

## IE4497 formula

$$\text{Gaussian Dist: } p(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$$

If  $X \sim \mathcal{N}(\cdot|\mu, \sigma^2)$ , then

$$aX \sim \mathcal{N}(\cdot|a\mu, a^2\sigma^2); X + c \sim \mathcal{N}(\cdot|\mu + c, \sigma^2); \text{ let } Z \sim \mathcal{N}(\cdot|\mu, \sigma^2), \text{ then } X = \sigma Z + \mu \sim \mathcal{N}(\cdot|\mu, \sigma^2)$$

if  $X$  and  $Y \sim \mathcal{N}(\cdot|0, v^2)$  are independent, then  $X + Y \sim \mathcal{N}(\cdot|0, \sigma^2 + v^2)$

$$\text{If } X \sim \mathcal{N}(\cdot|\mu, \sigma^2) \text{ Joint pdf: } P((X, Y) \in A) = \int_A p(x, y) dx dy, \quad p(x, y) = p(x)p(y|x) \quad (= p(x)p(y) \text{ if independent})$$

$$p(x) = \int p(x, y) dy (\text{marginal distribution}), \quad p(x|y) = \frac{p(x, y)}{p(y)} = \frac{p(y|x)p(x)}{p(y)} (\text{Bayes' Theorem})$$

$$\frac{p(a|b, c)}{p(a|c)} = \frac{p(b|a, c)}{p(b|c)}, \quad E[X] = \int xp(x) dx,$$

$$E[g(X)] = \int g(x)p(x) dx, \quad E[X + Y] = E[X] + E[Y], \quad E[X|Y = y] = \int xp(x|y) dx$$

$$E[E[X|Y]] = E[X], \quad E[f(X)g(Y)|Y] = E[f(X)|Y]g(Y)$$

$$\text{var}(X) = E[(X - E[X])^2] = E[X^2] - E[X]^2, \quad \text{cov}(X, Y) = E[(X - E[X])(Y - E[Y])]$$

$$\text{var}(X + Y) = \text{var}(X) + \text{var}(Y) + 2\text{cov}(X, Y). \quad \text{If } X = AZ + \mu : E[x] = \mu, \text{ cov}(X) = AA^T$$

$$E[AX] = AE[x], \quad \text{cov}(x) = \sum_{xx} = E[(x - E[x])(x - E[x])^T], \quad \text{cov}(Ax) = A \text{cov}(x) A^T$$

$$\text{cov}(x, y) = \sum_{xy} = E[(x - E[x])(y - E[y])^T], \quad \frac{\delta(a^T x)}{\delta x} = a, \quad \frac{\delta(x^T A x)}{\delta(x)} = (A + A^T)x, \quad \frac{\delta a^T X b}{\delta X} = ab^T$$

$$\frac{\delta \det(X)}{\delta X} = \det(X)(X^{-1})^T, \quad \frac{\delta a^T X^{-1} b}{\delta X} = -(X^{-1})^T ab^T (X^{-1})^T$$

$$\underbrace{p(\theta|x)}_{\text{posterior}} = \frac{p(x, \theta)}{p(x)} = \frac{p(x|\theta)p(\theta)}{p(x)} \propto \underbrace{p(x|\theta)}_{\text{likelihood}} \cdot \underbrace{p(\theta)}_{\text{prior}}$$

$$\text{Residual sum of squares } RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2, \quad \text{Root mean squared error } RMSE = \sqrt{\frac{1}{n} RSS}.$$

$$\text{Coefficient of determination } R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{RSS}{\sum_{i=1}^n (y_i - \bar{y})^2},$$

$$\text{where } \bar{y} = \frac{1}{n} \sum_i y_i. \quad \text{TSS is the total sum of squares or the empirical variance of the data.}$$

$$\text{- Binomial distribution: } \text{Bin}(x|\theta, n) = \binom{n}{x} \theta^x (1 - \theta)^{n-x}, \quad \binom{n}{x} = \frac{n!}{x!(n-x)!}$$

$$\text{- Beta distribution: } \text{Beta}(x|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1} \quad (*\text{Beta dist is a conjugate prior to Bin dist})$$

$$\text{- Categorical distribution: } \text{Cat}(x|\theta_1, \dots, \theta_K) = \theta_x$$

$$\text{- Multinomial distribution: } \text{Mult}(x|\theta, n) = \binom{n}{x_1, \dots, x_K} \prod_{k=1}^K \theta_k^{x_k}, \quad \binom{n}{x_1, \dots, x_K} = \frac{n!}{x_1! \dots x_K!}.$$

$$\text{- Dirichlet distribution: } \text{Dir}(x|\alpha_1, \dots, \alpha_K) = \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K x_k^{\alpha_k-1}, \quad (*\text{Dirich dist is a conjugate prior to Multi Nor dist})$$

$$\text{- Joint Gaussian distribution: } N(x|\mu, \Sigma) = \frac{1}{(2\pi)^{K/2} (\det \Sigma)^{1/2}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

$$\text{- Uniform distribution: } \text{Unif}(x|a, b) = \begin{cases} \frac{1}{b-a} & \text{if } a \leq x \leq b, \\ 0 & \text{otherwise.} \end{cases}, \quad \text{- Exponential distribution: } \text{Exp}(x|\lambda) = \begin{cases} \lambda e^{-\lambda x} & \text{if } x \geq 0, \\ 0 & \text{otherwise.} \end{cases}$$

$$\text{Bernoulli distribution: } p(D|\theta) = \theta^{N_1} (1 - \theta)^{N_0}, \quad N_k = \sum_{i=1}^n 1\{x_i = k\}. \quad \text{Exponential distribution: } p(D|\lambda) = \lambda^n \exp(-\lambda \sum_{i=1}^n x_i)$$

$$\text{Normal distribution: } p(D|\theta) = N(\mu, \sigma^2)^n = \frac{1}{(2\pi)^{n/2} \sigma^n} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right)$$

$$\text{Linear model: } y = w^\top x + \epsilon, \quad \text{Basis function expansion: } y = w^\top \phi(x) + \epsilon, \quad \text{MLE of linear model: } w_{ML} = (\Phi^\top \Phi)^{-1} \Phi^\top y$$

$$\text{Kullback-Leibler divergence: } D(p_{\theta_0} \| p_{\theta}) = E_{\theta_0} \left[ \log \frac{p(x|\theta_0)}{p(x|\theta)} \right], \quad \text{MAP estimate: } \theta_{MAP} = \max_{\theta} p(\theta|D)$$

$$\text{MMSE estimate: } \hat{\theta}(D) = \min_a E[(\theta - a)^2 | D] = E[\theta | D]$$

$$\text{Mixture model: } p(x|\theta) = \sum_{k=1}^K \pi[k] p_k(x|\theta), \text{ GMM: } p(x|\theta) = \sum_{k=1}^K \pi[k] \mathcal{N}(x|\mu_k, \Sigma_k), \quad \theta = (\pi[k], \mu_k, \Sigma_k)_{k=1}^K$$

$$\text{E step: } Q(\theta|\theta^{(m)}) = E_{y \sim p(\cdot|x, \theta^{(m)})} [\log p(y|\theta)|x, \theta^{(m)}], \text{ M step: } \theta^{(m+1)} = \arg \max_{\theta \in \Theta} Q(\theta|\theta^{(m)})$$

$$\text{GMM E step: } r_{ik}^{(m)} = p(z_i = k|x_i, \theta^{(m)}) = \frac{\pi^{(m)}[k] \mathcal{N}(x_i|\mu_k^{(m)}, \Sigma_k^{(m)})}{\sum_{k'} \pi^{(m)}[k'] \mathcal{N}(x_i|\mu_{k'}^{(m)}, \Sigma_{k'}^{(m)})}, n_k^{(m)} = \sum_{i=1}^n r_{ik}^{(m)},$$

$$Q(\theta|\theta^{(m)}) = -\frac{1}{2\sigma^2} \sum_{i=1}^n \|x_i - \mu_{k_i}^{(m)}\|^2 + \text{const.}$$

$$\text{GMM M step: } \pi^{(m+1)}[k] = n_k^{(m)}/n, \mu_k^{(m+1)} = (1/n_k^{(m)}) \sum_{i=1}^n r_{ik}^{(m)} x_i; \quad \Sigma_k^{(m+1)} = (1/n_k^{(m)}) \sum_{i=1}^n r_{ik}^{(m)} (x_i - \mu_k^{(m+1)})(x_i - \mu_k^{(m+1)})^\top.$$

$$\text{MAP estimate: } \theta^{\text{MAP}} = \arg \max_{\theta \in \Theta} (\log p(x|\theta) + \log p(\theta)). \quad \text{EM for MAP: } \text{Estep: } Q(\theta|\theta^{(m)}) = E_{y \sim p(\cdot|x, \theta^{(m)})} [\log p(y|\theta)|x, \theta^{(m)}]$$

$$\text{Mstep: } \theta^{(m+1)} = \arg \max_{\theta \in \Theta} (Q(\theta|\theta^{(m)}) + \log p(\theta))$$

$$\text{Markov property: } p(x_t|x_1, \dots, x_{t-1}) = p(x_t|x_{t-1}), \text{ Transition probability: } T(i, j) = p_{x_t|x_{t-1}}(j|i)$$

$$\text{Unigram model: } p(x_t = x), \text{ Bigram model: } p(x_t|x_{t-1}), \text{ n-gram model: } p(x_t|x_{t-1}, x_{t-2}, \dots, x_{t-n+1})$$

$$\text{PageRank score: } \pi_i = \sum_j T(j, i) \pi_j$$

$$\text{MLE for Markov model: } \log p(D|\pi, T) = \sum_{i=1}^n \log \pi(x_i, 0) + \sum_{i=1}^n \sum_{t=1}^{t_i} \log T(x_{i,t-1}, x_{i,t})$$

$$= \sum_{x=1}^M N_x \log \pi(x) + \sum_{x=1}^M \sum_{y=1}^M N_{xy} \log T(x, y)$$

$$N_x = \sum_{i=1}^n \mathbb{I}\{x_{i,0} = x\}, N_{xy} = \sum_{i=1}^n \sum_{t=1}^{t_i} \mathbb{I}\{x_{i,t-1} = x, x_{i,t} = y\}, \hat{\pi}(x) = \frac{N_x}{n}, \hat{T}(x, y) = \frac{N_{xy}}{\sum_z N_{xz}},$$

$$\text{HMM: } p(x_0, \dots, x_T, z_0, \dots, z_T|\theta) = \pi(z_0) p(x_0|z_0) \prod_{t=1}^T T(z_{t-1}, z_t) p(x_t|z_t)$$

$$\text{Baum-Welch algorithm: MLE - } \log p(D|\theta) = \sum_{i=1}^n \log \pi(z_{i,0}) + \sum_{i=1}^n \sum_{t=1}^{t_i} \log T(z_{i,t-1}, z_{i,t}) + \sum_{i=1}^n \sum_{t=0}^{t_i} \log p(x_{i,t}|\phi_{z_{i,t}})$$

$$\text{Steps for BW Algo: (1) Initialize } \theta^{(0)}. \text{ (2) E step: At iteration } m, \text{ use Forward-Backward Algorithm to compute}$$

$$\gamma_{i,t}(z) = p(z_{i,t} = z|x_{i,\cdot}, \theta^{(m)}) \propto \alpha_j(z) \beta_j(z), \quad \xi_{i,t}(z, z') = p(z_{i,t-1} = z, z_{i,t} = z'|x_{i,\cdot}, \theta^{(m)}) \\ \propto \alpha_{t-1}(z) p(x_{i,t}|z_{i,t} = z') \beta_t(z') p(z_{i,t} = z'|z_{i,t-1} = z). \quad \text{(3) M step: Find } (m+1).$$

$$\pi^*(z) = \frac{\sum_{i=1}^n \gamma_{i,0}(z)}{n}, \quad T^*(z, z') = \frac{\sum_{i=1}^n \sum_{t=1}^{t_i} \xi_{i,t}(z, z')}{\sum_u \sum_{i=1}^n \sum_{t=1}^{t_i} \xi_{i,t}(z, u)}, \quad \hat{\phi}_z = \text{emission prob. model parameters.}$$

$$\text{Viterbi algorithm: } z_0^*, \dots, z_T^* = \arg \max_{z_0, \dots, z_T} p(z_0, \dots, z_T|x_0, \dots, x_T)$$

$$\text{Cumulative distribution function: } P(X \leq x) = F(x), X = F^{-1}(U), \text{ where } U \sim \text{Unif}([0, 1]), u \leq F(x) \Leftrightarrow F^{-1}(u) \leq x$$

$$\text{Transformation method: } p_Y(y) = \sum_{k=1}^K \frac{p_X(x_k)}{|f'(x_k)|}, \text{ where } x_1, x_2, \dots, x_K \text{ are solutions to } f(x) = y$$

$$\text{Rejection sampling: } p(z) = \frac{1}{M} \tilde{p}(z), \text{ where } M \text{ is unknown, } kq(z) \geq \tilde{p}(z) \text{ for all } z$$

$$\text{Accept } z \sim q(z) \text{ if } u \sim \text{Unif}([0, kq(z)]) \leq \tilde{p}(z)$$

$$\text{Acceptance probability: } P(z \text{ accepted}) = \int P(z \text{ accepted}|z) q(z) dz = \int \frac{\tilde{p}(z)}{kq(z)} q(z) dz = \frac{M}{k}$$

$$\text{Rejection sampling for Bayesian Inference: } \tilde{p}(\theta) = p(D|\theta)p(\theta) \text{ and } q(\theta) = p(\theta) : k = \max_{\theta} \frac{\tilde{p}(\theta)}{q(\theta)} = \max_{\theta} p(D|\theta)$$

$$\text{Importance sampling: sample } z \text{ where } |f(z)|p(z) \text{ is large for better efficiency rather than from } p(z) \text{ directly.}$$

$$E_p[f(z)] = E_q \left[ \frac{p(z)}{q(z)} f(z) \right], \quad \tilde{w} = \frac{p(z)}{q(z)}, w(z) = \frac{\tilde{w}(z)}{\sum_{i=1}^n \tilde{w}(z_i)}, \quad E_p[f(z)] \approx \sum_{i=1}^n w(z_i) f(z_i)$$

$$\text{Tail sampling: } P(X > a) \approx \sum_{i=1}^n w(z_i), \text{ where } z_1, z_2, \dots \text{ are sampled from } q(z) \text{ with support } (a, \infty) \text{ and } w(z_i) = \frac{p(z_i)}{q(z_i)}$$

$$\text{Sampling importance resampling (SIR): 1. Sample } z_1, \dots, z_n \text{ from } q(z).$$

2. Compute weights  $w(z_1), \dots, w(z_n)$  where  $w(z_i) = \frac{\tilde{w}(z_i)}{\sum_{j=1}^n \tilde{w}(z_j)}$ .

3. Resample with replacement from  $\{z_1, \dots, z_n\}$  according to weights  $(w(z_1), \dots, w(z_n))$ .

SIR for Bayesian inference:  $w(z_i) = \frac{p(D|z_i)}{\sum_{j=1}^n p(D|z_j)}$ ; Sampling for EM:  $Q(\theta|\theta^{(m)}) \approx \frac{1}{n} \sum_{i=1}^n \log p(x, z_i|\theta)$

Stationary distribution:  $\sum_x \pi(x)T(x, y) = \pi(y)$ , Reversible MC:  $\pi(x)T(x, y) = \pi(y)T(y, x)$

Metropolis-Hastings algorithm: (1) Initialize  $x = Z_0$ . (2) For each  $m = 1, 2, \dots$  (3) Sample  $y \sim q(x, y)$ .

(4) Compute acceptance probability  $A(x, y) = \min \left( 1, \frac{\tilde{\pi}(y)q(y, x)}{\tilde{\pi}(x)q(x, y)} \right)$ .

if  $\pi(x) \propto \psi(x)h(x)$  and  $q(x, y) = h(y)$ , then  $A(x, y) = \min \left( 1, \frac{\psi(y)}{\psi(x)} \right)$ .

(5) With probability  $A(x, y)$ , set  $Z_m = y$ ; otherwise set  $Z_m = x$  (6) Update  $x = Z_m$ .

- random walk MH:  $q(x, y) = q(y - x)$ .  $y - x \sim \mathcal{N}(\cdot|0, \Sigma)$  - Gaussian centered at  $x$ .

$y - x \sim \text{Unif}[-\delta, \delta]^d$  - Uniform distribution centered at  $x$ .

Gibbs sampling:  $p(z_i|z_{-i}) = \frac{p(z_1, \dots, z_d)}{p(z_{-i})}$ . For each  $i$ , let  $z_{-i} = \{z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_d\}$ , i.e.,  $z_i$  removed.

(1) Initialize  $(z_1^{(0)}, \dots, z_d^{(0)})$ . (2) For each  $k$ : (3) sample  $z_1^{(k)}$  from  $p(\cdot|z_2^{(k-1)}, \dots, z_d^{(k-1)})$

(4) sample  $z_2^{(k)}$  from  $p(\cdot|z_1^{(k)}, z_3^{(k-1)}, \dots, z_d^{(k-1)})$  ... (5) sample  $z_j^{(k)}$  from  $p(\cdot|z_1^{(k)}, \dots, z_{j-1}^{(k)}, z_{j+1}^{(k-1)}, \dots, z_d^{(k-1)})$  ...

(6) sample  $z_d^{(k)}$  from  $p(\cdot|z_1^{(k)}, \dots, z_{d-1}^{(k)})$

Ising model: Given a noisy image  $y$ , want to recover  $z$ . compute the posterior  $p(z | y)$ .

$$\text{Likelihood} : p(y | z) = \prod_j p(y_j | z_j) = \prod_j \mathcal{N}(y_j | z_j, \sigma^2)$$

$$p(z_j | z_{-j}) \propto \prod_{s \in N_j} \psi(z_s, z_j), \text{ where } N_j \text{ is the neighborhood of pixel } z_j, \text{ and}$$

$$\psi(u, v) = \exp(Juv) \text{ with } J > 0 \text{ as the "coupling strength".}$$

$$(N + 2P - F) / \text{stride} + 1$$

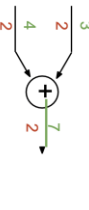
Input volume: **32x32x3**

**10 5x5** filters with stride 1, pad 2

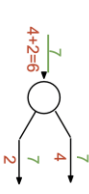
Number of parameters in this layer?

each filter has **5\*5\*3 + 1 = 76** params  
 $\Rightarrow 76 * 10 = 760$

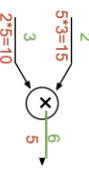
add gate: gradient distributor



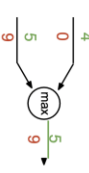
copy gate: gradient adder



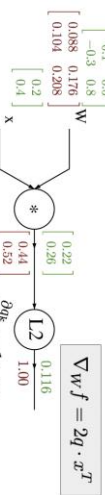
mul gate: "swap multiplier"



max gate: gradient router



A vectorized example:  $f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W \cdot x)_i^2$



$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \dots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \dots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \dots + q_n^2$$

Activation Functions

**Sigmoid**

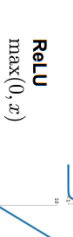
$$\sigma(x) = \frac{1}{1+e^{-x}}$$

**tanh**

$$\tanh(x)$$

**ReLU**

$$\max(0, x)$$



**Leaky ReLU**

$$\max(0, 1.02x)$$

**Maxout**

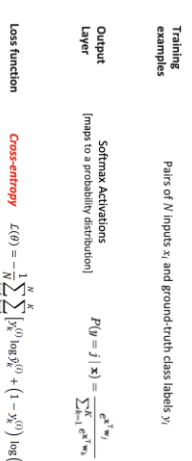
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

**ELU**

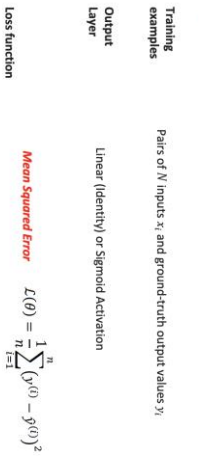
$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



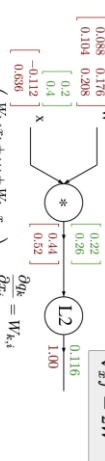
Classification tasks



Regression tasks



A vectorized example:  $f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W \cdot x)_i^2$



$$q = W \cdot x = \begin{pmatrix} W_{1,1}x_1 + \dots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \dots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \dots + q_n^2$$

filtering

smoothing

fixed-lag smoothing

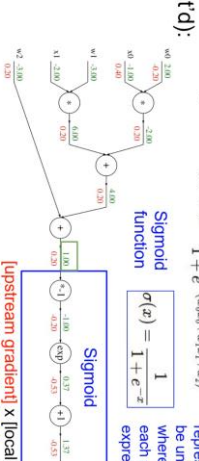
prediction



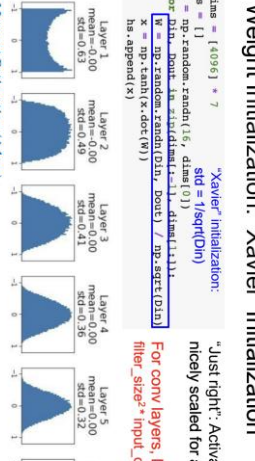
Batch Normalization



Another example



Weight Initialization: "Xavier" Initialization



AlexNet showed that you can use CNNs to train Computer Vision models.

ZFNet, VGG shows that bigger networks work better

GoogLeNet is one of the first to focus on efficiency using 1x1 bottleneck convolutions and global avg pool instead of FC layers

ResNet showed us how to train extremely deep networks

- Limited only by GPU & memory!

- Showed diminishing returns as networks got bigger

After ResNet: CNNs were better than the human metric and focus shifted to Efficient networks:

- Lots of tiny networks aimed at mobile devices: MobileNet, ShuffleNet

Neural Architecture Search can now automate architecture design

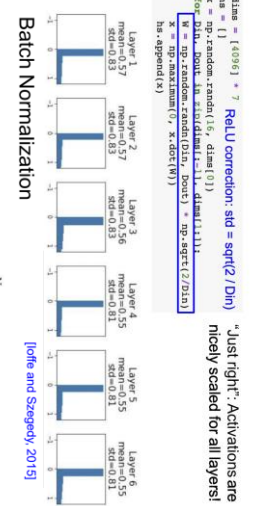
Computational graph



Backprop with Matrices

$$\frac{\partial L}{\partial x} = \left( \frac{\partial L}{\partial y} \right) w^T$$

Weight Initialization: He Initialization



Input:  $x: N \times D$

Learnable scale and shift parameters:

$$\gamma, \beta: D$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sigma_j}$$

$$\beta_j = \mu_j$$

$$\gamma_j = \sigma_j$$

Normalized  $x$ , Shape is  $N \times D$

Output, Shape is  $N \times D$