

ESSENTIALS

Matrix/Vector Basics

Multiplication: $C = AB \Leftrightarrow c_{ik} = \sum_{j=1}^m a_{ij} b_{jk}$
 Orthogonal Matrix: (full rank, square matrix with orthonormal columns) $A^{-1} = A^T$, $AA^T = A^TA = I$
 $\det(A) \in \{-1, 1\}$, $\det(A^TA) = 1$. preserves: inner product norm, distance, angle, rank, matrix orthogonality
 Inner product: $\langle x, y \rangle = x^T y = \sum_{i=1}^n x_i y_i$.
 $\langle x \pm y, x \pm y \rangle = \langle x, x \rangle \pm 2\langle x, y \rangle + \langle y, y \rangle$.
 $\langle x, y+z \rangle = \langle x, z \rangle + \langle x, y \rangle$. $\langle x, y \rangle = \|x\|_2 \|y\|_2 \cos(\theta)$
 Outer product: $UV^T \Leftrightarrow (UV^T)_{ij} = u_i v_j$
 Norms: $\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2} = \sqrt{\langle x, x \rangle}$, $\|M\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^m m_{ij}^2}$

Probability / Statistics Basics

$P(x) = \sum_{y \in Y} P(x,y)$ $P(x,y) = P(x|y)P(y)$
 $\forall y \in Y: \sum_{x \in X} P(x,y) = 1$
 $P(x|y) = \frac{P(x,y)}{P(y)}$ if $P(y) > 0$
 posterior $P(A|B) = \frac{\text{prior } P(A) \text{ likelihood } P(B|A)}{\text{evidence } P(B)}$

If i.i.d: $P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i)$

If X & Y independent: $P(x,y) = P(x) \Leftrightarrow P(y|x) = P(y)$

Expectation: $E[f(X)] = \sum_{x \in X} P(x)f(x) = \int P(x)f(x) dx$

Variance: $\text{Var}[X] = E[(X - \mu_X)^2] = \sum_{x \in X} (x - \mu_X)^2 P(x) = E[X^2] - E[X]^2$

Expectation Properties

Linearity: $E[X+Y] = E[X] + E[Y]$, $E[aX] = aE[X]$

For constant: $X = c$ then $E[X] = c$, so we can take the expectation of unrelated stuff, and get the same: $\log P_0(x^{(i)}) = E_{z \sim p_0}(\log P_0(x^{(i)}))$

Double: For any X : $E[E[X]] = E[X]$.

Convex Function

$\forall x_1, x_2 \in X, \forall t \in [0,1]:$ (also if $\forall x f''(x) \leq 0$)

$f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2)$

Example: Show that if f is convex, any local optimum is global. Assume there is a local optimum \hat{x} that is not the global optimum x^* , then if we choose t in the ball of the local optimum we know that: $f(t\hat{x} + (1-t)x^*) \geq f(\hat{x})$. Since $f(x^*) < f(\hat{x})$, we have $t \cdot f(\hat{x}) + (1-t)f^*(x) < f(\hat{x})$. So we get $f(t\hat{x} + (1-t)x^*) \geq f(\hat{x}) > t \cdot f(\hat{x}) + (1-t)f^*(x)$ which contradicts the convexity of f .

Jensen Inequality

For convex ϕ : $\phi(E[X]) \leq E[\phi(X)]$ and also $\forall \alpha_i \geq 0, \sum \alpha_i = 1$ and any $x_i > 0$: $\log(\sum \alpha_i x_i) \geq \sum \alpha_i \log(x_i)$

Matrix calculus (numerator layout)

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} \\ \vdots \\ \frac{\partial y_m}{\partial x_n} \end{bmatrix} \quad \frac{\partial \mathbf{X}}{\partial \mathbf{y}} = \begin{bmatrix} \frac{\partial x_1}{\partial y_1} \dots \frac{\partial x_n}{\partial y_1} \\ \vdots \\ \frac{\partial x_1}{\partial y_m} \dots \frac{\partial x_n}{\partial y_m} \end{bmatrix} \quad \frac{\partial \mathbf{y}}{\partial \mathbf{Z}} = \begin{bmatrix} \frac{\partial y_1}{\partial z_1} \dots \frac{\partial y_1}{\partial z_n} \\ \vdots \\ \frac{\partial y_m}{\partial z_1} \dots \frac{\partial y_m}{\partial z_n} \end{bmatrix}$$

$$\frac{\partial \mathbf{A}}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial a_{11}}{\partial x_1} \dots \frac{\partial a_{1n}}{\partial x_n} \\ \vdots \\ \frac{\partial a_{mn}}{\partial x_1} \dots \frac{\partial a_{mn}}{\partial x_n} \end{bmatrix} \quad \frac{\partial \mathbf{X}}{\partial \mathbf{A}} = \begin{bmatrix} \frac{\partial x_1}{\partial a_{11}} \dots \frac{\partial x_n}{\partial a_{11}} \\ \vdots \\ \frac{\partial x_1}{\partial a_{mn}} \dots \frac{\partial x_n}{\partial a_{mn}} \end{bmatrix}$$

1D Derivation Rules

$$\frac{\partial(f \cdot g)}{\partial x} = f' \cdot g + f \cdot g' \quad \frac{\partial f}{\partial x} = \frac{f \cdot g - f' \cdot g}{g^2} \quad \frac{\partial f(g(x))}{\partial x} = f'(g(x))g'(x)$$

$$\frac{\partial x^r}{\partial x} = rx^{r-1} \quad \frac{\partial}{\partial x} \frac{1}{f(x)} = -\frac{f'(x)}{(f(x))^2} \quad (f^g)' = f^g \left(f' \frac{\partial}{\partial x} + g' \ln(f) \right)$$

Multivariate Chain Rule

For function $g: \mathbb{R}^n \rightarrow \mathbb{R}, x \mapsto y$ and $f: \mathbb{R}^m \rightarrow \mathbb{R}, y \mapsto z$ we have: $\frac{\partial z_k}{\partial x_j} = \sum_{i=1}^m \frac{\partial z_k}{\partial y_i} \frac{\partial y_i}{\partial x_j}$
 $\frac{\partial z_k}{\partial x} = \frac{\partial z_k}{\partial y} \frac{\partial y}{\partial x}$ $\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$
 $\frac{\partial f(x, y, t)}{\partial t} = \frac{\partial f}{\partial x} \frac{\partial x}{\partial t} + \frac{\partial f}{\partial y} \frac{\partial y}{\partial t}$ x, y, t can be scalar, vector or matrix.

Derivatives:

$$\frac{\partial}{\partial x}(b^T x) = \frac{\partial}{\partial x}(x^T b) = b \quad \frac{\partial}{\partial x}(x^T x) = 2x \quad \frac{\partial}{\partial x}(x^T Ax) = (A+A^T)x$$

$$\frac{\partial}{\partial x}(b^T Ax) = A^T b \quad \frac{\partial}{\partial x}(c^T x b) = cb^T \quad \frac{\partial}{\partial x}(c^T x^T b) = bc^T$$

$$\frac{\partial}{\partial x}(\|x-b\|_2) = \frac{x-b}{\|x-b\|_2} \quad \frac{\partial}{\partial x}(\|x\|_2^2) = 2x \quad \frac{\partial}{\partial x}(\|x\|_F^2) = 2X$$

$$\frac{\partial}{\partial x} \log(x) = \frac{1}{x} \quad \frac{\partial}{\partial x} \frac{1}{f(x)} = -\frac{f'(x)}{(f(x))^2} \quad \frac{\partial}{\partial x} |x| = (-1)^{|x|} (x > 0)$$

Kullback - Leibler Divergence

$$D_{KL}(P || Q) = - \sum_{x \in X} P(x) \log \left(\frac{Q(x)}{P(x)} \right) = \sum_{x \in X} P(x) \log \left(\frac{P(x)}{Q(x)} \right)$$

$$D_{KL}(P || Q) = \int_{-\infty}^{\infty} P(x) \log \left(\frac{P(x)}{Q(x)} \right) dx$$

Is not symmetric: $D_{KL}(P || Q) \neq D_{KL}(Q || P)$

Entropy: $H(p) = -\sum_i p_i \log p_i$

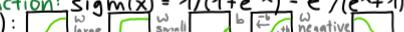
Cross entropy: $H(p, q) = E_p(-\log(q)) = H(p) + D_{KL}(p || q)$

Jensen-Shannon Divergence

$$D_{JS}(P || Q) = \frac{1}{2} D_{KL}(P || \frac{1}{2}(P+Q)) + \frac{1}{2} D_{KL}(Q || \frac{1}{2}(P+Q))$$

NEURAL NETWORK BASICS

Perceptron learning: Init $w^{(0)} = 0$, then while $\exists y$ that is misclassified do $\{w^{k+1} = w^k + n(y - \hat{y})x^k\}$ where $\hat{y} = (\langle w^k x^k \rangle > 0)\}$ if lin. sep. converges in finite time.

Sigmoid function: $\text{sigm}(x) = 1/(1+e^{-x}) = e^x/(e^x+1)$ for $\sigma(w^T x + b)$: 

Logistic regression (MLE): Given $D = \{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$

Model: $y^{(i)} \sim \text{Bernoulli}(\sigma(w^T x^{(i)}))$

$$\text{MLE: } w_{\text{MLE}} = \underset{w}{\operatorname{argmax}} P(D|w) = \underset{w}{\operatorname{argmax}} \prod_{i=1}^n P(y^{(i)} | x^{(i)}, w) = \underset{w}{\operatorname{argmax}} \prod_{i=1}^n [\sigma(w^T x^{(i)})]^{y^{(i)}} [1 - \sigma(w^T x^{(i)})]^{1-y^{(i)}}$$

Neg. log. likelihood: $-\log P(D|w) = \sum_i y^{(i)} \log \pi_i + (1-y^{(i)}) \log (1-\pi_i)$

Softmax function: $\text{softmax}(x_i) = e^{x_i} / \sum_j e^{x_j}$

Loss Functions

Squared Loss: $\frac{1}{2}(y - \hat{y})^2$

Cross entropy: $-y \log(\hat{y}) - (1-y) \log(1-\hat{y})$

Negative log-likelihood: $-\sum_i y_i \log \hat{y}_i$

Why not accuracy: is a step function so a small change in weights might not change accuracy. Gradient is either zero or undefined.

MLE: Write down proba. dist., decompose into per sample prob., minimize negative log likelihood.

Universal Approximation Theorem

$\exists g(x)$ as NN (with non-linear activation), $g(x) \approx f(x)$ and $|g(x) - f(x)| < \epsilon \quad \forall x \in C^\infty(\mathbb{R})$ (continuous functions) (given enough hidden units, one layer is enough)

Training NN: Init weights to small random values, init biases to 0 or small positive values, update the params with GD. (non-lin. activation make loss non-convex)

TRAINING NEURAL NETWORKS

Regularization: Any technique which aims to reduce generalization error (not training error)

Data Standardization: $X_s = (X - \mu)/\sigma$ where $\mu = \sum_{i=1}^n x_i$ and $\sigma = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2}$. Make sure to use μ and σ for the test set!

Batch Normalization

Apply Batch Norm on each layer while computing S^2, μ over the entire batch.
 $\mu = \frac{1}{n} \sum_{i=1}^n z_i$, $S^2 = \frac{1}{n-1} \sum_{i=1}^n (z_i - \mu)^2$, $\tilde{z}_i = \frac{z_i - \mu}{\sqrt{S^2 + \epsilon}}$, $\hat{z}_i = \gamma \tilde{z}_i + \beta$ with $\epsilon > 0$ and γ, β learnable parameters

- Bias in linear / conv layer becomes redundant w/o BN

- Makes weights in deeper layers more robust to changes in shallow layers - slight regularization effect

- Test time: μ, S^2 estimated using exponential average of training time values

Data Augmentation: Generally more data leads to better performance, but acquiring training data is tedious (especially labels). Can generate synthetic data with a model-based approach, but that can lead to discrepancies between the real and synthetic data. Can also expand the data set by transforming existing samples to produce new samples (rotate, scale, noise injection (Gaussian blur, Salt & Pepper noise, channelwise dropout)), need to ensure consistency between transformed data sample and label. For classification introduces invariance $M(T(I)) = M(I)$, for regression introduces equivariance $M(T(I)) = T(M(I))$

Early Stopping: Save parameter where the lowest validation loss occurs then stop if there is no improvement for p (patience) epochs. Then retrain model on entire training set for final test error (either for the same number of epochs (easiest & decent performance) or continue training).

As regularizer: Assuming $\|\nabla C(\theta)\| \leq K$ bounded, learning rate ϵ , initial parameters θ_0 , optimal number of training iterations i^* and param. θ^* learned in i^* steps: $\|\theta_0 - \theta^*\|_2 \leq K\epsilon i^*$.

Pseudocode early stopping

```
i,j,i* <- 0; v <- infinity; theta0 <- theta0
while j < p do:
    p: patience (how long to wait)
    theta <- train(n); i <- i+n train n steps
    v <- validate(theta)
    if v < v then
        j <- 0; theta* <- theta; i* <- i; v <- v
    else
        j <- j+1
return theta*, i*      optimal param. w.r.t. validation set
```

Dropout: Let $r_i^t \sim \text{Bernoulli}(p)$ for $p \in [0, 1]$, $y^t = r^t \cdot y^t$, $z_i^{t+1} = w_i^{t+1} y^t + b_i^{t+1}$, $y^{t+1} = f(z^{t+1})$ (non-lin. act. in

Test time: Would need to evaluate $\sum_r p(r)p(y^r, \theta_r)$ (sum over exponentially many terms), but can instead rescale weights: $W = Wp$

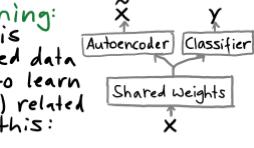
Bagging: Train many models on a bootstrapped dataset, then predict the test set on all models and aggregate (average, majority vote)

Differences to dropout: models are independent (not models sharing weights), train all models until convergence (not train only a small percentage of its models).

Batch Normalization: Problem: in SGD the distribution of x^t constantly changes (internal covariate shift), so we need to adapt constantly. Eliminating this speeds up training as no adaption is necessary.

Solution: Standardize input distribution to each layer i : $\tilde{x}_i = (x_i - \mathbb{E}[x_i]) / \sqrt{\text{Var}[x_i]}$. As this loses some representational power, we perform a linear transform: $\tilde{x}_i = \mu_i - \frac{1}{n} \sum_{j=1}^n x_j + \beta_i$ which makes the transform reversible (with learnable parameters μ and β). Note: sometimes $\sqrt{\text{Var}[x_i]} + \epsilon$ for small constant ϵ (division by 0)

Semi-supervised Learning: Acquiring labeled data is expensive, but unlabeled data usually plentiful. Want to learn from both. Assumes $P(x)$ related to $P(y|x)$. Example of this:



Multitask learning: for datasets with similar samples and related labels. Divide model into shared and task-specific parameters. Assumes that $P(y|x)$ related to $P(y_t|x)$.

Activation Functions:

- Sigmoid: $\sigma(x) = 1/(1+e^{-x})$, $\sigma'(x) = \sigma(x)(1-\sigma(x))$ + interpretation as the firing rate of a neuron. Finite range \rightarrow stable training
- saturates towards either end, so large activations get small gradients. Not zero-centered, always positive (cannot change "direction").

Tanh: $\tanh(x) = 2/\cosh(2x) - 1 = (e^x - e^{-x})/(e^x + e^{-x})$, $\tanh'(x) = 1 - \tanh^2(x)$ + has finite range, zero-centered, has stronger steeper gradients than σ . activations also saturate.

Rectified Linear Unit: $\text{ReLU}(x) = \max(0, x)$, $\text{ReLU}'(x) = 1 \text{ if } x > 0 \text{ else } 0$ + inexpensive operation, sparse activations, greatly accelerates convergence.

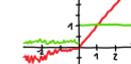
- can blow up the activation and destabilize training. "dying" ReLU problem: units with negative activations get no update

Leaky ReLU: with small pos. const. α : $L\text{ReLU}(x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$, $L\text{ReLU}'(x) = \begin{cases} 1 & \text{for } x \geq 0 \\ \alpha & \text{for } x < 0 \end{cases}$

+ Attempts to solve "dying" ReLU by introducing small negative slope for neg. activ.

Randomized Leaky ReLU: $\alpha \sim U(a, b)$, $f(x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$, $f'(x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$

At test time set $\alpha = (a+b)/2$.



ELBO Derivation: start with data log-likelihood:
 $\log p_{\theta}(x^{(i)}) = \mathbb{E}_{z \sim q_{\theta}(z|x^{(i)})} [\log p_{\theta}(x^{(i)})] = (\text{P}(x^{(i)}) \text{ does not depend on } z)$
 $E_{\mathbb{E}_z} [\log \frac{p_{\theta}(x^{(i)}|z)}{p_{\theta}(z)}] \cdot p_{\theta}(z) = E_{\mathbb{E}_z} [\log p_{\theta}(x^{(i)})] - E_{\mathbb{E}_z} [\log \frac{q_{\theta}(z|x^{(i)})}{p_{\theta}(z)}] + E_{\mathbb{E}_z} [\log \frac{q_{\theta}(z|x^{(i)})}{p_{\theta}(z)}]$
 $= E_{\mathbb{E}_z} [\log p_{\theta}(x^{(i)}|z)] - D_{KL}(q_{\theta}(z|x^{(i)}) || p_{\theta}(z)) + D_{KL}(q_{\theta}(z|x^{(i)}) || p_{\theta}(z|x^{(i)})) \geq 0$

Monte Carlo Gradient Estimator

Often encountered: $\nabla_{\theta} \mathbb{E}_{p_{\theta}(x)}[f(z)]$ with some differentiable function f . We also have a transformation $z = g(\epsilon, \varphi)$.
 $\nabla_{\theta} \mathbb{E}_{p_{\theta}(x)}[f(z)] = \nabla_{\theta} \int p_{\theta}(z) f(z) dz = \nabla_{\theta} \int p(\epsilon) f(g(\epsilon, \varphi)) d\epsilon$
 Note: this works because of the transformation and the differential area under the integral remains the same ($p_{\theta}(z) dz = p(\epsilon) d\epsilon$)
 $= \nabla_{\theta} \int p(\epsilon) f(g(\epsilon, \varphi)) d\epsilon = \nabla_{\theta} \mathbb{E}_{p(\epsilon)}[f(g(\epsilon, \varphi))] = \mathbb{E}_{p(\epsilon)}[\nabla_{\theta} f(g(\epsilon, \varphi))]$
 We can then estimate this using: $\frac{1}{S} \sum_{s=1}^S \mathbb{E}_{p(\epsilon)}[\nabla_{\theta} f(g(\epsilon^{(s)}, \varphi))]$

Or in practice draw samples from the noise distribution and average the gradients.

DKL Non-Negativity: $-D_{KL} = - \int_x p(x) \log \frac{p(x)}{q(x)} dx = \int_x p(x) \log \frac{p(x)}{p(x)} dx \leq \log \int_x p(x) \frac{p(x)}{p(x)} dx = \log \int_x q(x) dx = \log(1) = 0$
 Trick: $E(p(x)) \leq \psi(E(x))$ iff ψ is concave (\log is)

Reparameterization Trick: (to enable VAE backprop)

Given Gaussian with μ and $\sigma^2 \sim N(\mu, \sigma^2)$
 Assume Z underlying random variable $\epsilon \sim N(0, 1)$
 $Z = \mu + \sigma \epsilon = f(x, \epsilon, \theta) \Rightarrow$ can take derivative w.r.t μ, σ (tells us how changes in μ, σ affects output for fixed x)

Advantage of using N : when modelling $p_{\theta}(z)$ and $q_{\theta}(z|x^{(i)})$ with normal distributions with diagonal covariance matrix, it allows drawing easily from the prior & can calculate DKL analytically
Equation for analytical DKL with N :
 For any $p(z) = N(\mu_p, \sigma_p^2 I)$ and $q(z) = N(\mu_q, \sigma_q^2 I)$
 $\int p(z) \log(q(z)/p(z)) dz = -\frac{1}{2} \log(2\pi) - \frac{1}{2} \sum_j \log \sigma_j^2 - \frac{1}{2} \sum_j \frac{\sigma_{q,j}^2 + (\mu_{q,j} - \mu_{p,j})^2}{\sigma_{p,j}^2}$

VAE Backprop: $\mathbb{E}_z [\log p_{\theta}(x^{(i)}|z)] - D_{KL}(q_{\theta}(z|x^{(i)}) || p_{\theta}(z))$
 $= \mathbb{E}_z [\log p_{\theta}(x^{(i)}|z)] + \mathbb{E}_z [\log \frac{p_{\theta}(z)}{q_{\theta}(z|x^{(i)})}] = \mathbb{E}_z [\log \frac{p_{\theta}(z)}{q_{\theta}(z|x^{(i)})}] = \mathbb{E}_z [\log \frac{p_{\theta}(z)}{q_{\theta}(x^{(i)}, z)}] - \mathbb{E}_z [\nabla_{\theta} \log \frac{p_{\theta}(z)}{q_{\theta}(x^{(i)}, z)}]$
 $\approx \sum_k \mathbb{E}_z [\nabla_{\theta} \log(p_{\theta}(z, x^{(i)})) / q_{\theta}(x^{(i)}, z)]$

Limitations of VAEs: Tendency to generate blurry images (believed to be due to injected noise and weak inference models), recently better with inverse autoregressive flow
B-VAE: learn disentangled representation without supervision. For $\beta = 1$ same as VAE
 $\min: C_B = \mathbb{E}_{z \sim p_{\theta}(z)} [\log p_{\theta}(x|z)] + \beta D_{KL}(q_{\theta}(x|z) || p_{\theta}(x))$

AUTOREGRESSIVE MODELS

Autoregression: for timeseries, $x_i = b_0 + b_1 x_{i-1} + \dots + b_n x_n$ taking observations from previous timesteps as input.
 E.g. sequence models
Tabular approach: $p(x) = \prod_i p(x_i | x_1, \dots, x_{i-1})$ #params 2^{m+1} can also specify $p(x_i | x_1, \dots, x_{i-1}) = \text{Ber}(f(x_1, \dots, x_{i-1}))$
Fully visible belief networks: $x_i = p(x_i = 1 | x_1, \dots, x_{i-1})$ modeled via logistic regression: $f(x_1, \dots, x_{i-1}) = \sigma(a_0 + \dots + a_{i-1} x_{i-1})$, has $O(n^2)$ parameters

Neural Autoregressive Density Estimator (NADE): autoencoder-like network to learn $p(x_i = 1 | x_1, \dots, x_{i-1})$
 $h^{(i)} = \sigma(b + W_{i-1} X_{i-1})$, $x_i = \sigma(C_i + V_i h^{(i)})$ for binary data
 (we can use that $(b + W_{i-1} X_{i-1}) - (b + W_{i-1} X_{i-1}) = U_{i-1} X_{i-1}$)
 Train by max. $\frac{1}{T} \sum_{t=1}^T \log(p(x_t)) = \frac{1}{T} \sum_{t=1}^T \sum_{i=1}^T \log p(x_i^{(t)} | x_{i-1}^{(t)})$
Advantages: efficient computations in O(TD)

could make use of second order optimizers, easily extendable to other types of observations (reals, multinomials). Random input order is fine.
 Training: take data from training, Inference: Sampling
Extensions: RNADE: conditionals modelled by mixture of Gaussians, DeepNADE: a DNN trained to assign a cond. distr. to any given var. given any subset of others (order-less), ConvNADE

Masked Autoencoder Distribution Estimator (MADE): Constrain autoencoder s.t. output can be used as conditional $p(x_i | x_{<i})$. To fulfill autoregressive property: no computational path output unit x_d and any of the input units x_1, \dots, x_d must exist. Can also be seen as masks $m^{(i)}$: 1 if unit i of layer l is connected to unit j of layer $l-1$. Train with NLL of binary x . Computing $p(x)$ is just forward pass, sampling requires D forward passes. In practice very large hidden layers necessary.

Generative Models for Natural Images: Use an explicit density model & use chain rule to decompose likelihood of image x into prod. of 1D dists. $p(x) = \prod_{i=1}^n p(x_i | x_{<i}) \rightarrow$ need to define ordering of pixels, complex distributions \rightarrow param. via NN. Then max. LL.

PixelRNN: Generate pixel starting from the corner and dependency to prev. pixels modelled using LSTM
 Issue: sequential gen is slow due to explicit pixel dependencies

Pixel CNN: Still generate starting from corner, but model the dependency as CNN:  (masked convolutions)
 $p(x_i | x_{<i}) = p(x_{i,1} | x_{<i}) \cdot p(x_{i,2} | x_{i,1}, x_{<i}) \cdot p(x_{i,3} | x_{i,2}, x_{i,1}, x_{<i})$
 Stacks can be used to speed it up. Add information between vertical and horizontal stack (preserve pixel dependencies)
 $h_{k,k} = \tanh(W_{k,l} * h_k) \odot (W_{k+1,l} * h_k)$. Training is faster than PixelRNN (can parallelize convolutions) but still slow (sequential generation).

Advantages: explicit likelihood, likelihood of training data is good eval. metric, good samples.

Dilated convolution: Receptive field size  can increase exponentially with dilated layers.

WaveNet: Temporal convolutional networks uses dilated convolution. Cannot use strided conv, due to need to preserve resolution

Variational RNN (VRNN): RNNs are autoregressive models. Increase expressive power of RNNs by incorporating stochastic latent variables into hidden state of RNN. Include one VAE per timestep. Including a dynamic prior explicitly models temporal dependencies between time steps. Can include conditional (e.g. style and content of handwriting).

Stochastic Temporal Convolutional Networks (STCN): Combines computational advantages of TCNs with expressiveness of stochastic RNNs: integrate a hierarchical VAE with a TCN.

Self Attention and Transformers: Transforms input Embeddings to Keys, Values, Queries $K = XW_k \in \mathbb{R}^{T \times D}$, $V = XW_v \in \mathbb{R}^{T \times D}$, $Q = XW_q \in \mathbb{R}^{T \times D}$. Attention weights: $\alpha = \text{softmax}(QK^T / \sqrt{D} + \gamma) \in \mathbb{R}^{T \times T}$ where $\gamma = \begin{pmatrix} \gamma_1 & \gamma_2 & \dots & \gamma_T \end{pmatrix}$ masks out future influences
 Prediction: $x_t = \alpha V$

NORMALIZING FLOWS

Have analytic model density $p_{\theta}(x)$ and latent variable space

$$z \xleftarrow{f^{-1}} x = f_0(z) \quad z \xleftarrow{f_0} x$$

where f_0 is deterministic & invertible, preserves dimensionality, Jacobian is efficient (block diagonal)

Maps one distribution to another

$$\text{For a transformation } f(z) = x \\ p_x(x) = p_z(f^{-1}(x)) \cdot |\det(\frac{\partial f^{-1}(x)}{\partial x})| \\ = p_z(f^{-1}(x)) \cdot |\det(\frac{\partial f^{-1}(x)}{\partial z})|$$

Coupling Layer

$$\begin{array}{c} \text{forward: } (y^A, y^B) = (h(x^A, \beta(x^A)), x^B) \\ \text{backward: } (x^A, x^B) = h^{-1}(y^A, \beta(y^B)) \\ \frac{\partial (x^A, x^B)}{\partial (y^A, y^B)} = \begin{bmatrix} h' & h' \beta' \\ 0 & 1 \end{bmatrix} \end{array}$$

β can be arbitrary
 h is any elementwise function

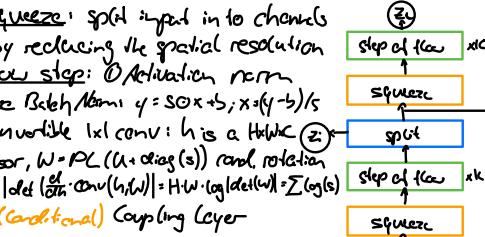
We get complex transforms via composition
 $x = f(z) = f_k \circ f_{k-1} \circ \dots \circ f_1$ (Training) (exact by (tilde) Jacobian)

$$\text{now } \log p_{\theta}(x) = \sum_{i=0}^k (\log p_{\theta}(f_i(z)) + \log |\det(\frac{\partial f_i(z)}{\partial x})|)$$

$$\text{Inference: } p_x(x) = p_z(f^{-1}(x)) \cdot |\det(\frac{\partial f^{-1}(x)}{\partial z})|$$

\Rightarrow Also conveniently via ODE: $\dot{x} = f_0(x(t), t)$

Multi-Scale Architecture (levels and k steps)



Applications in Computer Vision

Super Resolution: Normalizing flow learns a dist. of HR variants. Uses low resolution image as conditioning

Disentanglement: Conditions on learned attributes

C-Flow: Use two flows, condition one on the other. Each flow handles one representation, conditioning goes transfer

Θ Training is still slow

GENERATIVE ADVERSERIAL NETWORKS

GAN Intuition: Want to sample high dimensional data points from true data distribution (intractable), so draw from simple distribution and use a neural network to transform into output. Use generator $G: \mathbb{R} \rightarrow X$ and discriminator $D: X \rightarrow [0, 1]$

GAN Objective: Consider only the discriminator
 $L(D) = -\frac{1}{N} (\sum_{i=1}^N y_i \log(D(x^{(i)})) + \sum_{i=1}^N (1-y_i) \log(1-D(x^{(i)})))$
 for GANs we know 50% are $y=1$ and 50% are $y=0$. Our goal is to train D & G jointly:

$$= \frac{1}{2} (\mathbb{E}_{x \sim p_d} [\log D(x)] + \mathbb{E}_{z \sim p_g} [\log(1 - D(G(z)))])$$

$$= \frac{1}{2} (\mathbb{E}_{x \sim p_d} [\log D(x)] + \mathbb{E}_{x \sim p_g} [\log(1 - D(x))])$$

\Rightarrow a good G maximizes loss function of D.

$$\text{Value function: } V(G, D) = \mathbb{E}_{x \sim p_d} [\log(D(x))] + \mathbb{E}_{x \sim p_g} [\log(1 - D(x))]$$

$$\Rightarrow G^* \text{ will fool any } D: \quad G^* = \arg \max_{G} V(G, D)$$

Optimum of $V(G, D)$ if $p_d = p_g$: Let $D^* = \arg \max_D V(G, D)$

$$\text{If } p_d = p_g \text{ then } D^* = \frac{p_d}{p_d + p_g}, \text{ Proof: } V(G, D) = \int_x p_d(x) \log(D(x)) + p_g(x) \log(1 - D(x)) dx$$

$$\text{Optimal } D: \quad \text{Val. } D = \frac{p_d}{p_d + p_g} \Rightarrow y = \frac{p_d}{p_d + p_g} \cdot \frac{1}{1 - p_d} \cdot \frac{p_d}{p_d + p_g} \cdot \frac{1}{1 - p_d} \Rightarrow D = \frac{p_d}{p_d + p_g}$$

$$\text{Obj. min JS: } V(D^*, G) = \mathbb{E}_{x \sim p_d} [\log(\frac{p_d}{p_d + p_g})] + \mathbb{E}_{x \sim p_g} [\log(\frac{p_g}{p_d + p_g})]$$

Global min is achieved if $p_d = p_g$ and at

$$\min V(\cdot, \cdot) = -\log 4. \text{ Proof: } \mathbb{E}_{x \sim p_d} [\log \frac{p_d}{p_d + p_g}] + \mathbb{E}_{x \sim p_g} [\log \frac{p_g}{p_d + p_g}]$$

$$= -\log(2) + \mathbb{E}[\log(\frac{p_d}{p_d + p_g})] - \log(2) + \mathbb{E}[\log(\frac{p_g}{p_d + p_g})]$$

$$= -\log(4) + D_{KL}(p_d || \frac{p_d}{p_d + p_g}) + D_{KL}(p_g || \frac{p_g}{p_d + p_g})$$

$$= -\log(4) + 2 \cdot D_{JS}(p_d || p_g) \quad \text{if } p_d = p_g \quad \text{then } = 0$$

Note: last term ≥ 0

Jensen-Shannon Divergence: symmetric DKL

$$D_{JS}(p || q) = \frac{1}{2} D_{KL}(p || \frac{p+q}{2}) + \frac{1}{2} D_{KL}(q || \frac{p+q}{2})$$

$$\text{Obj: } \min_{\theta} \max_{\theta} \mathbb{E}_{x \sim p_d} [\log(D(x))] + \mathbb{E}_{x \sim p_g} [\log(1 - D(G(x)))]$$

GAN Training: min _{θ} max _{θ} $\mathbb{E}_{x \sim p_d} [\log(D(x))] + \mathbb{E}_{x \sim p_g} [\log(1 - D(G(x)))]$

Alternate: 1) max _{θ} $\mathbb{E}_{x \sim p_d} [\log(D(x))] + \mathbb{E}_{x \sim p_g} [\log(1 - D_G(x))]$ (both GD)
 2) min _{θ} , $\mathbb{E}_{x \sim p_d} [\log(1 - D_G(x))]$

GAN Training Problems: In this problematic region G produces bad samples and gradient is very flat. Can lead to overtraining of the discriminator (G cannot catch up). Solution: In 2. gradient ascend: max _{θ} $\mathbb{E}_{x \sim p_d} [\log(D_G(x))]$. Generally: need to find Nash-equilibrium in two-player game. This leads to mode collapse (no sample diversity). Happens if there is a saddle in dual energy landscape (the generator learns only variations of a single sample).

Issues with JS divergence: correlates badly with sample quality (do not know when to stop training). D tends to be optimal, JSD saturates \rightarrow bad gradients
 Solution: use Wasserstein distance

GAN Pseudocode

for #training iterations do

for k steps do

sample minibatch of m noise samples $\{z \sim p_g(z)\}$
 sample minibatch of examples $\{x\}$ from data gen. dist.
 update discriminator by ascending on its stochastic grad.
 $\nabla_{\theta} \frac{1}{m} \sum_{i=1}^m [\log(D(x^{(i)}) + \log(1 - D(G(z^{(i)})))]$

sample minibatch of m noise samples $\{z \sim p_g(z)\}$
 update the generator (ascent): $\nabla_{\theta} \frac{1}{m} \sum_{i=1}^m \log(D(G(z^{(i)}))$

GAN: Comparison to VAE: No variational bound is needed, only samples. Generally sharper images. GANs are asymptotically consistent (not proven for VAE)

DC-GAN: G is an up-sampling network with deconvolutions, D is a CNN. Replace pooling by strided convolutions in D and strided deconvolutions in G. Use batch norm for both. Remove fully connected hidden layers, use ReLU in G except tanh for output, use LeakyReLU in D for all layers. Can do "vector math": smiling φ - neutral φ + neutral θ = smiling θ
 Image translation
Pix2Pix: $L(G, D) = L_{DCGAN}(G, D) + \lambda L_{L1}(G)$ where $L_{DCGAN}(G, D) = \mathbb{E}_{x,y} [\log D(x, y)] + \mathbb{E}_{x,z} [1 - \log D(x, G(x, z))]$ and $L_{L1}(G) = \mathbb{E}_{x,y,z} [\|y - G(x, z)\|_1]$

