# MLPR
## final notes

**linear regression**

total square error: $E(w, b) = (y - f)^T (y - f)$

can add bias as extra column of 1's to $X$ and a $w_0$ to $\underline{w}$

RBS form: $\exp\left(-(\underline{x} - \underline{c})^T (\underline{x} - \underline{c}) / h^2\right)$
    (↑ center)    (↘ width $(h)$)

**regularization:**

$$E_\lambda(\underline{w}; y, \Phi) = (y - \Phi \underline{w})^T (y - \Phi \underline{w}) + \lambda \underline{w}^T \underline{w}$$

$$= (y' - \Phi' \underline{w})^T (y' - \Phi' \underline{w}) \quad \text{where } y' = \begin{bmatrix} y \\ 0_k \end{bmatrix}, \; \Phi' = \begin{bmatrix} \Phi \\ \sqrt{\lambda} \, \mathbb{I}_k \end{bmatrix}$$

**model evaluation:**

gen. error $E = E_{p(x,y)}\left[ L(y, f(\underline{x})) \right] = \int L(y, f(\underline{x})) \, p(x, y) \, dx \, dy$

use mean test error: $\frac{1}{M} \sum_{m=1}^{M} L\left(y^{(m)}, f(\underline{x}^{(m)})\right) \qquad \underline{x}^{(m)}, y^{(m)} \sim p(x, y)$

**univariate Gaussian:** $p(z) = N(z, \mu, \sigma^2) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma^2}(z - \mu)^2\right)$

**error bars:** $\mathrm{var}[\bar{x}] = \frac{\sigma^2}{N}$

typical deviation: $\mu \pm \frac{\hat{\sigma}}{\sqrt{N}} \quad \Rightarrow \text{ can apply this to errors in test set}$

**multivariate Gaussians**

$$\mathrm{cov}[y] = E[yy^T] - E[y]E[y]^T$$
$$= E[A\underline{x}\underline{x}^T A] - E[A\underline{x}]E[A\underline{x}]^T \quad \text{if } y = A\underline{x} \text{ and } x \sim N(0,1)$$
$$= A \, E[xx^T] \, A^T$$
$$= AA^T \qquad\qquad \text{note on determinants: } |\Sigma| = |AA^T| = |A||A^T| = |A|^2$$
$$= \Sigma$$

$$p(\underline{z}) = N(\underline{z}; \mu, \bar{\Sigma}) = \frac{1}{|\Sigma|^{1/2} (2\pi)^{D/2}} \exp\left(-\frac{1}{2}(z - \mu)^T \Sigma^{-1} (z - \mu)\right)$$

**classification**

Bayes classifier: $P(y = k | \underline{x}) \propto p(\underline{x} | y) P(y = k)$
$$\propto N(\underline{x}; \mu_k, \Sigma_k) \, \pi_k \qquad \text{where } \pi_k = \frac{\sum_n I(y^{(n)} = k)}{N}$$
$$= S_k$$

$$P(y = k | \underline{x}, \theta) = \frac{S_k}{\sum_{k'} S_{k'}}$$

"Naive Bayes": features are independent, i.e., $\Sigma_k$ is diagonal for all $k$
natural when features $\underline{x}$ are binary:

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

regression and gradients:

calculate square error: $r^\top r = (y - Xw)^\top (y - Xw)$

then take deriv: $\nabla_w r^\top r = -2X^\top y + 2X^\top X w$

now can do:

closed form: $w = (X^\top X)^{-1} X^\top y$   (after setting $\nabla_w r^\top r$ to zero)

iterative: $w \leftarrow w - \eta \nabla_w [r^\top r]$

logistic regression:

$$NLL = -\sum_{n=1}^{N} \log \left[ \sigma(w^\top x^{(n)})^{y^{(n)}} \left(1 - \sigma(w^\top x^{(n)})\right)^{1 - y^{(n)}} \right]$$

$$= -\sum_{n=1}^{N} \log \sigma(z^{(n)} w^\top x^{(n)})$$

$\dfrac{\partial \sigma(a)}{\partial a} = \sigma(a)(1 - \sigma(a))$

$$\nabla_w NLL = -\sum_{n=1}^{N} (1 - \sigma_n) z^{(n)} x^{(n)}$$

approx grad: $\dfrac{\partial f(w)}{\partial w} \approx \dfrac{f(w + \frac{\varepsilon}{2}) - f(w - \frac{\varepsilon}{2})}{\varepsilon}$   (i.e., $\frac{\Delta y}{\Delta x}$)

softmax / robust solutions:

$f_k = \dfrac{\exp(w^{(k)\top} x)}{\sum_{k'} \exp(w^{(k')\top} x)}$          $\nabla_{w^{(k)}} \log f_c = (y_k - f_k) x$

a robust model randomly sets $y$ to $\frac{1}{2}$ w/ probability $\varepsilon$

$\begin{cases} \exp: & \text{real} \to \text{pos} \\ \log: & \text{pos} \to \text{real} \\ \text{sigmoid:} & \text{real} \to (0,1) \\ \text{logit:} & (0,1) \to \text{real} \end{cases}$

neural nets

don't set weights to zero or to $N(0,1) \Rightarrow$ set as $\dfrac{0.1 \, N(0,1)}{\sqrt{k}}$

chain rule: $\dfrac{\partial f}{\partial r} = \dfrac{\partial f}{\partial x}\dfrac{\partial x}{\partial r} + \dfrac{\partial f}{\partial y}\dfrac{\partial y}{\partial r}$

$\dfrac{\partial f}{\partial u_i} = \sum_{d=1}^{D} \dfrac{\partial f}{\partial q_d}\dfrac{\partial x_d}{\partial u_i}$   (generalized)
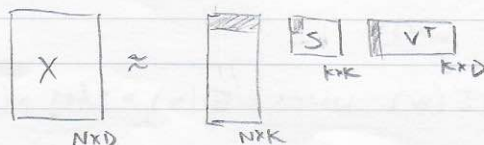
$\dfrac{\partial f}{\partial X_{ab}} = \overline{X}_{ab}$

automatic differentiation: $z = xy \Rightarrow \overline{X} = \overline{Z} Y^\top$ and $\overline{Y} = X^\top \overline{Z}$   (ordinary matrix operations: not $O(M^4)$!)

backprop the error signals: $\delta_i^{(l)} = \dfrac{\partial E}{\partial a_i^{(l)}} = \overline{a}_i^{(l)}$

## PCA

denoising autoencoder: essentially dropout

$X \approx USV^T$



▨ = operations for $1^{st}$ PC

$US \approx XV$ is data projected down to $k$ dims
$USV^T \approx XVV^T$ is data projected back up (but it's low-rank $k$)

## Bayesian regression:

' conjugate prior is one where the product of the prior and the likelihood combines to give a dist w/ the same functional form as the prior

$$p(\underline{w}|D) \propto p(D|\underline{w}) \, p(\underline{w})$$
$$\propto \prod_n N(y^{(n)}; \underline{w}^T \underline{\phi}(x^{(n)}), \sigma^2) \, N(\underline{w}; 0, \sigma^2_{prior})$$

assuming we're updating a Gaussian and have a Gaussian prior

## Bayesian inference and prediction:

prediction: $p(y|x, D) = \int p(y, \underline{w}|x, D) \, dw$

$$= \int p(y|x, \underline{w}) \underbrace{p(\underline{w}|D)}_{\text{see last section}} dw$$

(posterior)
$p(\underline{w}|D) = N(\underline{w}; w_N, V_N)$ if Gaussian

linear and Gaussian models can be solved in closed form

for linear: $\mu_f = E_w[x^T \underline{w}] = x^T w_N$

$var[f] = \underline{x}^T V_N \underline{x}$ $\qquad (= E[(x^T\underline{w} - x^T\underline{w}_N)^T (x^T\underline{w} - x^T\underline{w}_N)])$

$\qquad\qquad f - \mu_f \qquad f - \mu_f$

$p(y|D, x) = N(y; \underline{x}^T w_N, \underline{x}^T V_N \underline{x} + \sigma^2)$

add this for noise going from $f$ to $y$

## Bayesian model selection:

small feat
overly complex
models

$$p(y|X, M) = \int p(y|X, \underline{w}, M) \, p(\underline{w}|M) \, dw$$

$\frac{1}{\sqrt{\alpha}}$ ↗ "$\sigma$" here

can do marginal likelihood to estimate hyperparams (e.g., $\sigma_{prior}$ and $\sigma_y^2$)

$$p(\bar{y}|X, \alpha, \sigma) = \int p(y|X, \underline{w}, \sigma) \, p(\underline{w}|\alpha) \, dw \longleftarrow$$

maximize this to find $\alpha$ and $\sigma$

## GPs:

joint dist: $p\left(\begin{bmatrix} y \\ \underline{f}_* \end{bmatrix}\right) = N\left(\begin{bmatrix} y \\ \underline{f}_* \end{bmatrix}; \underline{0}, \begin{bmatrix} K + \sigma_n^2 I & K_* \\ K_* & K_{**} \end{bmatrix}\right)$

can compute hyperparams by max likelihood

$\max \log(y|X, \theta) = \max \log N(\cdot)$

$k(x^{(i)}, x^{(j)}) = \sigma_f^2 \exp\left(-\frac{1}{2} \sum_{d=1}^{D} (x_d^{(i)} - x_d^{(j)})^2 / l_d^2\right)$

$\sigma_f$ ↳ typical vertical height

"typical dist. b/t turning points"

or by integrating $p(\underline{f}_*|y, X) =$ marginalize out $\theta$...

Bayesian logistic regression and Laplace approx

MAP: $\underline{w}^* = \arg\max_{\underline{w}} [\log p(\underline{w}|D)] = \arg\max_{\underline{w}} [\log p(D|\underline{w}) - \frac{1}{2\sigma_w^2} \underline{w}^T \underline{w}]$

  likelihood   "$\lambda$"

$\rightsquigarrow$ basically regularized MLE

Laplace approx

 1. find $\underline{w}^*$: $\underline{w}^* = \arg\min_{\underline{w}} E(\underline{w})$ where $E(\underline{w}) = -\log p(\underline{w}, D)$

   energy

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

 3. set $p(\underline{w}|D) \simeq N(\underline{w}; \underline{w}^*, H^{-1})$

   sometimes not so great

variation inference:

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

usually want $D_{KL}(q\|p)$ where $q$ approximates $p$

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

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

  cross-entropy··· small = good    (neg) entropy; want entropy to be large

$$= \mathbb{E}_q [\log p(D|\underline{w})] + \cdots + \log p(D) \quad \text{(spread out)}$$

   $J(q)$

only the neg log likelihood term can't be computed in closed form

$$-\mathbb{E}_{N(\underline{w}; \underline{w}, V)}[\log p(D|w)]$$

So use MC

Gaussian mixture models:

$$p(\underline{x}^{(n)} | \theta) = \sum_k p(\underline{x}^{(n)}, z^{(n)} = k | \theta)$$

     $\pi_k$

$$= \sum_k p(\underline{x}^{(n)} | z^{(n)} = k, \theta) P(z^{(n)} = k | \theta)$$

$$= \sum_k \pi_k N(\underline{x}^{(n)}; \mu^{(k)}, \Sigma^{(k)})$$

use EM to take NLL of $p(D|\theta)$

E: set soft responsibilities: $r_k^{(n)} = p(z^{(n)} = k | x^{(n)}, \theta) = \frac{\pi_k N(\cdot)_k}{\sum_\ell \pi_\ell N(\cdot)_\ell}$

M: update params $\theta = \{\pi, \{\mu^{(k)}, \Sigma^{(k)}\}$

$$\mu^{(k)} = \frac{1}{r_k} \sum_{n=1}^N r_k^{(n)} x^{(n)}$$

$$\Sigma^{(k)} = \cdots$$