

MLPR (Non-Bayesian)

Linear Regression: $f = Xw + b$ or $f(x) = \sum_k w_k \phi_k(x)$

RBF: $\phi(x) = \exp(-(x-c)^T(x-c)/h^2)$

Sigmoid: $\sigma(x) = 1/(1 + \exp(-v^T x - b))$

L2: $\hat{E}(w) = E(w) + \lambda w^T w$; or $y' = \begin{bmatrix} y \\ 0_k \end{bmatrix}$; $\phi' = \begin{bmatrix} \phi \\ \sqrt{\lambda} I_k \end{bmatrix}$

Bias: for μ , ^{var.} scales w/ $1/N$

Gaussian: $p(x) \propto e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$; $p(x) \propto e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$

+ve definite: $z^T \Sigma z \geq 0 \quad \forall z \in \mathbb{R} \neq 0$

1-hot: $y = [0 \ 0 \dots 1 \dots 0]^T$

Naive Bayes:

→ binary $P(x|y=k, \theta) = \prod_d \theta_{d,k}^{x_d} (1-\theta_{d,k})^{1-x_d}$

→ cond. indep. given class

Gaussian Generative: $P(x|y=k) = \mathcal{N}(x, \mu^{(k)}, \Sigma^{(k)})$

$\Rightarrow P(y=k|x) \propto P(x|y=k) \pi_k \rightarrow \pi_k \sim \frac{\sum_n \mathbb{I}(y^{(n)}=k)}{N}$

SGD vs. Analytic (faster)

↓
 $w^{(t+1)} \leftarrow w^{(t)} - \eta \nabla_w [r^T r] \rightarrow r = y - Xw$

Logistic: $f(x; w) = \sigma(w^T x)$

$L(w) = P(y|X, w)$; $LL(w) = -\sum_n \log \sigma(z^{(n)} x^{(n)} w^T)$

softmax: $f_k = \frac{e^{x^T w^{(k)}}}{\sum_{k'} e^{x^T w^{(k')}}}$

Robust

$p(y=1|x, w, m) = \begin{cases} \sigma(w^T x) & m=1 \\ \frac{1}{2} & m=0 \end{cases}$; $p(m) = \begin{cases} 1-\epsilon & m=1 \\ \epsilon & m=0 \end{cases}$

Convex: $C(\alpha \underline{u} + (1-\alpha) \underline{u}') \leq \alpha C(\underline{u}) + (1-\alpha) C(\underline{u}')$

Neural Nets

$$h^{(l)} = p^{(l)}(w^{(l)} h^{(l-1)} + \underline{b}^{(l)}), \quad h^{(0)} = \underline{x}, \quad l=1 \dots N$$

Backprop

FW mode
 $\frac{\partial}{\partial \underline{u}} \rightarrow \star$
 $\frac{\partial}{\partial \underline{u}} \rightarrow \text{scalar}$

$$\frac{\partial f}{\partial u_i} = \sum_{d=1}^D \frac{\partial f}{\partial x_d} \frac{\partial x_d}{\partial u_i}$$

BW mode (Reverse)
 $\frac{\partial z}{\partial \underline{u}} \rightarrow \text{error wgt. } \star$

$$C = AB, \Rightarrow \bar{A} = \bar{C} B^T, \quad \bar{B} = A^T \bar{C}$$

$$C = p(A) \Rightarrow \bar{A} = p'(A) \odot \bar{C}$$

• $A \times B$ is $O(LMN)$
 $L \times M \quad M \times N$

Autoencoders

NN, hidden layer $k \ll D$ for dim reduction
 $L \rightarrow$ sparse; if $k=2$, visualize

PCA

$$\Sigma = \frac{1}{N} X^T X; \quad V = \text{eig}(\Sigma); \quad \text{project \& profit}$$

\nwarrow center first

$$VV^T = I$$

SVD

$$X \approx U S V^T$$

U $N \times k$ eigenvectors of XX^T
 S $k \times k$ diagonal
 V^T $k \times D$ eigenvectors of $X^T X$

PCA is SVD: $X \approx X V V^T$

\rightarrow probabilistic
PPCA

- assume k -dim Gaussian $\underline{z} \sim N(\underline{0}, I_k)$
- project onto D $\underline{x} = W \underline{z}$ (W is $D \times k$)
- Then $\underline{x} \sim N(\underline{0}, W W^T + \sigma^2 I)$ add some noise (for a cheeky encode)

MLP K (Bayesian)

• Linear Regression

$$p(y|x, w) = \mathcal{N}(y; f(x; w), \sigma_y^2)$$

pinch of noise

1) w/ MLE \equiv regular LR

2) + prior $p(w) = \pi(w) = \mathcal{N}(w; 0, \sigma_w^2 \mathbb{I})$

$$p(w|D) \propto p(D|w) p(w) \rightarrow \text{Bayesian}$$

conspiracy

likelihood

$$\propto \mathcal{N}(w; 0, \sigma_w^2 \mathbb{I}) \mathcal{N}(y; \phi w, \sigma_y^2 \mathbb{I})$$

b/p of Gaussian

w/ w_T, V_T Predictions

3) ~~infer~~ over posterior: $p(y|x) = \int p(y|x, w) p(w|D) dw$

test

test

posterior

4 (optional) Model choice: $p(y|x, M) = \int p(y|x, w, M) p(w|M) dw$

over prior

train

prior

5) test functions:

$p(D|M)$

\hat{y} is a point estimate; optimize $E_{p(y|D)} [L(y, \hat{y})]$

6) Cheeky - breeky (if linear)

$$p(y|D, x) = \mathcal{N}(y; x^T w_T, x^T V_T x + \sigma_y^2)$$

7) for linear $(p(D|M))^{-1} = \frac{p(w|D)}{p(w) p(D|w)}$

triple Gauss

Gaussian Processes

1) function f as vector, $\text{cov}[f_i, f_j] = k(x^{(i)}, x^{(j)})$

2) $f \sim GP$ iH $p(f) = \mathcal{N}(f; 0, k)$

3) Pinch a noise + test points

$$p\left(\begin{bmatrix} y \\ f^* \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} y \\ f^* \end{bmatrix}; 0, \begin{bmatrix} K(x, x) + \sigma_y^2 & K(x, x^*) \\ K(x^*, x) & K(x^*, x^*) \end{bmatrix}\right)$$

noise

sample

test

Kernels

linear: $k(x^{(i)}, x^{(j)}) = \sigma_w^2 \phi(x^{(i)})^T \phi(x^{(j)}) + \sigma_b^2$

Gaussian: $k(\underline{x}^{(i)}, \underline{x}^{(j)}) = \sigma_f^2 \exp\left(-\frac{1}{2} \sum_{d=1}^D (\underline{x}_d^{(i)} - \underline{x}_d^{(j)})^2 / \ell_d^2\right)$

function scale length btw turning points

$p(y | \underline{x}, \theta) \rightarrow$ marginal LH, hyperparameters

Logistic Regression

MAP: $\underline{w}^* = \arg\max [\log p(\underline{w} | D)] = \dots$ Bayes

MLE: no prior \underline{w}

$p(\underline{w} | D) \propto \prod \sigma(\underline{w}^T \underline{x}^{(n)}) \mathcal{N}(\underline{w}; \underline{0}, \sigma_w^2 \mathbb{I})$

↑ approximated by

↓ to find
Laplace Approx.

1) Define $E(\underline{w}) = -\log p(\underline{w}, D)$

2) $\underline{w}^* = \arg\min E(\underline{w})$

3) $H_{ij} = \frac{\partial^2 E(\underline{w})}{\partial w_i \partial w_j} \Big|_{\underline{w} = \underline{w}^*}$

4) $p(\underline{w} | D) \approx \mathcal{N}(\underline{w}; \underline{w}^*, H^{-1})$

5) predictions

same as $p(\underline{w} | D)$ up to a const.

At $\underline{w} = \underline{w}^*$

$p(D) \approx \frac{p(\underline{w}^*, D)}{|H|^{1/2}}$

approximates $p(\underline{w} | D)$

$$p(y | \underline{x}, D) \approx \int \sigma(\underline{w}^T \underline{x}) \mathcal{N}(\underline{w}; \underline{w}, V) d\underline{w}$$

$$= E_{\mathcal{N}(\underline{w}; \underline{w}, V)} [\sigma(\underline{w}^T \underline{x})]$$

$$= \int \sigma(a) \mathcal{N}(a; \underline{w}^T \underline{x}, \underline{x}^T V \underline{x}) da$$

Variational Inference

1) Define $q(\underline{w}; \alpha = \{\underline{w}, V\}) = \mathcal{N}(\underline{w}; \underline{w}, V)$

$D_{KL}(p || q) = \int p(\underline{z}) \log \frac{p(\underline{z})}{q(\underline{z})} d\underline{z} \geq 0$ (Gibbs)

if $p(\underline{w} | D)$ intractable

q is high where $p \neq 0$

high entropy, spread out

2) Minimize

$D_{KL}(q(\underline{w}; \alpha) || p(\underline{w} | D)) = - \int d\underline{w} q(\underline{w}; \alpha) \log p(\underline{w} | D) + \int d\underline{w} q(\underline{w}; \alpha) \log q(\underline{w}; \alpha)$

Variational Inference (ctnd.)

$$p(\underline{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\underline{w})p(\underline{w})}{p(\mathcal{D})}$$

$$D_{KL}(q||p) = \underbrace{\mathbb{E}_q[\log q(\underline{w})] - \mathbb{E}_q[\log p(\mathcal{D}|\underline{w})] - \mathbb{E}_q[\log p(\underline{w})]}_{J(q)} + \log p(\mathcal{D}) \rightarrow \text{const. w.r.t. } q$$

- Bound on $p(\mathcal{D})$ / marginal LH

$$D_{KL}(q||p) \geq 0, \Rightarrow \log p(\mathcal{D}) \geq \underline{J(q)}$$

ELBO: evidence lower bound

- Hardest term is $\mathbb{E}_q[\log p(\mathcal{D}|\underline{w})] =$

$$= \sum_{n=1}^N \mathbb{E}_q[\log p(y^{(n)} | x^{(n)}, \underline{w})]$$

- Minimize $J(\underline{m}, V)$ w.r.t. \underline{m}, V

- $V = LL^T$, $L = \begin{cases} \log(L_{ii}) & i=j \\ L_{ij} & i \neq j \end{cases}$

- Reparameterization trick for

$$\mathbb{E}_{N(\underline{w}; \underline{m}, V)}[f(\underline{w})] = \mathbb{E}_{N(\underline{v}; 0, \mathbb{I})}[f(\underline{m} + L\underline{v})]$$

sample \underline{w} , by $\underline{v} \sim N(0, \mathbb{I})$

$$\underline{w} = \underline{m} + L\underline{v}$$

- Monte-Carlo Estimate

$$\dots \approx \frac{1}{S} \sum_{s=1}^S f(\underline{m} + L\underline{v}^{(s)}), \quad \underline{v}^{(s)} \sim N(0, \mathbb{I})$$

$$\approx f(\underline{m} + L\underline{v}) \quad \bullet \text{ single unbiased estimate}$$

- Compute $\nabla_{\underline{m}}[\dots]$, $\nabla_L[\dots]$ & SGD

Gaussian Mixture Models

$z^{(n)} \in \{1, \dots, K\} \rightarrow$ classes

• MLE: $p(\underline{x}^{(n)} | \theta) = \sum_k \pi_k N(\underline{z}^{(n)}; \mu^{(k)}, \Sigma^{(k)})$

Expectation Maximization

$\sigma_k^{(n)} = \delta_{z^{(n)}, k} \rightarrow$ responsibility for cluster k

$$\pi_k = \frac{\sigma_k}{N}, \quad \sigma_k = \sum_{n=1}^N \sigma_k^{(n)}, \quad \mu^{(k)} = \frac{1}{\sigma_k} \sum_{n=1}^N \sigma_k^{(n)} \underline{x}^{(n)}$$

$$\Sigma^{(k)} = \frac{1}{\sigma_k} \sum_{n=1}^N \sigma_k^{(n)} \underline{x}^{(n)} \underline{x}^{(n)T} - N^{(k)} \mu^{(k)} \mu^{(k)T}$$

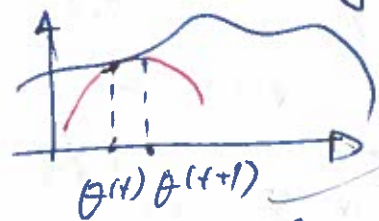
• E-step: $\sigma_k^{(n)} = p(z^{(n)} = k | \underline{x}^{(n)}, \theta) = \frac{\pi_k N(\underline{x}^{(n)}; \mu^{(k)}, \Sigma^{(k)})}{\sum_l \pi_l N(\underline{x}^{(n)}; \mu^{(l)}, \Sigma^{(l)})}$

• M-step: update θ

L1: $c(\underline{w}) = E(\underline{w}) + \lambda \|\underline{w}\|_1$

Newton: $\underline{w}^{(t+1)} \leftarrow \underline{w}^{(t)} - H^{-1} f$

Bound-based:



needs tight bound

Ensembling:

$$p(y | \underline{x}, D) \approx \frac{1}{S} \sum_{s=1}^S p(y | \underline{x}, \underline{w}^{(s)}), \quad \underline{w}^{(s)} \sim p(\underline{w} | D)$$

Bagging:
• sample w w/ replacement from N points S times
• fit S models, then profit

Mixture of experts: $p(y | \underline{x}, \theta) = \sum_k p(y | \underline{x}, z=k, \theta) \underbrace{p(z=k | \underline{x}, \theta)}_{N.N.}$

Boosting: $F_{m+1}(x) = F_m(x) + h(x)$ weak model

fit $F_{m+1}(x)$ to $y - F_m(x) = h(x)$