

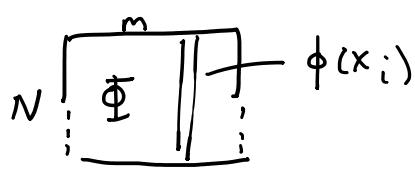
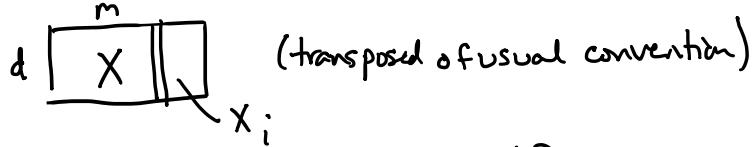
# Gaussian Processes

Friday, March 27, 2020 3:05 PM

For ML and beyond - not in classical PAC framework since it's Bayesian

Notation:  $x \in \mathbb{R}^d$ ,  $S = ((x_i, y_i))_{i=1}^m$ , as usual

feature map  $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^N$  (or, later,  $\rightarrow$  Hilbert space, possibly  $\infty$ -dim)



GP use kernels (like SVM),  
and as w/ kernel-trick,  
only need to know  $\langle \phi(x), \phi(x') \rangle = k(x, x')$   
and don't need to know  $\phi$  explicitly

Idea of GP:

we predict w/ a function  $f$  ( $f: X \rightarrow \mathbb{R}$  for regression  
and instead of constraining  $f \in \mathcal{F}$ ,  $f: X \rightarrow \{\pm 1\}$  for classification)

we'll instead encode prior knowledge via

a Bayesian prior on  $f$

i.e.  $f$  is a random function = stochastic process  
(generalizes random variable)

then, given observations  $S$ ,

we update, i.e., compute the posterior distribution using Bayes' Rule

### 1.1 A Pictorial Introduction to Bayesian Modelling

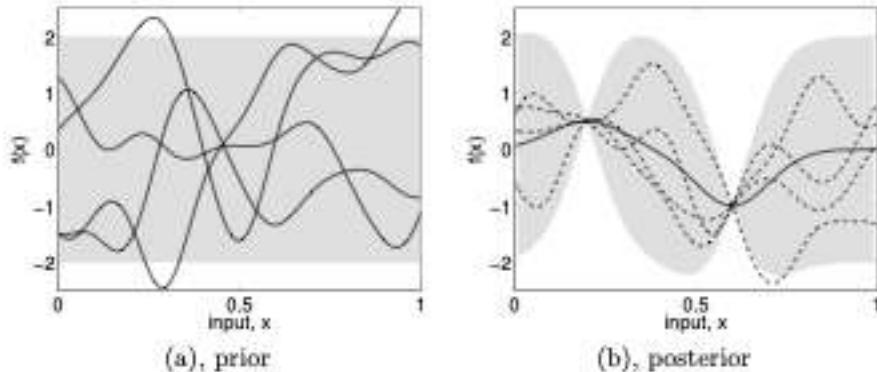
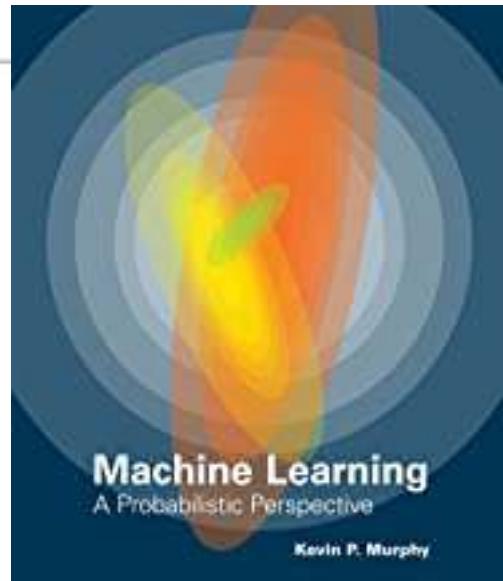


Figure 1.1: Panel (a) shows four samples drawn from the prior distribution. Panel (b) shows the situation after two datapoints have been observed. The mean prediction is shown as the solid line and four samples from the posterior are shown as dashed lines. In both plots the shaded region denotes twice the standard deviation at each input value  $x$ .



## Classification

is a bit more complicated (we'll need to "squash" the output to the range  $[0, 1]$  using, eg, the logistic function, and this complicates things)

## Regression

is cleaner, so will focus on this.

We essentially always do the kernel trick, so really focus on linear regression

$$\text{our model is } f(x) = \langle w, x \rangle \text{ or } f(x) = \langle w, \Psi(x) \rangle \\ w \in \mathbb{R}^d \qquad \qquad \qquad w \in \mathcal{H} \leftarrow \text{Hilbert space} \\ \Psi: \mathbb{R}^d \rightarrow \mathcal{H}$$

So a prior on  $f$  really means a prior on  $w$

For least-squares, we assume error term is  $\varepsilon \sim N(0, \sigma_m^{-2})$

$$y = f(x) + \varepsilon \quad \text{for some true } f.$$

i.e.,  $y = \langle w, x \rangle + \varepsilon$  for some true but unknown  $w$

If we know  $w, x$ , then

$$y | w, x \sim N(\langle w, x \rangle, \sigma_m^2)$$

Gaussian facts:

- i) Sum of independent Gaussians is Gaussian (i.e. convolve PDFs)

Thus, given a dataset  $X = d \boxed{m}$ ,

the likelihood is

$$p(\vec{y} | X, w) \sim N(X^T w, \sigma_m^2 I)$$

Adding PDF's is different - that is a mixture model

2) product of Gaussian r.v. is a Gaussian r.v

3) conditional, marginal of multivariate normal is still normal  
cf. A.2 [RW]

We encode the prior on  $w$  as

$$w \sim N(0, \Sigma_p)$$

↑ Normal distr.  
aka Gaussian  $N(\mu, \Sigma)$

today, " $\Sigma$ " is a matrix usually  
except  $\sum_{i=1}^m$  is a sum

Given data  $(y, X)$ , how to update the distribution on  $w$ ?

Bayes' Rule

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}}, \quad p(w|y, X) = \frac{p(y|X, w) p(w)}{p(y|X)}$$

← independent of  $w$ ,  
so can ignore

known Gaussian

known Gaussian

product of two Gaussians is a Gaussian  
so with some computation, the posterior of  $w$  is

$$p(w|X, y) \propto N(\bar{w}, A^{-1})$$

$$\text{with } A := \frac{1}{\sigma_m^2} X X^T + \Sigma_p^{-1}$$

$$\bar{w} := \frac{1}{\sigma_m^2} A^{-1} X y$$

So we have a full distribution on  $w$  (more than usual PAC learning)!

If we need to pick out a single  $w$  to get a "point estimate",  
then usual idea is to pick the median,  
ie. maximum a posteriori (MAP) estimate.

For a Gaussian, median = mean, so MAP is  $\bar{w}$

To make a prediction at a new point  $x_*$ , again we get a full distribution

$$p(f(x_*) | X_*, X, y) = \int p(f(x_*) | x_*, w) \cdot p(w | X, y) dw$$

$$\sim N(\langle \bar{w}, x_* \rangle, \langle x_*, A^{-1} x_* \rangle)$$

and again, to get a point estimate, take the median/mean,

$$\text{so } \langle \bar{w}, x_* \rangle$$

alternative derivation: if we posit that  $f(x_*)$  is jointly Gaussian

with  $y$ ,  $\begin{bmatrix} y \\ f(x_*) \end{bmatrix} \sim N\left(\begin{bmatrix} \mu \\ \mu_{x_*} \end{bmatrix}, \begin{bmatrix} \Sigma & \Sigma_{x_*} \\ \Sigma_{x_*}^T & \Sigma_{x_* x_*} \end{bmatrix}\right)$

unknown

then  $f(x_*) | x_*, X, y \sim N(\mu_x, \tilde{\Sigma})$

w/ expressions for  $\mu_x$  and (see eq. 15.7  
for  $\tilde{\Sigma}$  via Schur complement to 15.9 Murphy)

repeating this for kernel version

replace  $X$  of  $\boxed{m}$  with  $\Phi$  of  $\boxed{m}$

$f(x) = \langle w, \phi(x) \rangle$  is model

almost everything stays the same, just replace  $X$  with  $\Phi$

i.e., now  $A = \frac{1}{\sigma_m^2} \underbrace{\Phi \Phi^T}_{\dim(H) \times \dim(H)} + \Sigma_p^{-1}$

$$\bar{w} = \frac{1}{\sigma_m^2} A^{-1} \Phi y \in H$$

$$f(x_*) | x_*, X, y \sim N(\langle \bar{w}, \phi(x_*) \rangle, \langle \phi(x_*), A^{-1} \phi(x_*) \rangle)$$

... though computing  $A^{-1}$  is awkward since it is  $\dim(H) \times \dim(H)$

trick is, again, the matrix inversion lemma

$$\Rightarrow \bar{w} = \sum_p \Phi (K + \sigma_m^{-2} I)^{-1} y$$

$$A^{-1} = \sum_p -\sum_p \Phi (K + \sigma_m^{-2} I)^{-1} \Phi^T \Sigma_p$$

$$\text{where } K := \Phi^T \sum_p \Phi \in \mathbb{R}^{m \times m}$$

and again we can use the kernel trick: all the  $\Phi$  or  $\phi$  terms combine to be in the form

$$\Phi^T \sum_p \Phi \text{ or } \phi(x_*)^T \sum_p \Phi \text{ or } \phi(x_*)^T \sum_p \phi(x_*)$$

So... define the kernel / covariance function

$$K(x, x') = \langle \phi(x), \sum_p \phi(x') \rangle \\ = \langle \psi(x), \psi(x') \rangle \text{ if } \psi(x) = \sum_p \phi(x)$$

Recall that our main choices are

(1) prior distribution on  $w$

(and we've fixed this to be a centered normal,  
so main choice is  $\sum_p$ )

(2) nonlinear map  $\phi$

... so you can see that equivalently, just specify the Kernel

## Function Space pt. of view

Def A Gaussian Process (GP) is a collection → ex: in time-series, indexed by  $\mathbb{R}$   
of random variables, any finite number of which have a joint  
Gaussian distribution → in spatial-stat, indexed by  $\mathbb{R}^2$  or  $\mathbb{R}^3$   
(implicitly imposes a consistency requirement or "marginalization property")

$f \sim GP(m, K)$      $f: X \rightarrow \mathbb{R}$  a function  
 $m: X \rightarrow \mathbb{R}$  the mean function  
 $K: X \times X \rightarrow \mathbb{R}$  is a pos.semidef kernel  
(covariance)

⇒ (1)  $\forall x, \mathbb{E} f(x) = m(x)$   
(typically set  $m(x) \equiv 0$ )

(2)  $\text{Cov}(f(x), f(x')) = K(x, x')$

(3)  $\forall m, \forall \{x_1, \dots, x_m\} \subseteq X,$   
 $(f(x_1), \dots, f(x_m)) \sim N(\vec{\mu}, K)$

$$\vec{\mu}_i := m(x_i), K_{ij} = k(x_i, x_j)$$

Given  $K$ , we saw above we can characterize the prediction,  
just need some  $O(m^3)$  computation (invert a matrix ...)

In practice, usually want to tune hyperparameters of the kernel

e.g. interpretation as smoothness of  $f$ ,  
or decorrelation length-scale

You can solve for hyperparameters in a principled way (but issues w/ nonconvexity) in  $O(n^3)$  too

## Comparison to Kernel-SVM

for simplicity, use trivial kernel,  $\Psi = I$ ,  
 $K(x, x') = \langle x, x' \rangle$

$$(\text{Soft-})\text{SVM} \text{ solves } \min_w \frac{1}{2} \|w\|^2 + \lambda^{-1} \sum_{i=1}^m [1 - y_i \langle w, x_i \rangle]_+$$

GP does more than just give a single pt. estimate, but  
if you used GP to get the MAP estimate, it is

$$\arg \min_w \frac{1}{2} \|w\|^2 + \sum_{i=1}^m -\log(p(y_i | \langle w, x_i \rangle))$$

where  $p(y | \langle w, x \rangle)$  is the likelihood that we already discussed

and usual choice is  $y | w, x \sim N(\langle w, x \rangle, \sigma_m^2)$

You can change this, though there is no likelihood that recovers the hinge loss (and hence SVM)

So, both GP & SVM can work w/ kernels, and both give

$$\text{estimates of the form } \arg \min_w \frac{1}{2} \|w\|^2 + \sum_i \text{loss}(w, y_i, x_i)$$

- SVM is based on minimization
- GP minimization can be done in closed form, and  
you recover an entire posterior distribution  
(but using GP for classification is messier)

Both Rasmussen & Williams' book and Murphy's book discuss more connections

- e.g., - GP vs neural networks
  - Smoothing splines as special case of GP,  
usually for 1D or 2D data