# Partition Function

In UGM $E(x) = -\log \tilde{p}(x)$ is unnormalized. $p(x, \theta) = \frac{1}{Z(\theta)} \tilde{p}(x, \theta)$. $Z(\theta)$ normalization constant or **Partition Function**. Integral over unnormalized probability of all states, often intractable. Many deep models have it tractable or don't require computation of $p(x)$, others directly address an intractable $Z(\theta)$

## Log Likelihood Gradient for Energy Models

$Z$ depends on parameters → NLL gradient has term of $Z$ gradient $\quad \frac{\partial \log p(x, \theta)}{\partial \theta} = -\frac{\partial E(x)}{\partial \theta} - \frac{\partial \log Z(\theta)}{\partial \theta} \quad Z$ intractable → $\nabla Z$ intractable.

- Monte Carlo Approximation → $\frac{\partial}{\partial \theta} \log Z = E_{x \sim p(x)} \frac{\partial}{\partial \theta} \log \tilde{p}(x)$
  because math

- $-\frac{\partial \log p(x, \theta)}{\partial \theta} = \frac{\partial E(x)}{\partial \theta} - E_{x \sim p(x)} \frac{\partial}{\partial \theta} E(x)$
  - **Positive phase** → push energy down on positive contributions
  - **Negative phase** → pushes energy up everywhere, proportional to current mass

- On a minimum the two terms cancel out

## Stochastic ML and Contrastive Divergence

- Expectation computed with markov chains burn-in every time we need gradient, if SGD → once per step, computationally infeasible
- Balance between pushing up (on model) where data occurs and pushing down where model samples occur.. **here positive phase assumed tractable**
                    max log p              min log Z
- Negative phase approximations make it cheaper to compute but also push down in wrong locations, are points model currently believes in strongly, its incorrect beliefs about the world. Hallucinations / dreams

**Contrastive Divergence**: main cost is burnin → let's draw samples from data distribution to initialize markov chain, free because they are already available. Initially negative phase not accurate because model and data distribution diverge, then better, more accurate.
        **weakness!** fails to suppress spurious modes → high prob regions far from training examples because MCMC initialized on training points won't go there

- CD is biased for RBMs, small bias. Use CD to initialize more expensive MCMC methods  • CD not great for initializing deep models right away because it's difficult
- CD like penalizing a markov chain changing output rapidly when data is from the input.                           to sample hidden units given visible samples, hidden are not in the data
- Useful for pretraining of shallow model                                                                              → we'll need burnin

## Stochastic ML / Persistent CD

Initializes the markov chains at each step with their state from previous steps. Short SGD steps → $m_t \approx m_{t-1}$, so previous samples are fair. Shortens mixing time. SML is best. **Weakness**: if $\epsilon$ too small or $\epsilon$ too large → if SGD moves too fast wrt markov mixing rate. No formal way to check for this but empirically look from negative phase sample variance. **When drawing sample / generative use**: reset markov chain from random, because samples used for training might distort performance. SML has higher variance than CD because different training points in pos/neg phases

- MCMC methods generally cool because they allow decomposition of $\log \tilde{p}$ and $\log Z$ terms → can combine with other methods providing a lower bound on $\log \tilde{p}$
                                                                                                    for positive phase

## Pseudolikelihood

idea! let's avoid computing the partition function altogether. ratio of unnormalized probabilities cancels partitions out.

$p(a|b) = \frac{\tilde{p}(a,b)}{\sum_{a,c} \tilde{p}(a,b,c)}$  a vars we want conditional of, b vars to condition on, c irrelevants • move c into b to reduce cost: $\sum_i^N \log p(x_i | x_{-i})$

- Reduction from $m^n$ to $k \times n$. Ok for large datasets. **Generalized Pseudolikelihood!** m sets of variables appearing together left of conditioning bar
                        $\sum_m \log p(x_s | x_{-s})$              • m=n ; $S = \{i\}$ → pseudolikelihood vanilla
- Poor where we need good model of full joint → density estimation  • Good where data has structure allowing S to capture most correlations ; ie images
- Cannot be used with variational inference or other lower bound techs because has $\tilde{p}$ at denominator → lower bound on denominator is upper bound on expression. Maximizing upper bound makes no sense.
- Still useful to train single layer model    • Per step cost is larger than SML because all conditionals computed; generalized PL can have similar cost
- Implicit prior: all states have more than one variable different from training examples

# SCORE / RATIO MATCHING

ALSO AVOIDS COMPUTING $Z$ AND DERIVATIVES. MINIMIZES EXPECTED SQUARE DIFFERENCE BETWEEN DERIVATIVES OF MODEL LOG PDF AND DATA LOG PDF WRT INPUT.

- $\theta^* = \min_\theta J(\theta) = \frac{1}{2} E_x \| \nabla_x \log p_{MODEL}(x,\theta) - \nabla_x \log p_{DATA}(x) \|_2^2$  REQUIRES KNOWLEDGE OF $p_{DATA}$ → $\hat{J}(\theta) = \frac{1}{M} \sum_{i}^{M} \sum_{j}^{N} \frac{\partial^2}{\partial x_j^2} \log p_{MODEL}(x,\theta) + \frac{1}{2}\left(\frac{\partial}{\partial x_{ij}} \log p_{MODEL}(x,\theta)\right)^2$

  ↳ GSM ALLOWS IT
- DERIVATIVES WRT $x$ → NOCAN USE ON DISCRETE DATA   • NOT COMPATIBLE WITH VARIATIONAL INFERENCE    | $x_1 .. x_M$ TRAINING DATA, DIM $N$
- OK FOR SHALLOW MODELS, NOT FOR DEEP   • LIKE CD WITH NON-GIBBS MARKOV CHAIN MOVING ON GRADIENT
  OR PRETRAINING

**RATIO MATCHING:** FOR BINARY DATA  $J^{RM}(\theta) = \frac{1}{M} \sum_{i}^{M} \sum_{j}^{N} \left( \frac{1}{1 + \frac{p_{MODEL}(x^i,\theta)}{p_{MODEL}(f(x_i),\theta)}} \right)^2$  $j$ IS BIT AT POSITION $j$. $f(x_i)$ FLIPS THE $j$TH BIT.

- RATIO → PARTITION FCN CANCELS OUT
- $N$ TIMES HIGHER COST THAN SML
- SAME PL IMPLICIT PRIOR → HAMMING DISTANCE, ETC   • USEFUL FOR HIGH-DIM SPARSE DATA IE WORD VECTORS

# DENOISING SCORE MATCHING

REGULARIZING SCORE MATCHING WITH FITTING $p_{SMOOTH}(x) = \int p_{DATA}(x+y) q(y|x) dy$ INSTEAD OF $p_{DATA}$. BECAUSE WE ONLY HAVE SAMPLES FROM IT AND GIVEN CAPACITY ANY ESTIMATOR WILL DEGENERATE TO SES OF DIRAC IMPULSES ON TRAINING POINT. SMOOTH WITH $q$ NORMALLY DISTRIBUTED NOISE

- SOME DENOISING AUTOENCODERS CORRESPOND TO ENERGY MODELS WITH **DSM** BUT AE IS LESS EXPENSIVE TO CRUNCH   • POSSIBLE TO DERIVE AE FOR ANY EBM ON REAL DATA

# NOISE-CONTRASTIVE ESTIMATION

**IDEA:** MODEL REPRESENTED AS  $\log p_{MODEL}(x) = \log \tilde{p}_{MODEL}(x,\theta) + c$. $c$ IS APPROX OF $-\log(Z(\theta))$. TREATED AS ONE MORE PARAM, OPTIMIZED AT SAME TIME AND W/SAME ALGO AS $\theta$
NOT A DISTRIBUTION BUT GETS BETTER AS $c$ CONVERGES. • NO CAN DO WITH MLE

- UNSUPERVISED ESTIMATION OF $p(x)$ → **SUPERVISED PROBLEM** → INTRODUCE NOISE $p_{NOISE}(x)$ EASY TO EVAL/SAMPLE → INTRODUCE SWITCH VAR $y$ → JOINT MODEL
- SWITCH DETERMINES WHETHER WE SAMPLE FROM $p_{DATA}(x) / p_{MODEL}(x)$ OR $p_{NOISE}(x)$ → MLE FOR FITTING $p_{JOINT-MODEL}$ TO $p_{TRAIN}$.
- → IS LOGISTIC REGRESSION APPLIED TO LOG-PROBS DIFFERENCE OF MODEL AND NOISE $p_{JOINT-MODEL}(y=1|x) = \sigma(\log p_{MODEL}(x) - \log p_{NOISE}(x))$
- GOOD ON FEW RANDOM VARS → USED FOR MODELING CONDITIONAL WORD DISTRIBUTION GIVEN CONTEXT
- DOES NOT WORK WITH VARIATIONAL BOUNDS / METHODS

# PARTITION FUNCTION ESTIMATION

FOR REALZ. WE NEED IT TO COMPUTE NORMALIZED LIKELIHOOD, MODEL EVALUATION, MONITORING PERFORMANCE, COMPARISON, ETC...

**IDEA:** TO COMPARE MODELS WE USE LIKELIHOOD RATIO → NOT STRICTLY NECESSARY TO HAVE $Z(\theta_M)$, BUT ONLY THEIR RATIO. $\frac{Z(\theta_A)}{Z(\theta_B)}$ → WE NEED RATIO AND ONE OF THE $Z$s AND WE CAN GET THE REST

- $\sum_t \ln \frac{\hat{p}_A(x^t,\theta_A)}{\hat{p}_B(x^t,\theta_B)} - N_{TEST} \ln \frac{Z(\theta_A)}{Z(\theta_B)} > 0 \to A > B$   • $Z_1 = \int \tilde{p}_1(x) dx = Z_0 \int p_0(x) \frac{\tilde{p}_1(x)}{p_0(x)} dx \to \frac{Z_1}{Z_0} \approx \frac{1}{M} \sum_i^M \frac{\tilde{p}_1(x)}{p_0(x)}, x \sim p_0$   • SAMPLE FROM $p_0$, WEIGHT ON RATIO AT SAME VALUE

  $KL(p_0 \| p_1)$ SMALL        IMPORTANCE SAMPLING   MONTECARLO APPROXIMATION
                               $p_0$ PROPOSAL
- WORKS IF $p_1$ AND $p_0$ ARE CLOSE; BUT $p_1$ IS OFTEN MESSY ( MULTIMODAL, HIGH DIM). IF $p_1$ AND $p_0$ NOT CLOSE → SAMPLES WILL MAKE NEGLIGIBLE CONTRIBUTIONS TO
WAT DO? FIND INTERMEDIATE DISTS BETWEEN $p_0$ AND $p_1$                                   SUM
                                                             • **OBS:** $Z(p_0)$ IS KNOWN !!

# ANNEALED IMPORTANCE SAMPLING

- LET'S INTRODUCE INTERMEDIATE DISTRIBUTIONS IN SEQUENCE $p_0 \to p_1$, THE RATIO THEN IS $\frac{Z_1}{Z_0} = \frac{Z_{\eta_1}}{Z_0} \cdots \frac{Z_1}{Z_{\eta_{N-1}}} = \prod_0^{N-1} \frac{Z_{\eta_{j+1}}}{Z_{\eta_j}}$
- LET'S ESTIMATE EACH FACTOR WITH IMPORTANCE SAMPLING → GET THE FINAL RATIO
- SEQUENCE IS DESIGNED TO SUIT THE PROBLEM   • **POPULAR CHOICE:** WEIGHTED GEOMETRIC AVERAGE: $p_{\eta_j} \propto p_1^{\eta_j} p_0^{1-\eta_j}$
- TO DO SAMPLING DEFINE MARKOV CHAIN TRANSITION FCN $T_{\eta_j}(x',x)$ TRANSITION PROBS SO TO LEAVE $p_{\eta_j}$ INVARIANT. $p_{\eta_j}(x) = \int p_{\eta_j}(x') T_{\eta_j}(x',x) dx'$
  → GIBBS, MH, ...
  → SAMPLE FROM $p_0$, USE TRANSITION CHAINS TO SAMPLE FROM INTERMEDIATES UNTIL WE GET TO $p_1$: $x_{\eta_1} \sim p_0(x)$, $x_{\eta_2} \sim T_{\eta_1}(x_{\eta_1},x)$ ...
  → IMPORTANCE WEIGHTS FOR JUMPS ARE PARTIAL RATIOS OF TRANSITIONS $\frac{\tilde{p}_{\eta_1}(x_{\eta_1})}{\tilde{p}_0(x_0)} \cdots \frac{\tilde{p}_1(x_1)}{\tilde{p}_{\eta_{N-1}}(x_{\eta_{N-1}})} = w^k$
  → FINAL RATIO: $\frac{Z_1}{Z_0} \approx \frac{1}{M} \sum_i^M w^k$
- EQUIVALENT TO SIMPLE IMPORTANCE SAMPLING ON EXTENDED STATE SPACE   • MOST COMMON WAY OF ESTIMATING UGM PARTITION FCNS, AT THE MOMENT

# BRIDGE SAMPLING
RELIES ON INTERPOLATION DISTRIBUTION $p_*$   $\frac{Z_1}{Z_0} \approx \sum_i^M \frac{\tilde{p}_*(x_0)}{\tilde{p}_0(x_0)} / \sum_i^M \frac{\tilde{p}_*(x_1)}{\tilde{p}_1(x_1)}$  OPTIMAL DISTRIBUTION IS $\propto \frac{\tilde{p}_0(x) \tilde{p}_1(x)}{R \tilde{p}_0(x) + \tilde{p}_1(x)}$, $R$ IS $\frac{Z_1}{Z_0}$ !!! → RECURSIVE ESTIMATE FROM COARSE START

- AIS > BRIDGE IF $KL(p_0 \| p_1)$ IS LARGE
- **LINKED IMPORTANCE SAMPLING:** USE BRIDGE TO 'INTERPOLATE' AIS SEQUENCE
- **PARTITION FCN TRACKING:** BRIDGE SAMPLING ESTIMATE OF RATIOS OF PARTITION FCNS OF NEIGHBORING PARALLEL TEMPERING CHAINS COMBINED WITH AIS ESTIMATES OVER TIME → LOW VARIANCE $Z$ ESTIMATE AT EVERY ITERATION