

# Training Products of Experts by Minimizing Contrastive Divergence

Geoffrey E. Hinton

*presented by Frank Wood*



# Goal

- Learn parameters for probability distribution models of high dimensional data
  - (Images, Population Firing Rates, Securities Data, NLP data, etc)

## Mixture Model

$$p(\vec{d} | \theta_1, \dots, \theta_n) = \sum_m \alpha_m f_m(\vec{d} | \theta_m)$$

Use EM to learn parameters

## Product of Experts

$$p(\vec{d} | \theta_1, \dots, \theta_n) = \frac{\prod_m f_m(\vec{d} | \theta_m)}{\sum_{\vec{c}} \prod_m f_m(\vec{c} | \theta_m)}$$

Use Contrastive Divergence to learn parameters.



# Take Home

- Contrastive divergence is a general MCMC gradient ascent learning algorithm particularly well suited to learning Product of Experts (PoE) and energy-based (Gibbs distributions, etc.) model parameters.
- The general algorithm:
  - Repeat Until "Convergence"
    - Draw samples from the current model *starting from the training data*.
    - Compute the expected gradient of the log probability w.r.t. all model parameters over both samples and the training data.
    - Update the model parameters according to the gradient.



# Sampling - Critical to Understanding

- Uniform
  - rand()
  - Linear Congruential Generator
    - $x(n) = a * x(n-1) + b \bmod M$   
0.2311 0.6068 0.4860 0.8913 0.7621 0.4565 0.0185
- Normal
  - randn()
  - Box-Mueller
    - $x_1, x_2 \sim U(0,1) \rightarrow y_1, y_2 \sim N(0,1)$ 
      - $y_1 = \sqrt{-2 \ln(x_1)} \cos(2 \pi x_2)$
      - $y_2 = \sqrt{-2 \ln(x_1)} \sin(2 \pi x_2)$
- Binomial(p)
  - if(rand()) < p
- More Complicated Distributions
  - Mixture Model
    - Sample from a Gaussian
    - Sample from a multinomial (CDF + uniform)
  - Product of Experts
    - Metropolis and/or Gibbs



BROWN

# The Flavor of Metropolis Sampling

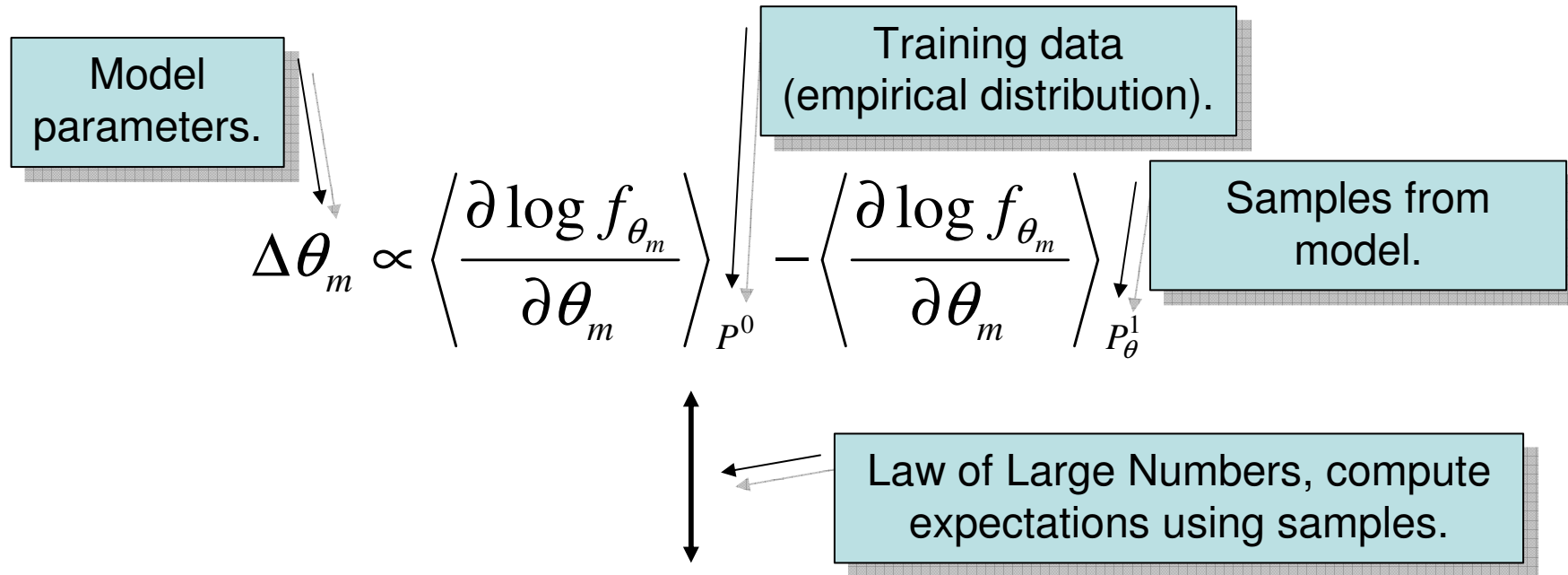
- Given some distribution  $p(\vec{d} | \theta)$ , a random starting point  $\vec{d}_{t-1}$ , and a symmetric proposal distribution  $J(\vec{d}_t | \vec{d}_{t-1})$ .
- Calculate the ratio of densities  $r = \frac{p(\vec{d}_t | \theta)}{p(\vec{d}_{t-1} | \theta)}$  where  $\vec{d}_t$  is sampled from the proposal distribution.
- With probability  $\min(r, 1)$  accept  $\vec{d}_t$ .
- Given sufficiently many iterations

$$\{\vec{d}_n, \vec{d}_{n+1}, \vec{d}_{n+2}, \dots\} \sim p(\vec{d} | \theta)$$

Only need to know the distribution up to a proportionality!



# Contrastive Divergence (Final Result!)



$$\Delta \theta_m \propto \frac{1}{N} \sum_{d \in D} \frac{\partial \log f_{\theta_m}(d)}{\partial \theta_m} - \frac{1}{N} \sum_{c \sim P_\theta^1} \frac{\partial \log f_{\theta_m}(c)}{\partial \theta_m}$$

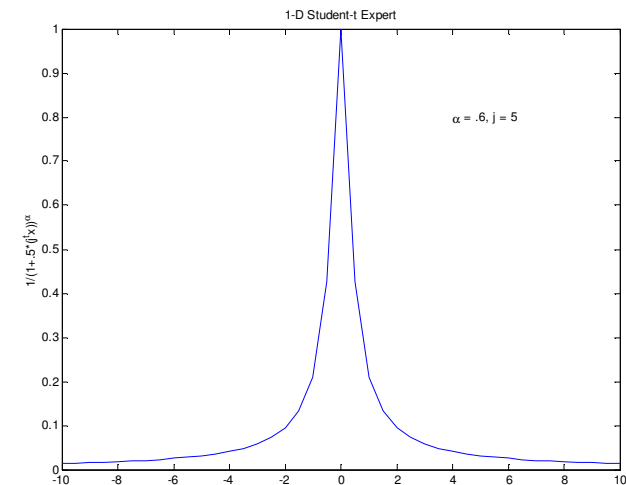


*Now you know how to do it, let's see why this works!*

# But First: The last vestige of concreteness.

- Looking towards the future:
  - Take  $f$  to be a Student-t.

$$f_{\theta^m}(\vec{d}) = f_{\alpha_m, \vec{j}_m}(\vec{d}) = \frac{1}{\left(1 + \frac{1}{2}(\vec{j}_m^T \vec{d})\right)^{\alpha_m}}$$



- Then (for instance)

$$\frac{\partial \log f_{\alpha_m, \vec{j}_m}(\vec{d})}{\partial \alpha_m} = \frac{-\partial \alpha_m \log\left(1 + \frac{1}{2}(\vec{j}_m^T \vec{d})\right)}{\partial \alpha_m} = -\log\left(1 + \frac{1}{2}(\vec{j}_m^T \vec{d})\right)$$

Dot product  $\Leftrightarrow$  Projection  $\Leftrightarrow$  1-D Marginal



# Maximizing the training data log likelihood

- We want maximizing parameters

$$\arg \max_{\theta_1, \dots, \theta_n} \log p(D | \theta_1, \dots, \theta_n) = \arg \max_{\theta_1, \dots, \theta_n} \log \prod_{\vec{d} \in D} \left( \frac{\prod_m f_m(\vec{d} | \theta_m)}{\sum_{\vec{c}} \prod_m f_m(\vec{c} | \theta_m)} \right)$$

Standard PoE form

Over all training data.

Assuming d's drawn independently from p()

- Differentiate w.r.t. to all parameters and perform gradient ascent to find optimal parameters.
- The derivation is somewhat nasty.





# Maximizing the training data log likelihood

$$\begin{aligned}\frac{\partial \log p(\mathbf{D} | \theta_1, \dots, \theta_n)}{\partial \theta_m} &= \frac{\partial \log \prod_{\vec{d} \in \mathbf{D}} p(\vec{d} | \theta_1, \dots, \theta_n)}{\partial \theta_m} \\&= \sum_{\vec{d} \in \mathbf{D}} \frac{\partial}{\partial \theta_m} \log p(\vec{d} | \theta_1, \dots, \theta_n) \\&= N \left\langle \frac{\partial \log p(\vec{d} | \theta_1, \dots, \theta_n)}{\partial \theta_m} \right\rangle_{P_\theta^\infty} \star\end{aligned}$$

Remember this  
equivalence!



# Maximizing the training data log likelihood

$$\begin{aligned}\frac{1}{N} \frac{\partial \log p(\mathbf{D} | \theta_1, \dots, \theta_n)}{\partial \theta_m} &= \frac{1}{N} \frac{\partial}{\partial \theta_m} \log \prod_{\vec{d} \in \mathbf{D}} \frac{\prod_m f_m(\vec{d} | \theta_m)}{\sum_{\vec{c}} \prod_m f_m(\vec{c} | \theta_m)} \\&= \frac{1}{N} \sum_{\vec{d} \in \mathbf{D}} \frac{\partial \log f_m(\vec{d} | \theta_m)}{\partial \theta_m} - \frac{1}{N} \sum_{\vec{d} \in \mathbf{D}} \frac{\partial \log \sum_{\vec{c}} \prod_m f_m(\vec{c} | \theta_m)}{\partial \theta_m} \\&= \frac{1}{N} \sum_{\vec{d} \in \mathbf{D}} \frac{\partial \log f_m(\vec{d} | \theta_m)}{\partial \theta_m} - \frac{\partial \log \sum_{\vec{c}} \prod_m f_m(\vec{c} | \theta_m)}{\partial \theta_m}\end{aligned}$$



BROWN

# Maximizing the training data log likelihood

$$\begin{aligned}
 &= \frac{1}{N} \sum_{d \in D} \frac{\partial \log f_m(\vec{d} | \theta_m)}{\partial \theta_m} - \frac{\partial \log \sum_{\vec{c}} \prod_m f_m(\vec{c} | \theta_m)}{\partial \theta_m} \\
 &= \left\langle \frac{\partial \log f_m(\vec{d} | \theta_m)}{\partial \theta_m} \right\rangle_{P^0} - \frac{\partial \log \sum_{\vec{c}} \prod_m f_m(\vec{c} | \theta_m)}{\partial \theta_m} \\
 &= \left\langle \frac{\partial \log f_m(\vec{d} | \theta_m)}{\partial \theta_m} \right\rangle_{P^0} - \frac{1}{\sum_{\vec{c}} \prod_m f_m(\vec{c} | \theta_m)} \frac{\partial \sum_{\vec{c}} \prod_m f_m(\vec{c} | \theta_m)}{\partial \theta_m}
 \end{aligned}$$

log(x)' = x'/x



BROWN

# Maximizing the training data log likelihood

$$\begin{aligned}
 &= \left\langle \frac{\partial \log f_m(\vec{d} | \theta_m)}{\partial \theta_m} \right\rangle_{P^0} - \frac{1}{\sum_{\vec{c}} \prod_m f_m(\vec{c} | \theta_m)} \frac{\partial \sum_{\vec{c}} \prod_m f_m(\vec{c} | \theta_m)}{\partial \theta_m} \\
 &= \left\langle \frac{\partial \log f_m(\vec{d} | \theta_m)}{\partial \theta_m} \right\rangle_{P^0} - \frac{1}{\sum_{\vec{c}} \prod_m f_m(\vec{c} | \theta_m)} \frac{\sum_{\vec{c}} \prod_{j \neq m} f_j(\vec{c} | \theta_j) \partial f_m(\vec{c} | \theta_m)}{\partial \theta_m} \\
 &= \left\langle \frac{\partial \log f_m(\vec{d} | \theta_m)}{\partial \theta_m} \right\rangle_{P^0} - \frac{1}{\sum_{\vec{c}} \prod_m f_m(\vec{c} | \theta_m)} \frac{\sum_{\vec{c}} \prod_m f_m(\vec{c} | \theta_m) \partial \log f_m(\vec{c} | \theta_m)}{\partial \theta_m}
 \end{aligned}$$

$\log(x)' = x'/x$



# Maximizing the training data log likelihood

$$\begin{aligned}
 &= \left\langle \frac{\partial \log f_m(\vec{d} | \theta_m)}{\partial \theta_m} \right\rangle_{P^0} - \frac{1}{\sum_{\vec{c}} \prod_m f_m(\vec{c} | \theta_m)} \frac{\sum_{\vec{c}} \prod_m f_m(\vec{c} | \theta_m) \partial \log f_m(\vec{c} | \theta_m)}{\partial \theta_m} \\
 &= \left\langle \frac{\partial \log f_m(\vec{d} | \theta_m)}{\partial \theta_m} \right\rangle_{P^0} - \sum_{\vec{c}} \left( \frac{\prod_m f_m(\vec{c} | \theta_m)}{\sum_{\vec{c}} \prod_m f_m(\vec{c} | \theta_m)} \frac{\partial \log f_m(\vec{c} | \theta_m)}{\partial \theta_m} \right) \\
 &= \left\langle \frac{\partial \log f_m(\vec{d} | \theta_m)}{\partial \theta_m} \right\rangle_{P^0} - \sum_{\vec{c}} p(\vec{c} | \theta_1, \dots, \theta_n) \frac{\partial \log f_m(\vec{c} | \theta_m)}{\partial \theta_m}
 \end{aligned}$$



# Maximizing the training data log likelihood

$$\begin{aligned} &= \left\langle \frac{\partial \log f_m(\vec{d} | \theta_m)}{\partial \theta_m} \right\rangle_{P^0} - \sum_c p(c | \theta_1, \dots, \theta_n) \frac{\partial \log f_m(\vec{c} | \theta_m)}{\partial \theta_m} \\ &= \left\langle \frac{\partial \log f_m(\vec{d} | \theta_m)}{\partial \theta_m} \right\rangle_{P^0} - \left\langle \frac{\partial \log f_m(\vec{c} | \theta_m)}{\partial \theta_m} \right\rangle_{P_{\theta^\infty}} \end{aligned}$$

Phew! We're done! So:

$$\begin{aligned} \frac{\partial \log p(D | \theta_1, \dots, \theta_n)}{\partial \theta_m} &\Leftrightarrow N \left\langle \frac{\partial \log p_\theta^\infty(\vec{d} | \theta_m)}{\partial \theta_m} \right\rangle_{P^0} \\ &\propto \left\langle \frac{\partial \log f_m(\vec{d} | \theta_m)}{\partial \theta_m} \right\rangle_{P^0} - \left\langle \frac{\partial \log f_m(\vec{c} | \theta_m)}{\partial \theta_m} \right\rangle_{P_{\theta^\infty}} \end{aligned}$$



BROWN

# Equilibrium Is Hard to Achieve

- With:

$$\frac{\partial \log p(D | \theta_1, \dots, \theta_n)}{\partial \theta_m} \propto \left\langle \frac{\partial \log f_m(\vec{d} | \theta_m)}{\partial \theta_m} \right\rangle_{P^0} - \left\langle \frac{\partial \log f_m(\vec{c} | \theta_m)}{\partial \theta_m} \right\rangle_{P_{\theta}^{\infty}}$$

we can now train our PoE model.

- But... there's a problem:
  - $P_{\theta}^{\infty}$  is computationally infeasible to obtain (esp. in an inner gradient ascent loop).
  - Sampling Markov Chain must converge to target distribution. Often this takes a *very* long time!



BROWN

# Solution: Contrastive Divergence!

$$\frac{\partial \log p(D | \theta_1, \dots, \theta_n)}{\partial \theta_m} \propto \left\langle \frac{\partial \log f_m(\vec{d} | \theta_m)}{\partial \theta_m} \right\rangle_{P^0} - \left\langle \frac{\partial \log f_m(\vec{c} | \theta_m)}{\partial \theta_m} \right\rangle_{P_\theta^1}$$

- Now we don't have to run the sampling Markov Chain to convergence, instead we can stop after 1 iteration (or perhaps a few iterations more typically)
- Why does this work?
  - Attempts to minimize the ways that the model distorts the data.



BROWN



# Equivalence of $\operatorname{argmax} \log P()$ and $\operatorname{argmax} \text{KL}()$

$$\begin{aligned} P^0 \| P_\theta^\infty &= \sum_{\vec{d}} P^0(\vec{d}) \log \frac{P^0(\vec{d})}{P_\theta^\infty(\vec{d})} \\ &= \sum_{\vec{d}} P^0(\vec{d}) \log P^0(\vec{d}) - \sum_{\vec{d}} P^0(\vec{d}) \log P_\theta^\infty(\vec{d}) \\ &= H(P^0) - \left\langle \log P_\theta^\infty(\vec{d}) \right\rangle_{P^0} \end{aligned}$$

$$\frac{\partial P^0 \| P_\theta^\infty}{\partial \theta_m} = - \left\langle \frac{\partial \log P_\theta^\infty(\vec{d})}{\partial \theta_m} \right\rangle_{P^0}$$

This is what  
we got out of  
the nasty  
derivation!



# Contrastive Divergence

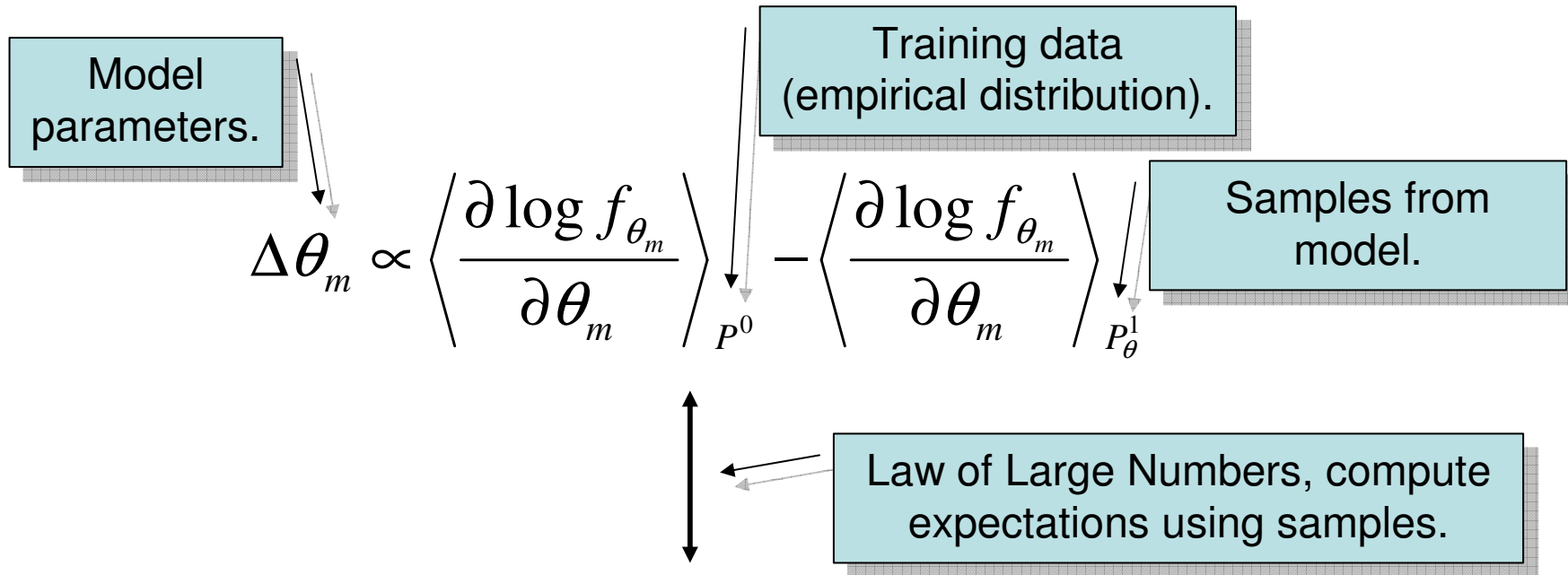
- We want to “update the parameters to reduce the tendency of the chain to wander away from the initial distribution on the first step”.

$$\begin{aligned}
 \frac{\partial}{\partial \theta_m} (P^0 \| P_\theta^\infty - P_\theta^1 \| P_\theta^\infty) &= - \left\langle \frac{\partial \log P_\theta^\infty(\vec{d})}{\partial \theta_m} \right\rangle_{P^0} - \left\langle \frac{\partial \log P_\theta^\infty(\vec{d})}{\partial \theta_m} \right\rangle_{P_\theta^1} \\
 &\propto \left\langle \frac{\partial \log f_m(\vec{d} | \theta_m)}{\partial \theta_m} \right\rangle_{P^0} - \left\langle \frac{\partial \log f_m(\vec{c} | \theta_m)}{\partial \theta_m} \right\rangle_{P_\theta^\infty} - \left\langle \frac{\partial \log f_m(\vec{d} | \theta_m)}{\partial \theta_m} \right\rangle_{P_\theta^1} + \left\langle \frac{\partial \log f_m(\vec{c} | \theta_m)}{\partial \theta_m} \right\rangle_{P_\theta^\infty} \\
 &\propto \left\langle \frac{\partial \log f_m(\vec{d} | \theta_m)}{\partial \theta_m} \right\rangle_{P^0} - \left\langle \frac{\partial \log f_m(\vec{d} | \theta_m)}{\partial \theta_m} \right\rangle_{P_\theta^1}
 \end{aligned}$$



BROWN

# Contrastive Divergence (Final Result!)



$$\Delta \theta_m \propto \frac{1}{N} \sum_{d \in D} \frac{\partial \log f_{\theta_m}(d)}{\partial \theta_m} - \frac{1}{N} \sum_{c \sim P_\theta^1} \frac{\partial \log f_{\theta_m}(c)}{\partial \theta_m}$$



*Now you know how to do it and why it works!*