

Gaussian process classifiers:

Recall logistic regression:

$$p(y=+1|x, w) = \sigma(x^T w)$$

$$p(y_i/x_i, w) = \sigma(y_i f_i) \quad \left\{ \begin{array}{l} \text{then} \\ \text{for LR,} \\ \text{other sigmoids possible} \end{array} \right. \frac{1}{1 + \exp(-w^T x)}$$

Now write $p(w) = N(0, \Sigma_p)$

then $\log p(w|x, y) = -\frac{1}{2} w^T \Sigma_p w + \sum_i \log \sigma(y_i f_i) + \text{const.}$

Our version so far is maximum likelihood

$$\hat{=} \text{choose } w = \text{argmax } \sum \log \sigma(y_i f_i)$$

penalized LR is maximum a posteriori

$$\hat{=} \text{choose } w = \text{argmax} \left[\sum_i \log \sigma(y_i f_i) - \frac{1}{2} w^T \Sigma_p w \right]$$

Alternative:

$$p(y_* = 1 | x_*, D) = \int p(y_* = 1 | w, x_*) p(w | D) dw$$

i.e. average the prediction.

Gaussian process classification:

- replace $w^T x$ with a GP
- notice we don't know $f(x)$ at the training points - nuisance parameters.
- classifying a point x_*
 - 1): find f_*

$$p(f_* | X, y, x_*) = \int p(f_* | X, \cancel{y}, x_*, f) \cdot p(f | X, y) \cdot df$$

2): prediction

$$\pi_x = p(y_* = 1 | X, y, x_*) = \int p(y_* = 1 | X, y, x_*, f_*) \cdot p(f_* | X, y, x_*) \cdot df_*$$

	unknown	
	↓	
obs	f_i	obs
y_i	?	x_i
1	?	...
0	?	...
⋮		etc

Notice $p(f^*/X, x_*, f)$ is straight forward from linear regression

$$N \left(K(x_*, x) [K(x, x)]^{-1} f, K(x_*, x_*) - K(x_*, x) [K(x, x)]^{-1} K(x, x_*) \right)$$

but we need $p(f/X, y)$.

Laplace approx:

- approx $p(f/X, y)$ with a Gaussian
- mean is at $\text{argmax } p(f/X, y) = f_0$
- covar is $-H_f \log p(f/X, y) \Big|_{f=f_0}$

$$p(f|x, y) \propto p(f|x) \cdot p(y|f)$$

\uparrow GP prior \uparrow from sigmoid

Now to find $f_0 = \operatorname{argmax} p(f|x, y)$

$$\begin{aligned} \psi(f) &\stackrel{\Delta}{=} \log p(y|f) + \log p(f|x) \\ &= \log p(y|f) - \frac{1}{2} f' K^{-1} f - \frac{1}{2} \log |K| \\ &\qquad\qquad\qquad - \frac{n}{2} \log 2\pi \end{aligned}$$

$$\nabla_f \psi = \nabla_f \log p(y|f) - K^{-1} f$$

$$H_f \psi = H_f \log p(y|f) - K^{-1} = -W - K^{-1}$$

\uparrow this will be diagonal

Note $K = K[x, x]$

work with logistic

$$p(y_i | f_i) = -\log(1 + \exp(-y_i f_i))$$

$$\frac{\partial}{\partial f_i} \log p(y_i | f_i) = t_i - \pi_i$$

$$\left(\frac{1 + y_i}{2}\right)$$

$$p(y_i = 1 | f_i) = -\log(1 + \exp(-f_i))$$

$$\frac{\partial^2}{\partial f_i^2} \log p(y_i | f_i) = -\pi_i [1 - \pi_i]$$

Notice that at f_0 ,

$$\underline{f_0} = K \left[\nabla \log p(y | \underline{f_0}) \right] \quad (\text{from gradient})$$

Use Newton's method to find $\underline{f_0}$

then

$$p(f | X, y) \approx q(f | X, y) = N(f_0,$$

in what sense?



$$(K^{-1} + W)^{-1}$$

- integration

Prediction:

want posterior mean for f_* .

$$\begin{aligned} E[f_* | X, y, x_*] &= \int \left[\int f_* p(f_* | X, y, x_*, f) df_* \right] \times \\ &\quad p(f | X, y, x_*) df \\ &= \int E[f_* | X, y, x_*, f] \cdot p(f | X, y, x_*) df \end{aligned}$$

now approx p with q from above

$$E_q[f_* | X, y, x_*] = K[x_*, x] [K(X, X)]^{-1} f_0$$

Variance:

$$\begin{aligned} &E \left[\left(f_* - E(f_*) \right)^2 \right] \\ &= E \left[\left\{ \left(f_* - E(f_* | f) \right) + \left(E(f_* | f) - E(f_*) \right) \right\}^2 \right] \end{aligned}$$

not a function of f

so these two are indep

So:

$$K \text{ Var} = E_{P(f_* | X, X_*, f)} \left[(f_* - E(f_* | X, X_*, f))^2 \right] + \frac{E}{q} \left[(E(f_* | X, X_*, f) - E(f_* | X, y, X_*))^2 \right]$$

• First term from GP regression

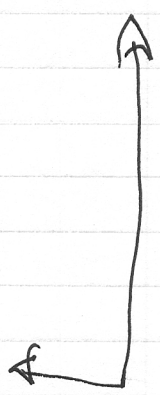
$$K(X_*, X_*) - K(X_*, X) [K(X, X)]^{-1} K(X, X_*)$$

• Second term from Gaussian approx to q

$$K(X_*, X) [K(X, X)]^{-1} (K^{-1} + w)^{-1} K[X, X]^{-1} K(X, X_*)$$

• now our predictions are

$$\begin{aligned} \tilde{\pi}_* &= E_q [\pi_* | X, y, X_*] \\ &= \int \sigma(f_*) q(f_*) df_* \end{aligned}$$



• Now choose parametric

$$q(\theta) \sim N(m_\theta, v_\theta I)$$

now build

$$q_{i+1}(\theta) \quad \text{by:}$$

$$q_0(\theta) = t_0(\theta)$$

$$\hat{p}_{i+1}(\theta) = \frac{t_{i+1}(\theta) q_i(\theta)}{\int_{\Theta} t_{i+1}(\theta) q_i(\theta) d\theta}$$

and $q_{i+1}(\theta)$ chosen to minimize

$$D(\hat{p}_{i+1}(\theta) \parallel q_{i+1}(\theta))$$

if the q_i are spherical gaussian

$$m_i, v_i$$

Alternative strategy: Expectation-propagation

must approx $p(f/x, y) = \frac{1}{Z} \cdot p(f/x) \cdot \prod_{i=1}^n p(y_i / f_i)$

posterior on function value

strategy: approx

$$p(y_i / f_i) \approx \approx t_i(f_i / \tilde{z}_i, \tilde{\mu}_i, \tilde{\sigma}_i^2)$$

$$\triangleq \tilde{z}_i N(f_i / \tilde{\mu}_i, \tilde{\sigma}_i^2)$$

site params $\tilde{z}_i, \tilde{\mu}_i, \tilde{\sigma}_i^2$

this gives

$$q(f/x, y) \triangleq \frac{1}{Z_{EP}} p(f/x) \cdot \prod_i t_i$$

$$= N(\mu, \Sigma)$$

How to choose $\tilde{z}_i, \mu_i, \sigma_i^2$?

- Start from some approximation
- leave out t_i term
 \rightarrow cavity dist q_{-i}
- now choose site i params to best approx $q_{-i} p(y_i | f_i)$

iterate:

cavity dist

$$q_{-i}(f_i) = \int \left[p(f | x) \prod_{j \neq i} t_j(f_j) \right] df_{-i}$$

this is gaussian, so marginal is easy.

$$q_{-i}(f_i) = N(\mu_{-i}, \sigma_{-i}^2)$$

