

Recap

• $\text{Rad}_S(H)$ = $\mathbb{E}_{\sigma} \sup_{h \in H} \sigma \cdot \begin{bmatrix} h(x_1) \\ \vdots \\ h(x_m) \end{bmatrix}$
Rademacher complexity

$$\sigma = [\sigma_1, \dots, \sigma_m]^T \in \{-1, 1\}^m$$

with $P_{\sigma} \gamma_2, \sigma_i = 1$ iid

$$S = \{(x_i, y_i)\}_{i \in [m]}$$

• Several notions of function complexity $\nearrow R(h)$

$$\text{generalization gap} = \underbrace{\mathbb{E}_{(x,y) \sim \mathcal{D}} \ell(x, y, h(x))}_{\hat{R}_S(h)} - \frac{1}{m} \sum_{i=1}^m \ell(x_i, y_i, h(x_i))$$

$$\leq \text{Rad}_S(H) + O\left(\sqrt{\frac{\log 1/\delta}{2m}}\right)$$

with $P_{\sigma} \geq 1 - \delta$ (over $S \sim \mathcal{D}^m$)

• $\text{Rad}_S(H) \leq \frac{\Lambda \|G\|_F}{m}$
RKHS

(Proof : last time)

$$\Lambda = \sup_{\alpha} \|\alpha\|$$

m : # samples

G : Gram matrix evaluated on S .

- kernel regression: why is it effective?
- Regularization term: $Y^* Y$ and its effect on Fourier coefficients of $f \in \mathcal{H}$.

- Complexity of \mathcal{H} ($\text{Rad}_S(\mathcal{H})$)

- $x \mapsto \kappa(x, \cdot)$ (RKHS map)

$x \mapsto \{\sqrt{\lambda_j} \psi_j(x)\}_j$ (Mercer map)

dimension of feature space

\Leftrightarrow dim of kernel regression
small when $\lambda_j \rightarrow 0$

Bayesian methods Gaussian process regression

Want to learn h s.t. $h(x_i) \approx y_i$
 $i \in [m]$

$$P(h | X, y)$$

What is the probab. of finding h given
 $S = \{(x_i, y_i)\}_{i \in [m]}$

- Beyond pointwise prediction, want uncertainties
- logistic regression

Bayes' rule

$$P(h | X, y) = \frac{\text{Prior}(h) \times P(X, y | h)}{P(X, y)}$$

$P(X, y | h)$: like likelihood
model evidence

$$P(X, y) = \int P(X, y | h) d\text{Prior}(h)$$

- $X, y = \{(x_i, y_i)\}_{i=1}^m$ iid

$$P(X, y | h) = \prod_{i=1}^m P(y_i - h(x_i))$$

- how to choose a prior on h ?

Gaussian process: is a function $x \rightarrow h(x)$

s.t. for every m , $\{h(x_1), h(x_2), \dots, h(x_m)\}$

is multivariate normally distributed.

$$\begin{bmatrix} h(x_1) \\ h(x_2) \end{bmatrix} \sim \frac{1}{\sqrt{2\pi} \sqrt{\det C}} e^{-\frac{1}{2} [a \ b]^T C^{-1} [a \ b]}$$

$$a = h(x_1) - \mu_1$$

$$b = h(x_2) - \mu_2$$

$$\begin{bmatrix} h(x_1) \\ \vdots \\ h(x_m) \end{bmatrix} \sim \frac{1}{(2\pi)^{m/2} \sqrt{\det C_m}} e^{-\frac{1}{2} [a_1 \dots a_m]^T C_m^{-1} [a_1 \dots a_m]}$$

Ex: C_m is PD.

$$C_{ij} = \mathbb{E}(h(x_1) h(x_2)) - \mathbb{E}h(x_1) \mathbb{E}h(x_2)$$

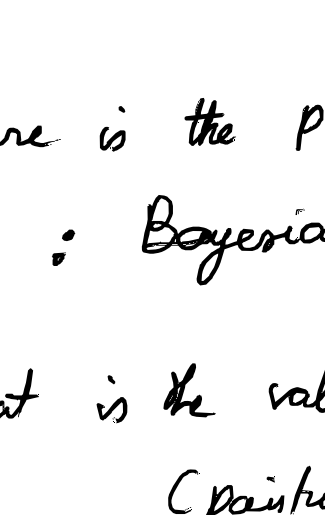
Want to sample $S = \{(x_i, y_i)\}_{i=1}^m$

$$P([h(x_1), h(x_2), \dots, h(x_m)] | S)$$

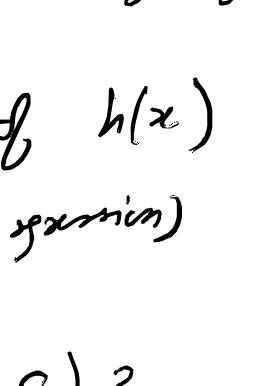
Bayes' rule

$$= \frac{P(S | h(x_1), \dots, h(x_m)) \times P(h(x_1), \dots, h(x_m))}{P(S)}$$

Under Gaussian process assumption
if h is GP.



$$x_i = [x^{(1)}(x_i) \ x^{(2)}(x_i)]$$



$$(x, y) \sim \mathcal{D}$$

Posterior

$$P([h(x_1), \dots, h(x_m)] | S) = \frac{P(S | h(x_1), \dots, h(x_m))}{P(S)} \times P(h(x_1), \dots, h(x_m))$$

$$P(h(x_1), \dots, h(x_m)) = \frac{e^{-\frac{1}{2} (h(x_1), \dots, h(x_m))^T G^{-1} (h(x_1), \dots, h(x_m))}}{(2\pi)^{m/2} \sqrt{\det G}}$$

$$\text{e.g. } C(h(x_1), h(x_2)) = K(x_1, x_2)$$

Where is the posterior maximized?

: Bayesian view of regression

What is the value of $h(x)$ at a new x ?

(pointwise regression)

What is $\max_h P(h(x) | S)$?

$$P(h(x) | S) = \frac{P(S | h) P(h)}{P(S)}$$

$$P\left(\begin{bmatrix} h(x_1) \\ \vdots \\ h(x_m) \end{bmatrix}\right) = \text{M.V. G.}$$

$$P(S | h) = \prod_{i=1}^m P(y_i - h(x_i))$$

e.g. Gaussian likelihood

$$y_i = h(x_i) + \epsilon_i$$

(labels/observations)

$$\epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

$$\max_h P(h | S) = \max_h P(S | h) P(h)$$

$$\arg\max_h P(h(x) | S) : \text{Maximum a posteriori estimate}$$

log ↑

$$\Rightarrow \arg\max_h P(h | S) = \arg\max_h \log P(h | S)$$

$$= \arg\max_h \log P(S | h) + \log P(h)$$

Under GP assumption for prior and Gaussian assumption on likelihood,

$$= \arg\max_h - \sum_{i=1}^m \frac{(y_i - h(x_i))^2}{2\sigma^2} - \frac{1}{2} (h(x_1), \dots, h(x_m))^T G^{-1} (h(x_1), \dots, h(x_m))$$

$$= \arg\min_{h \in \mathcal{H}} \sum_{i=1}^m \frac{(y_i - h(x_i))^2}{2\sigma^2} + \frac{1}{2} (h(x_1), \dots, h(x_m))^T G^{-1} (h(x_1), \dots, h(x_m))$$

Use representer theorem:

$$h(x) = \sum_{i=1}^m \alpha_i K(x, x_i)$$

$$\Rightarrow \begin{bmatrix} h(x_1) \\ \vdots \\ h(x_m) \end{bmatrix} = G \alpha \quad \alpha \in \mathbb{R}^m$$

Applying same old kernel trick

$$\min_{\alpha \in \mathbb{R}^m} \frac{1}{2\sigma^2} \|Y - G\alpha\|^2 + \alpha^T G \alpha$$

$$Y = [y_1, \dots, y_m]^T$$

$$\alpha = (G + \sigma^2 I)^{-1} Y$$

Bayesian view \Rightarrow automatic regularization

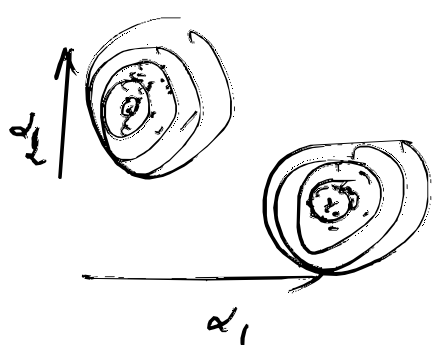
Solving for MAP estimate:

New want: samples from Bayesian posterior

• Bayesian inference:

$$\text{Posterior}(\alpha | S) = \frac{\text{Prior}(\alpha) \times P(S|\alpha)}{P(S)}$$

$$\nabla \log \text{Posterior}(\alpha | S) = \nabla \log \text{prior}(\alpha) + \nabla \log P(S|\alpha)$$



Score function:
 $\nabla \log \text{density} : \mathcal{X} \rightarrow \mathbb{R}^d$

Setting of B.I: (typically)

Score function of target: known

• Settings for sampling from target distributions

Generative modeling:

Want samples from prob. dist. Π .

Given: $\alpha_1, \dots, \alpha_m \sim \Pi$.

m samples



...

Algorithms for sampling and generative modeling that use kernels.

Generative modeling: GAN

MMD GAN

Sampling: Stein Variational Gradient descent

Integral probability metrics

P, Q are two probabi. distributions

$$Q, P: \mathcal{X} \rightarrow \mathbb{R}^+$$

$$d_f(P, Q) = \sup_{f \in F} \left| \mathbb{E}_{X \sim P} f(X) - \mathbb{E}_{X \sim Q} f(X) \right|$$

- F : Lipschitz functions with Lipschitz constant 1.

i.e.

$$\|f(x) - f(y)\| \leq \|x - y\| \quad \forall f \in F.$$

$d_f(P, Q)$: Wasserstein-1 distance

- F : RKHS
and $\|f\|_H \leq 1$
- $d_f(P, Q) :=$
Maximum
mean discrepancy