.

This is conventionally written as

$$dx = \mu(x,t)dt + \sigma(x,t)dB \quad (2)$$

where $B$ denotes a Weiner process (simple diffusion, aka Brownian motion)

I find (1) more intuitive than (2) but they are the same thing.

# The SGD SDE

We now consider SGD

$$\Phi_{i+1} = \Phi_i - \eta \hat{g}_i$$

We consider $\Phi_i$ and $\Phi_{i+N}$ with $N$ small enough that

$$\Phi_{i+N} \approx \Phi_i$$

.

# Gradient Noise

$$\hat{g} = g(\Phi) + (\hat{g} - g(\Phi))$$

$\hat{g} - g(\Phi)$ has zero mean.

$$\Phi_{i+N} \approx \Phi_i - \eta N g(\Phi) - \eta \sum_{j=1}^{N} (\hat{g}_i - g(\Phi))$$

We pick $N$ large enough that $\sum_{j=1}^{N} (\hat{g}_i - g(\Phi))$ is approximately Gaussian.

# Gradient Noise

$$\Phi_{i+N} \approx \Phi_i - \eta N g(\Phi) - \eta \sum_{j=1}^{N} (\hat{g}_i - g(\Phi))$$

$$\approx \Phi_i - \eta N g(\Phi) - \eta\sqrt{N}\epsilon, \quad \epsilon \sim \mathcal{N}(0, \Sigma)$$

Now define $\Delta t = N\eta$ or $N = \Delta t/\eta$.

$$\Phi(t + \Delta t) \approx \Phi(t) - g(\Phi)\Delta t + \eta\epsilon\sqrt{\Delta t/\eta}, \quad \epsilon \sim \mathcal{N}(0, \Sigma)$$

$$= \Phi(t) - g(\Phi)\Delta t + \sqrt{\eta}\epsilon\sqrt{\Delta t}, \quad \epsilon \sim \mathcal{N}(0, \Sigma)$$

# The SGD SDE

$$\Phi(t + \Delta t) \approx \Phi(t) - g(\Phi)\Delta t + \sqrt{\eta}\epsilon\sqrt{\Delta t}, \quad \epsilon \sim \mathcal{N}(0, \Sigma)$$

$$= \Phi(t) - g(\Phi)\Delta t + \sqrt{\eta}\sigma(\Phi)\epsilon\sqrt{\Delta t}, \quad \epsilon \sim \mathcal{N}(0, I)$$

Here the matrix $\sigma(\Phi)$ is the square root of the covariance matrix $\Sigma(\Phi)$.

# The SGD SDE in One Dimension

$$\Phi(t + \Delta t) = \Phi(t) - g(\Phi)\Delta t + \sqrt{\eta}\sigma(\Phi)\epsilon\sqrt{\Delta t}$$

In one dimension, if the gradient noise $\sigma(\Phi)$ is constant, then the SGD SDE has the same form as Langevin dynamics and we get.

$$p(x) = \frac{1}{Z}\exp\left(\frac{-2\mathcal{L}(x)}{\eta\sigma^2}\right)$$

This is Gibbs and provides an interpretation of the learning rate as temperature.

# The SGD SDE in Higher Dimension

$$\Phi(t + \Delta t) = \Phi(t) - g(\Phi)\Delta t + \sqrt{\eta}\sigma(\Phi)\epsilon\sqrt{\Delta t}$$

This is almost the general case of an SDE.

Here $g(\Phi)$ is the gradient of a scalar function. This is not true for a general SDE.

But the matrix $\sigma(\Phi)$ is arbitrary.

Here the learning rate $\eta$ controls the level of noise but we do not in general have a Gibbs distribution.

# The SGD SDE, A Counter Example

If we have two dimensions $x$ and $y$ where the loss separates as $\mathcal{L}(x, y) = \mathcal{L}(x) + \mathcal{L}(y)$, and the matrix $\sigma(\Phi)$ is constant and diagonal, each dimension behaves as an independent one dimensional SGD and we get.

$$p(x, y) = \frac{1}{Z} \exp \left( \frac{-2\mathcal{L}(x)}{\eta \sigma_x^2} + \frac{-2\mathcal{L}(y)}{\eta \sigma_y^2} \right)$$

This is not Gibbs.

# Langevin-Adaptive SGD

Consider SGD where the update dirction is determined by a matrix $D$.

$$\Phi_{i+1} = \Phi_i - \eta D \hat{g}_i$$

$D$ defines a linear map from dual vectors to primal vectors.

The function defined by $D$ has a meaning indepent of the choice of coordinates.

# Coordinate Independent Formulation of Gradient Noise

We can define the covariance matrix of the noise as

$$\Sigma(\Phi) = E_{\hat{g}} \ (\hat{g} - g)(\hat{g}_i - g)^\top$$

The gradient noise covariance matrix $\Sigma(\Phi)$ defines a linear map from the primal vectors to dual vectors (independent of coordinates).

$$\Sigma(\Phi)\Delta\Phi = E_{\hat{g}} \ (\hat{g} - g)((\hat{g} - g)^\top \Delta\Phi)$$

# Solving for $D$ to Get Langevin

$$\Phi_{i+1} = \Phi_i - \eta D \hat{g}_i$$

Setting $\Delta t = N\eta$ we get

$$\Phi(t + \Delta t) = \Phi(t) - Dg\Delta t + \sqrt{\eta}D\epsilon\sqrt{\Delta t}, \quad \epsilon \sim \mathcal{N}(0, \Sigma(\Phi))$$

Here the noise vector $\epsilon$ is a dual vector.

# Solving for $D$

For a given probability density $p(\Phi)$ over the parameters $\Phi$ the flows are

$$\text{gradient flow} = -pDg$$

$$\text{diffusion flow} = -\frac{1}{2}\eta D\Sigma(\Phi)D\nabla_\Phi p$$

These are vectors in parameter space that are independent of the choice of coordinates.

# Solving for $D$

$$\text{gradient flow} = -pDg$$

$$\text{diffusion flow} = -\frac{1}{2}\eta D\Sigma(\Phi)D\nabla_\Phi p$$

Detailed Balance:

$$\frac{1}{2}\eta D\Sigma(\Phi)D\nabla_\Phi\, p = -pD\nabla_\Phi\mathcal{L}$$

# Solving for $D$

$$\frac{1}{2}\eta D\Sigma(\Phi)D\nabla_\Phi\, p = -pD\nabla_\Phi\mathcal{L}$$

$$\frac{1}{2}\eta D\Sigma(\Phi)D\frac{\nabla_\Phi\, p}{p} = -D\nabla_\Phi\mathcal{L}$$

$$\frac{1}{2}\eta D\Sigma(\Phi)D\left(\nabla_\Phi \ln p\right) = -D\nabla_\Phi\mathcal{L}$$

Setting $D = \Sigma(\Phi)^{-1}$ gives

$$\frac{1}{2}\eta\Sigma(\Phi)^{-1}\left(\nabla_\Phi \ln p\right) = -\Sigma(\Phi)^{-1}\nabla_\Phi\mathcal{L}$$

# The Gibbs distribution

$$\frac{1}{2}\eta\Sigma(\Phi)^{-1}\left(\nabla_\Phi \ln p\right) = -\Sigma(\Phi)^{-1}\nabla_\Phi\mathcal{L}$$

The factors of $\Sigma(\Phi)^{-1}$ now cancel (we can multiply both sides by $\Sigma(\Phi)$) and we get

$$\frac{1}{2}\eta\left(\nabla_\Phi \ln p\right) = -\nabla_\Phi\mathcal{L}$$

This equation is independent of coordinates.

# The Gibbs Distribution

$$\frac{1}{2}\eta \left( \nabla_\Phi \ln p \right) = -\nabla_\Phi \mathcal{L}$$

$$p(\Phi) = \frac{1}{Z} \exp \left( \frac{-2\mathcal{L}(\Phi)}{\eta} \right)$$

# The Gibbs Distribution

For the adaptive update

$$\Phi_{i+1} = \Phi_i - \eta \Sigma(\Phi)^{-1} \hat{g}_i$$

we have a stationary distribution

$$p(\Phi) = \frac{1}{Z} \exp\left(\frac{-2\mathcal{L}}{\eta}\right)$$

# Relating Langevin Dynamics and the SGD-SDE

Consider a Bayesian model:

$$p_{\text{Bayes}}(\Phi|\text{Train}) \;=\; \frac{1}{Z}\, p_{\text{Bayes}}(\Phi, \text{Train}) \quad Z = p_{\text{Bayes}}(\text{Train})$$

$$=\; \frac{1}{Z}\, p_{\text{pri}}(\Phi)P(\text{Train}|\Phi)$$

Consider the loss function:

$$\mathcal{L}(\Phi) \;=\; -\ln P_{\text{Bayes}}(\Phi|\text{Train})$$

We will give an SGD algorithm and learning rate $\eta$ such that:

$$p_{\text{SGD}}(\Phi|\mathcal{L}, \eta) = p_{\text{Bayes}}(\Phi|\text{Train})$$

# The Loss Function

$$p_{\text{Bayes}}(\Phi|\text{Train}) = \frac{1}{Z} \, p_{\text{Bayes}}(\Phi, \text{Train}) \quad Z = p_{\text{Bayes}}(\text{Train})$$

$$= \frac{1}{Z} \, p_{\text{pri}}(\Phi) P(\text{Train}|\Phi)$$

$$= \frac{1}{Z'} \, p_{\text{pri}}(\Phi) \sum_{(x,y)\in\text{Train}} p_{\Phi}(y|x)$$

$$\mathcal{L}(\Phi) = -\ln P_{\text{Bayes}}(\Phi|\text{Train})$$

$$= E_{(x,y)\sim\text{Train}}\left[-\ln(P_{\Phi}(y|x))\right] + \frac{1}{N_{\text{Train}}}\left(-\ln p_{\text{pri}}(\Phi)\right) + \ln Z''$$

# The Gibbs Distribution

For the adaptive update

$$\Phi_{i+1} = \Phi_i - \eta \Sigma(\Phi)^{-1} \hat{g}_i$$

we have a stationary distribution

$$p_{\text{SGD}}(\Phi) = \frac{1}{Z} \exp\left(\frac{-2\mathcal{L}}{\eta}\right)$$

$$= p_{\text{Bayes}}(\Phi | \text{Train}) \quad \text{for } \eta = 2$$

# END

# Appendix: Diffusion Flow

We consider the one dimensional case where we have a function $x(t) \in \mathbb{R}$. We consider a very small time step $\Delta t$ and consider only the diffusion flow.

$$x(t + \Delta t) \;=\; x(t) + \sigma \epsilon \sqrt{\Delta t}, \quad \epsilon \sim \mathcal{N}(0,1)$$

We assume a density $p_x$ of values of $x$ and let $p_\epsilon(\epsilon)$ be the normal distribution $\mathcal{N}(0,1)$ on $\epsilon$.

The quantity of mass transfer in the time interval $\Delta t$ from values above $x$ to values below $x$ is

$$\int_{z=0}^{\infty} p_x(x+z) \, p_\epsilon(\sigma \epsilon \sqrt{\Delta t} \leq -z) dz$$

$$= \int_{z=0}^{\infty} p_x(x+z) \, p_\epsilon\left(\epsilon \leq \frac{-z}{\sigma\sqrt{\Delta t}}\right) dz$$

$$= \int_{z=0}^{\infty} p_x(x+z) \, \Phi\left(\frac{-z}{\sigma\sqrt{\Delta t}}\right) dz$$

where $\Phi$ is the cummulative function of the Gaussian.

# Appendix: Diffusion Flow

The quantity of mass transfer in the time interval $\Delta t$ from values above $x$ to values below $x$ is

$$\int_{z=0}^{\infty} p_x(x+z) \ \Phi\left(\frac{-z}{\sigma\sqrt{\Delta t}}\right) dz$$

By a change of variables $u = z/(\sigma\sqrt{\Delta t})$ we get

$$\int_{u=0}^{\infty} p_x(x + \sigma\sqrt{\Delta t}\ u) \ \Phi(-u)\sigma\sqrt{\Delta t}\ du$$

As $\Delta t \to 0$ we can use the first order Taylor expansion of the density.

$$\sigma\sqrt{\Delta t} \int_{u=0}^{\infty} \left(p_x(x) + \sigma\sqrt{\Delta t}\ u\frac{dp_x(x)}{dx}\right) \ \Phi(-u)\ du$$

# Appendix: Diffusion Flow

$$\sigma\sqrt{\Delta t}\int_{u=0}^{\infty}\left(p_x(x) + \sigma\sqrt{\Delta t}\, u\frac{dp_x(x)}{dx}\right)\ \Phi(-u)\ du$$

$$= \ \sigma\sqrt{\Delta t}\, p_x(x)\left(\int_{u=0}^{\infty}\Phi(-u)\ du\right) + \sigma^2\Delta t\frac{dp_x(x)}{dx}\left(\int_{u=0}^{\infty}u\Phi(-u)du\right)$$

A similar alanysis shows that the mass transfer from lower values to higher values is

$$= \ \sigma\sqrt{\Delta t}\, p_x(x)\left(\int_{u=0}^{\infty}\Phi(-u)\ du\right) - \sigma^2\Delta t\frac{dp_x(x)}{dx}\left(\int_{u=0}^{\infty}u\Phi(-u)du\right)$$

The net mass transfer in the positive $x$ direction is the second minus the first or

$$= \ -2\sigma^2\Delta t\frac{dp_x(x)}{dx}\left(\int_{u=0}^{\infty}u\Phi(-u)du\right)$$

# Appendix: Diffusion Flow

The net mass transfer in the positive $x$ direction is

$$-2\sigma^2 \Delta t \frac{dp_x(x)}{dx} \left( \int_{u=0}^{\infty} u \Phi(-u) du \right)$$

Note that the mass transfer is proportional to $\Delta t$. Dividing by $\Delta t$ gives the flow per unit time.

$$\text{Diffusion flow} \quad = -\alpha \sigma^2 \frac{dp_x(x)}{dx} \qquad \alpha = 2 \int_{u=0}^{\infty} u \Phi(-u) du$$

$\alpha$ can be calculated using integration by parts.

$$
\begin{aligned}
\alpha &= 2 \int_0^{\infty} u \Phi(-u) du \\
&= \int_0^{\infty} \Phi(-u) du^2 \\
&= u^2 \Phi(-u)|_0^{\infty} + \int_0^{\infty} u^2 \phi(-u) du \quad \text{where } \phi \text{ is the Gaussian density} \\
&= \int_0^{\infty} u^2 \phi(-u) du \\
&= \frac{1}{2}
\end{aligned}
$$

# Appendix: Diffusion Flow

We now have that the diffusion flow is

$$\text{Diffusion flow} \quad = -\frac{1}{2} \sigma^2 \frac{dp_x(x)}{dx}$$

For dimension larger than 1 we have

$$\text{Diffusion flow} \quad = -\frac{1}{2} \Sigma \nabla_x p_x$$

Where $\Sigma = E\left(\hat{g} - g\right)(\hat{g} - g)^\top$ is the covariance matrix of the gradient noise.

Here we have derived this from first principle but it also follows from the Fokker–Planck equation (see Wikipedia).

# END