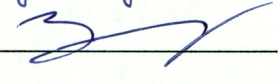


Optimize
 $L_1 \Leftrightarrow$ Laplacian error
 $L_2 \Leftrightarrow$ Gaussian.

Examination Aid Sheet
Faculty of Applied Science & Engineering

Both sides of the sheet may be used;
 must be printed on 8.5" x 11" paper.

Notation: S : sample space e.g. $S = \{1,2,3\}$, $S = \mathbb{R}$
 E : Event: $E \subseteq S$ e.g. $E = \{1,2\}$, $E = [0,1)$
 X : Random Variable: $X: S \rightarrow \mathbb{R}$
 $X=6 \Rightarrow$ eq. value S st $X(S)=6$ is the solⁿ
 P : Probability: $P(\emptyset) = 0, P(E) \geq 0$

Subject: ECE52
 Candidate's name: King Fung
 Candidate's signature: 

Definitions
 CDF: $\int_{-\infty}^x p(x) dx$
 $CDF(x) = 1$
 $E[X] = \int_{-\infty}^{\infty} x P(x) dx$
 $Var[X] = \int_{-\infty}^{\infty} (x - E[X])^2 P(x) dx$

Properties: $\int_{-\infty}^{\infty} P(x) dx = \int_{-\infty}^{\infty} P_c(x) dx = 1$ (PDE)
 $= E[X^2] - E[X]^2$
 $Var[N] = \sigma^2$

Function: $N(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\{-\frac{(x-\mu)^2}{2\sigma^2}\}$
 Beta(α, β): $\frac{x^{\alpha-1} (1-x)^{\beta-1}}{B(\alpha, \beta)}$ $x \in [0,1)$

Joint Probability: $P(x,y)$
 Marginal $P(x), P(y) \Rightarrow P(x) = \int_{-\infty}^{\infty} P(x,y) dy$ Conditional
 $P(x|y) = \frac{P(x,y)}{P(y)}$
 Bayes Theorem: $P(y|x) = \frac{P(x|y)P(y)}{P(x)} = \frac{P(x|y)P(y)}{\sum_x P(x|y)P(y)}$

Total Probability Thm: $P(x) = \sum_{y \in S_y} P(y) P(x|y) = \int_{-\infty}^{\infty} P(y) P(x|y) dy$

Models: kNN, GPR, Lin/Log reg, Neural Learning:
 for some passes over dataset (epochs):
 for some iterations:
 refine search dir
 more in search dir
 best refinement:
 steepest descent: $W \leftarrow W - \eta \frac{\partial L}{\partial W}$ Convexity
 \rightarrow Local min = global min.
 \rightarrow obeys Jensen's inequality: $f(\alpha W_1 + (1-\alpha)W_2) \leq \alpha f(W_1) + (1-\alpha)f(W_2)$

Loss fun: $L_p = \|y - \hat{y}\|_p = \sum_{i=1}^n (y_i - \hat{y}_i)^p$
 $p=1$: Taxicab $p=2$: Euclidean

kNN: Given M data pts $\{(x^{(m)}, t^{(m)})\}$ find the nearest k data pts $L_2(x^* - x^{(k)})$ output the avg/majority regressed class

Momentum (Christal) $\Delta W_{iM_0} = -\eta \frac{\partial L(t)}{\partial W}$
 $\Delta W_{CM_0}(t) = -\eta \frac{\partial L(t)}{\partial W} + \alpha \Delta W(t-1)$ Nesterov's $\alpha \frac{\partial L(t-1)}{\partial W} \rightarrow -\eta \frac{\partial L(t)}{\partial W}$
 $\Delta W_{NM_0} = -\eta \frac{\partial L(t)}{\partial W} + \alpha \Delta W(t-1)$

Stochastic (kMini batch) Gradient Descent:
 $L(w)$ expensive w.r.t $M \uparrow \uparrow$ Choose M i.i.d samples:
 $\frac{\partial L(w)}{\partial W} \approx \frac{1}{M} \sum_{m=1}^M \frac{\partial L_m(x^{(m)}, y^{(m)}, w)}{\partial W}$
 $\frac{\partial L(w)}{\partial W} \approx \frac{1}{M} \sum_{m=1}^M \frac{\partial L_m(x^{(m)}, y^{(m)}, w)}{\partial W}$

MLE VS MAP
 $\frac{\partial P}{\partial \text{param}} = 0$ $P(\text{param}|\text{data}) = \frac{P(\text{data}|\text{param})P(\text{param})}{P(\text{data})}$
 or $\frac{\partial \log P}{\partial \text{param}} = 0$ "requires a prior"

MLE for $N(x|\mu, \sigma^2)$:
 $\ln N(x|\mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln \pi$
 $\frac{d \ln N}{d \mu} = \frac{1}{\sigma^2} \sum_{n=1}^N x_n - N = 0 \Rightarrow \mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n$

Multivariate N_0 : $X = (x_1, \dots, x_N)^T \rightarrow$ i.i.d, gaussian
 $N_D(x|\mu, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\}$
 Odds $D \times D$ $E[X] = \mu$
 $COV[X] = \Sigma$

$\frac{d \ln N}{d \sigma^2} = \frac{1}{2(\sigma^2)^2} \sum_{n=1}^N (x_n - \mu_{ML})^2 - \frac{N}{2\sigma^2} = 0 \Rightarrow \sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2$

MLE of N_D : $\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n$ $\Sigma_{ML} = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})(x_n - \mu_{ML})^T$
 $LL: \ln N_D = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N (x_n - \mu)^T \Sigma^{-1} (x_n - \mu)$

Precision Matrix & Joint & Disj: $X \Rightarrow (x_a, x_b)$
 $\Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}$ $\mu \Rightarrow (\mu_a, \mu_b)$

Least Squares Error: $\frac{1}{2} \sum_{i=1}^n (x_i^T w - t_i)^2$
 Learn a Param Model $\text{Max}_w \prod_{n=1}^M P(x^{(n)}, y^{(n)}|w)$ $\Rightarrow \text{Max}_w \log P$
 $\Rightarrow \text{argmin}_w -\log P$

data $\Lambda = \Sigma^{-1}$ $\Lambda = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}$
 COND: $P(x_a|x_b) = N(x_a|\mu_a, \Sigma_{a|b})$
 $\Sigma_{a|b} = \Lambda^{-1} = \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}$

Regularizing $E = E_D(w) + E_w(w) \lambda$
 $= \frac{1}{2} L_2 + \frac{\lambda}{2} \|w\|_2^2$
 $q=1$: lasso \Rightarrow sparse
 $q=2$: quad \Rightarrow std.
 Dist: $P(x,y) = P(y|x, \mu) P(x)$
 $\text{Gen } P(x,y|w) = P(y|x, w) P(x)$
 Dist \Rightarrow outputs $\text{Gen} \Rightarrow$ inputs find
 The MAP \Rightarrow L_2 regularization
 MLE of gaussian: $w = (X^T X)^{-1} X^T Y$
 $\rightarrow \text{argmin}_{w,b} -\log P$
 MAP: Assume $w \in N(w_0, \sigma^2)$
 $\rightarrow \log P(w) = -\sum_{n=1}^N \frac{\lambda}{2} \|w_n\|_2^2 + c$

Parametric Dist: $P(x,y|w)$: likelihood of x,y given w .
 Marg: $p(x_a) = N(x_a|\mu_a, \Sigma_{aa})$

Loss Fun: $L_2: E[L_2] = \int \int (y - g(x))^2 P(x,y) dx dy$
 \rightarrow min w.r.t $g(x) = E[y|x]$
 Min kovsky $L_q: E[L_q] = \int \|y - g(x)\|_q^q P(x,y) dx dy$
 min im: $q=0 \Rightarrow$ mode $q=2$ mean $q=1$ median ($q=2 \checkmark$ noise)

Bias/Variance tradeoff:
 ↳ regularize: $\lambda \uparrow$ bias \uparrow
 $\lambda \downarrow$ var \uparrow

$$E[L_2] = \int (\hat{y}(x) - E[y|x])^2 P(x) dx$$

$$= \int (E[\hat{y}(x)] - E[y|x])^2 P(x) dx + \int (\hat{y}(x) - E[\hat{y}(x)])^2 P(x) dx \Rightarrow \text{complex eqn.}$$

Linear Basis Funct Model

$$y(x, w) = \sum_{j=0}^n w_j \phi_j(x)$$

e.g. $\phi_0(x) = x^j$
 $\phi_j(x) = \exp(-\frac{(x-\mu)^2}{2\sigma^2})$
 $\phi_j(x) = \sigma \frac{(x-\mu)^j}{j!}$
 $\phi_j(x) = \sigma \frac{(x-\mu)^j}{j!}$
 $\phi_j(x) = \frac{1}{1 + \exp(-x)}$

iteration: $w^{(i+1)} = w^{(i)} - \eta \nabla E = w^{(i)} - \eta (\sum_{n=1}^N \phi(x_n) \phi(x_n)^T) \phi(x_n)$

MLE: $w_{MLE} = (\Phi^T \Phi)^{-1} \Phi^T \epsilon$

W/quad penalization: $w^* = (\lambda I + \Phi^T \Phi)^{-1} \Phi^T \epsilon$

Classification: give a sample x_1 & x_2 class, count density

$$P(k|x) = \frac{p(x|c_k) p(c_k)}{\sum p(x|c_k) p(c_k)} = \frac{p(x|c_k) p(c_k)}{p(x)}$$

simple highest prob? see

$P(x|c_1)$ $P(x|c_2)$

\Rightarrow can multiply err cost by cost to weigh e.g. $\begin{bmatrix} 0 & 1000 \\ 1 & 0 \end{bmatrix}$

Decision Theory:

cost of opt equal weight, but not for: e.g. cancer normal

cancer	normal
0	1000
1	0

\Rightarrow Hard to figure out the outcome

Log regression:

W/ L_2 loss, b. $\hat{y}^{(m)} = \sigma(W^T x^{(m)} + b)$

$$\frac{\partial L}{\partial w} = \sum_m (\hat{y}^{(m)} - t^{(m)}) \sum_j y_j^{(m)} (1 - \hat{y}^{(m)}) x_j^{(m)}$$

also, $\frac{\partial \sigma}{\partial z} = \sigma(z)(1 - \sigma(z))$

decision base check: off class: $w \uparrow c_1 (+)$
 soft c_2 , soft $c_2 \Rightarrow w \rightarrow 0$

also $c_1 \Rightarrow$ bad MLE
 \hookrightarrow use MAP to constrain

Other loss: CROSS ENTROPY

KL distance in classification

$$KL(Q||P) = \sum Q(x) \log \frac{Q(x)}{P(x)}$$

↳ if $P(t=k|x) = \hat{y}(x)$, $KL(Q||P) \Rightarrow -e \log \hat{y} - (1-e) \log(1-\hat{y})$
 $\hat{z} = Q(t=k|x) \hat{e}$

Soft max & Multiclass, $t \in \{0, 1, \dots, k\} = \sum_{k=1}^k I(k, t) \log \phi(e^{z_k})$
 $\hookrightarrow P(t=k|x) = \frac{e^{z_k}}{\sum_j e^{z_j}}$ where z stays out.

Neurons

sign: $\begin{cases} z > 0 \\ z < 0 \end{cases} \Rightarrow \frac{\partial \sigma}{\partial z} = \begin{cases} 1 \\ 0 \end{cases}$

$z = \sum w_n x_n + b$

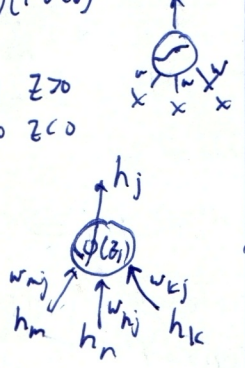
$\sigma(z) = \frac{1}{1 + e^{-z}}$ $\frac{\partial \sigma}{\partial z} = \sigma(z)(1 - \sigma(z))$

$f(z) = z$ $\frac{\partial f}{\partial z} = 1$

ReLU = $\max(0, z)$ $\frac{\partial \text{ReLU}}{\partial z} = \begin{cases} 1 & z > 0 \\ 0 & z < 0 \end{cases}$

tanh \times $1 - \tanh^2 x$

In general: a neuron.



Leary Multilayer:

$$z^1 = w^1 x + b^1 \Rightarrow h^1 = \phi(z^1)$$

$$z^2 = w^2 h^1 + b^2 \Rightarrow h^2 = \phi(z^2)$$

$$z^N = w^N h^{(N-1)} + b^N \Rightarrow \hat{y} = f(z^N)$$

Back Prop:

$$\frac{\partial L}{\partial z^{(n)}} = \frac{\partial y}{\partial z^{(n)}} \frac{\partial L}{\partial y} \frac{\partial z^{(n)}}{\partial w^{(n)}} = \frac{\partial L}{\partial z^{(n+1)}} \frac{\partial z^{(n+1)}}{\partial z^{(n)}}$$

$$\frac{\partial L}{\partial z^{(n)}} = \frac{\partial h^{(n)}}{\partial z^{(n)}} (w^{(n+1)})^T \frac{\partial L}{\partial z^{(n+1)}}$$

$$\frac{\partial L}{\partial w^{(n+1)}} = \frac{\partial L}{\partial z^{(n)}} h^{(n)}$$

Strategies:

Hyperparams: 2-3 loops w/ 500 units
 \rightarrow validate ReLU best practice.

Generalization: wider shall > narrow, do

Init: (identity or $N(x)$) $\frac{\partial \sigma}{\partial z}$

\hookrightarrow better

\hookrightarrow or use pretrain w/ eq Lts

Overfitting:
 \hookrightarrow stop once validation \uparrow

Forward Prop:

$$h_j = \phi(z_j) = \phi(\sum_n w_{nj} h_n + b_j)$$

Back Prop

$$\frac{\partial L}{\partial w_{nj}} = \frac{\partial L}{\partial h_j} \frac{\partial h_j}{\partial z_j} \frac{\partial z_j}{\partial w_{nj}}$$

input $\rightarrow \frac{\partial L}{\partial z_j} h_n$!!!

Dropout only deep neural for full forward pass: $E[w]$ wsh.

CNN: local context, w/ DEPTW

RNN: watch see next for multiple layers & when