# Gaussian process classification model

$\in \mathbb{R}$

- 'training' data $D = \{(x_1, t_1), (x_2, t_2), \ldots, (x_N, t_N)\} \in [\mathbb{R} \times \{0,1\}]^N$

  $\hookleftarrow$ class label $\in \{0,1\}$

- 'test' data: $\tilde{x}$

- model: $y(x) \sim$ GP with $m(x) = 0$, kernel $k(x, x')$  (latent process)

  observation: $t_i = \text{Bernoulli}(\sigma(\overbrace{y(x_i)}^{y_i}))$ $\leftarrow$ link fn (sigmoid) $\cdot \mathbb{R} \to [0,1]$

  (i.e., $p(t|y_i) = \sigma(y_i)^t (1 - \sigma(y_i))^{1-t}$) $\leftarrow$ $t(x) = \begin{cases} 0 & \text{wp } \frac{1}{1+e^{y(x)}} = \sigma(-y(x)) \\ 1 & \text{wp } \frac{e^{y(x)}}{1+e^{y(x)}} = \sigma(y(x)) \end{cases}$

- prior: with $K_D, k, c$ as in GP regression

  on $y(x)$

$$(y_1, y_2, \ldots, y_N, \tilde{y}) \sim \mathcal{N}\left(0, \begin{bmatrix} K_D & k \\ k^\mathsf{T} & c \end{bmatrix}\right)$$

- posterior: how do we compute $p(\tilde{y}|D)$?

$K_D = \{k(x_i, x_j)\}$, $k = \{k(x_i, \tilde{x})\}$

$c = k(\tilde{x}, \tilde{x})$

## posterior

- 'training' data $D = \{(x_1, t_1), (x_2, t_2), \ldots, (x_N, t_N)\} \in [\mathbb{R} \times \{0,1\}]^N$
- model: $\underline{y(x)} \sim$ GP with $m(x) = 0, k(x, x')$, $t_i =$ Bernoulli$(\sigma(y(x_i)))$
- likelihood given $y_i = y(x_i)$

$$\log p(t|y(x)) = t^T y + \sum_{i=1}^{N} \log(1 + e^{y_i})$$

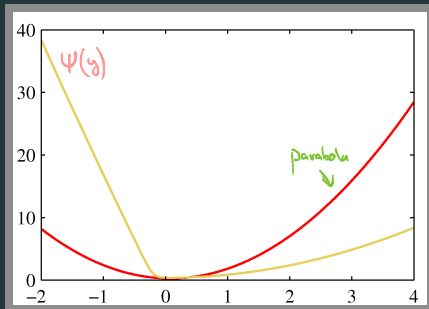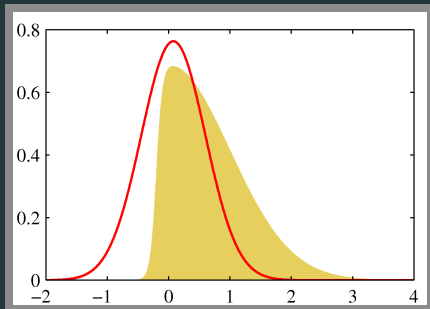- negative log of posterior $-\log p(y|t)$ ← approx this as a quadratic

$$\psi(y) = \frac{1}{2} y^T K^{-1} y + \frac{1}{2} \log|K| \bullet \left( t^T y + \sum_{i=1}^{N} \log(1 + e^{y_i}) \right) + \text{const}$$

$\underbrace{\qquad\qquad\qquad}$ $-\log$ of prior

$\underbrace{\qquad\qquad\qquad\qquad}$ -ive log of likelihood (prev page)

const → normalizata

prior $P(y) = (2\pi)^{-N/2} (|K|)^{-1/2} \exp\left( \frac{(y-m)^T K^{-1} (y-m)}{2} \right)$

# the Laplace approximation

approximate posterior as a multivariate Gaussian

ie - 'match the mode and the Hessian'



Want $P(y|D) \sim N(\tilde{\mu}, \tilde{\Sigma}) \Rightarrow -\log p(y|D) \approx \frac{1}{2}(y-\tilde{\mu})^T \tilde{\Sigma}^{-1}(y-\tilde{\mu})$

$+\frac{1}{2}\log|\tilde{\Sigma}|$

Laplace approx - set $\tilde{\mu} = \nabla \Psi(y)\big|_{\tilde{\mu}} = 0$, $\tilde{\Sigma} = \nabla\nabla\Psi(y)$

# Laplace approximation for GP classification

Final output

$$p(y|D) \sim N\left(y^*, H^{-1}\right)$$

where
$$K^{-1}y^* = t - \sigma_N(y^*)$$

$$H^{-1} = K^{-1} + W^*$$

$$\tilde{y}(\tilde{x}) \sim N(\cdot, \cdot)$$

$$\mathbb{E}\left[\tilde{y}(\tilde{x})|D\right] = k^T K^{-1} y^* = k^T \left(t - \sigma_N(y^*)\right)$$

$$Var\left(\tilde{y}(\tilde{x})|D\right) = c - k^T H k = c - k^T\left(K^{-1} + W^*\right)^{-1} k$$

and $t(\tilde{x}) \sim Bernoulli\left(\sigma(\tilde{y})\right)$

Q: Given $\hat{x}$, what do we want to predict?

A: Depends on loss fn · · ·

Typical setting (0 loss): $L(t(\hat{x})) = \mathbb{E}\left[ \mathbb{1}\{t(\hat{x}) \neq \tilde{t}\} \right]$

'true class label' for $\hat{x}$



$\sigma(y)$

output of classifier

· 'Bayes classifier' $\quad 1 - \mathbb{1}\{\tilde{t}=1\}$

· $L(t(\hat{x})) = p(\hat{t}=1) \cdot \mathbb{1}\{t(\hat{x})=0\}$

$\quad + p(\hat{t}=0) \cdot \mathbb{1}\{t(\hat{x})=1\}$

$= \underbrace{p(\hat{t}=1)}_{\text{const, control}} + \mathbb{1}\{t(\hat{x})=0\} \underbrace{\left( p(\hat{t}=1) - p(\hat{t}=0) \right)}_{\substack{\text{+ve, set } t(\hat{x})=1 \\ \text{else set } t(\hat{x})=0}}$

ie - Bayes classifier is the MAP estimator for $t(\hat{x})$

For Bayes classifier, need to compute

$$P\left[t(\tilde{x})=1|D\right] = p\left(\tilde{t}|D\right) = \int \sigma(\tilde{y}) \, p(\tilde{y}|D) \, d\tilde{y}$$

$\frac{e^{\tilde{y}}}{1+e^{\tilde{y}}}$   complicated fn :(

$N\left(\mu_0, \Sigma_0\right)$

known in closed form!
(Sec 4.5 of Bishop)

How can we evaluate this?

1) Variational Approx — replace $\sigma(\tilde{y})$ by 'probit link fn'
$\phi(\tilde{y}) = $ 'tail prob of the Gaussian'

2) Monte Carlo simulation

# Laplace approximation: model selection

<u>Idea</u> - Maximize marginal likelihood $p(t|\theta)$ (ie, min $-\log p(t|\theta)$)
(Bayesian Occam's Razor)
$\uparrow$
hyper params

$$p(t|\theta, x) = \underbrace{p(y|\theta)}_{\text{Prior}} \underbrace{p(t|y)}_{\text{likelihood}} \leftarrow \tilde{L}(y) = N(y^*, H^{-1})$$
$$\tilde{q}(t|\theta, x) \qquad \underbrace{p(y|\theta, D)}_{\text{posterior}} \leftarrow N(\tilde{\mu}_D, \tilde{\Sigma}_D)$$
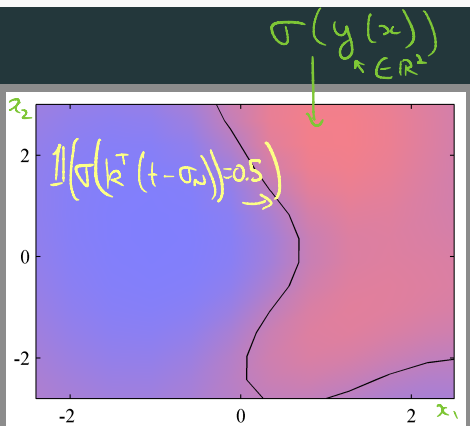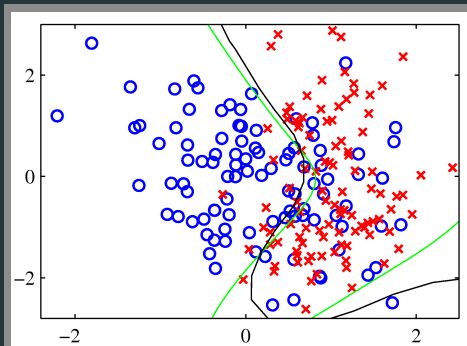
Laplace approx

- posterior, likelihood not available in chosed form — use Laplace approx instead

$$-\log\left(\tilde{q}(t|\theta, x)\right) = \frac{1}{2} y^{*T} K^{-1} y^* + \frac{1}{2}\log|K| + \frac{1}{2}\log|K^{-1}+W|$$
$\uparrow$
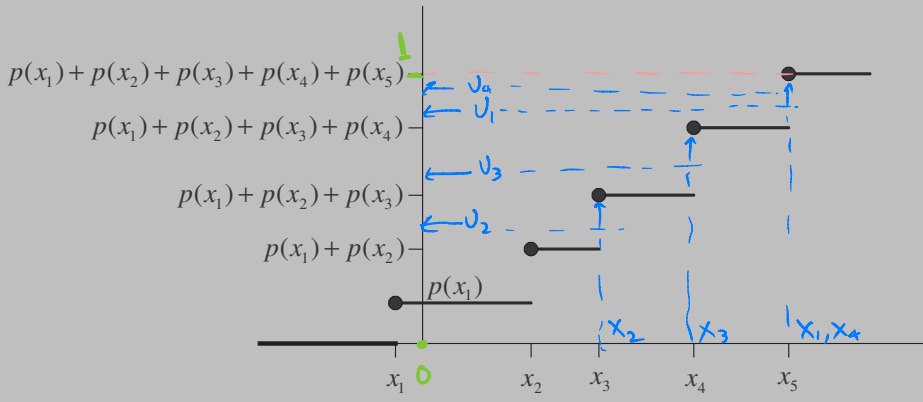Laplace approx for latent vars
$$-\log\left[p(t|y^*)\right]$$
$$\underbrace{\sigma(y^*)^t (1-\sigma(y^*))^{1-t}}$$

$\sigma\big(y(x)\big)$ $x \in \mathbb{R}^2$

$\mathbb{1}\big(\sigma\big(k^\top(t-\sigma_N)\big)=0.5\big)$
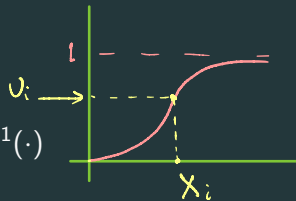
# basic monte carlo techniques

$X$ takes values $x_1 \leq x_2 \leq \ldots \leq x_5$, $\mathbb{P}[X = x_i] = p(x_i)$

## the inversion method

$X$ continuous r.v. with pdf $f$ and c.d.f. $F(\cdot)$

- want to generate samples of $X$.
- $F(\cdot)$ non-decreasing $\implies$ can define inverse $F^{-1}(\cdot)$
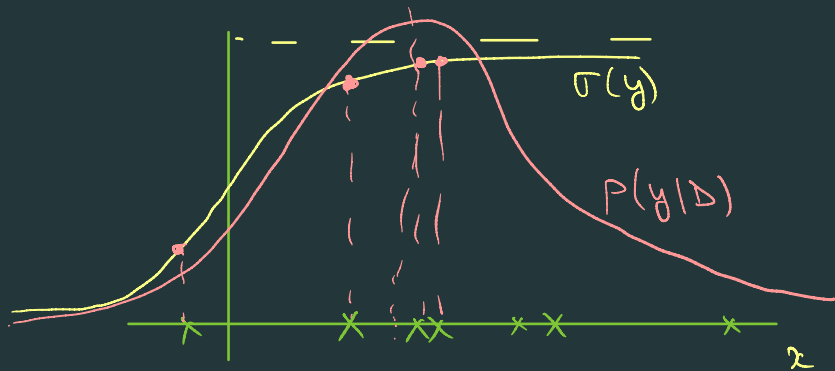- $F(x) = u \iff F^{-1}(u) = x$



### inversion method

given desired cdf $F$ (continuous, increasing), generate sample $X_0 \sim F$ as:

1. generate $U \sim U[0, 1]$.
2. return $X_o = F^{-1}(U)$.
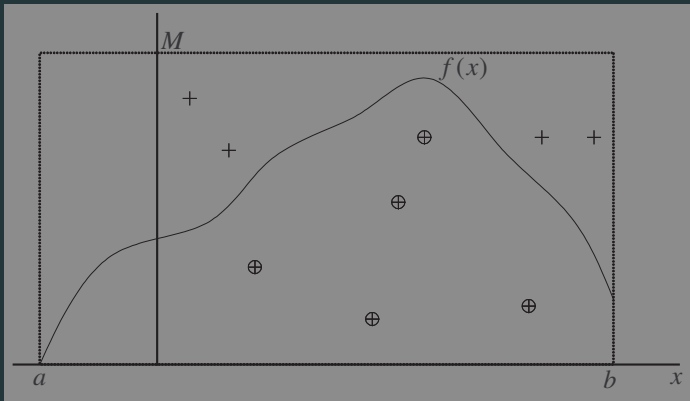
# Application – Compute integrals



$\sigma(y)$

$P(y|D)$

$x$

$$I = \mathbb{E}_{x \sim F}\left[\sigma(x)\right] \simeq \frac{1}{N}\sum_{i=1}^{N}\sigma(x_i)$$

## rejection sampling
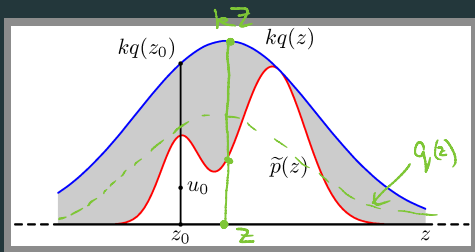
want samples of a rv $X \in [a, b]$, with pdf $f(x) \leq M$

1. Generate $U_1, U_2 \sim U[0, 1]$, and set $Z_1 = a + (b - a)U_1$, $Z_2 = MU_2$

2. if $Z_2 \leq f(Z_1)$, return $X_o = Z_1$; else, reject and repeat

# generalized rejection sampling



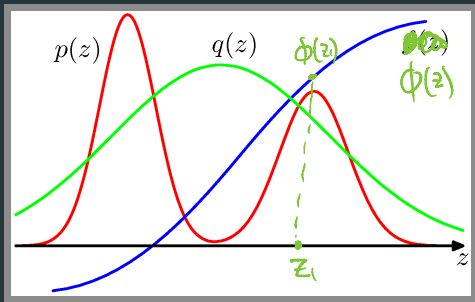- Given a 'sampler' $Z \sim Q$
- Want samples $X \sim P$

- find $k$ s.t $k q(x) \geqslant p(x) \; \forall x$ $\left( ie, k \geqslant \max \frac{p(x)}{q(x)} \right)$

- Generate $Z \sim Q$

- Accept (ie set $X = Z$) w.p $\dfrac{p(z)}{k q(z)}$, else repeat

# importance sampling (for estimating integrals)

- given function $\phi(\cdot)$, want $\mathbb{E}[\phi(X)]$ where $X \sim P$
- can generate samples $Z \sim Q$

## importance sampling

1. generate $Z_1, Z_2, \ldots, Z_L \sim Q$
2. compute $\mathbb{E}[\phi(X)] = \frac{1}{L} \sum_{i=1}^{L} w_i \phi(Z_i)$, where $w_i = p(Z_i)/q(Z_i)$



$$\mathbb{E}_{X \sim P}[\phi(x)] = \int \phi(x) \, p(x) \, dx$$

$$= \int \phi(x) \left( \frac{p(x)}{q(x)} \right) q(x) \, dx$$

$$= \mathbb{E}_{Z \sim Q}\left[ \phi(z) \, W(z) \right]$$

$$\underset{p(z)/q(z)}{\uparrow}$$

$$= \frac{1}{L} \sum_{i=1}^{L} w_i \phi(z_i)$$