

SUPERVISED LEARNING

$\underline{x}, \underline{y} = (x_i, y_i), i=1, \dots, N$ observations (input, output)

Assume $p(y|x) = p(y|x, \theta)$ parametric

MAXIMUM LIKELIHOOD:

$$\theta_{ML} = \arg \max_{\theta} p(\underline{y} | \underline{x}, \theta)$$

LINEAR REGRESSION

Assume $p(y|x, \theta) = \mathcal{N}(y | f(x, \omega), \beta^{-1})$, i.e. $y = f(x, \omega) + \epsilon, \epsilon \sim \mathcal{N}(0, \beta^{-1})$

and $f(x, \omega) = \omega_0 \phi_0(x) + \dots + \omega_{M-1} \phi_{M-1}(x)$, $\Phi = (\phi_0, \dots, \phi_{M-1})$ basis func. $\left\{ \begin{array}{l} \text{monomials} \\ \text{gaussian RBF} \\ \text{sigmoid} \\ \vdots \end{array} \right.$

As $\log p(y|x, \theta) \propto -E_D(\omega) = -\frac{1}{2} \sum_{n=1}^N (y_n - \omega^T \phi(x_n))^2$, and

$$\nabla_{\omega} E_D(\omega) = \sum_{n=1}^N (y_n - \omega^T \phi(x_n)) \phi^T(x_n) = 0 \Rightarrow \omega_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \underline{y}, \quad \Phi_{ij} = \phi_j(x_i)$$

REGULARIZATION

solve $\arg \min_{\omega} E_D(\omega) + \lambda E_{\omega}(\omega)$

$$\text{where } E_{\omega}(\omega) = \begin{cases} \frac{1}{2} \|\omega\|_2^2 & (\text{Ridge}) \\ \frac{1}{2} \|\omega\|_1 & (\text{Lasso}) \end{cases}$$

- REGULARIZATION IS BIAS

→ PRIOR BELIEFS \equiv BIASING

$$p(w|\alpha) = \mathcal{N}(w|0, \alpha I)$$

force parameters to be small
LIKELIHOOD

PRIOR

$$p(w|\underline{x}, \underline{y}, \alpha, \mathcal{N}) = \frac{p(\underline{y}|\underline{x}, w, \alpha, \mathcal{N}) p(w|\alpha)}{p(\underline{y}|\underline{x}, \alpha, \mathcal{N})}$$

POSTERIOR

MARGINAL LIKELIHOOD

$$\log p(w|\underline{x}, \underline{y}, \alpha, \mathcal{N}) = -\frac{\mathcal{N}}{2} \sum_{i=1}^{\mathcal{N}} (y_i - \underbrace{w^T \phi(x_i)}_{\text{quadratic in } w})^2 - \alpha \underbrace{w^T w}_{\text{quadratic in } w} + \text{const}$$

$$P(\omega | \underline{x}, \underline{y}, \alpha, \beta) = \mathcal{N}(\omega | m_N, S_N)$$

$$m_N = \beta S_N \Phi^T \underline{y}$$

$$S_N^{-1} = \alpha I + \beta \Phi^T \Phi$$

IF we use a general prior

$$P(\omega | m_0, S_0)$$

↳

$$P(\omega | \underline{x}, \underline{y}, \alpha, \beta) = \mathcal{N}(\omega | m_N, S_N)$$

$$m_N = S_N [S_0^{-1} m_0 + \beta \Phi^T \underline{y}]$$

$$S_N^{-1} = S_0^{-1} + \beta \Phi^T \Phi$$

CINER
BAYESIAN REGRESSION

GAUSSIAN PRIOR

↳

GAUSSIAN POSTERIOR

↳

ANALYTICALLY
COMPUTABLE!

PREDICTIVE DISTRIBUTION

TAKE A NEW x .

$$P(y | x, \underline{x}, \underline{y}, \alpha, \beta) = \int \underbrace{P(y | x, \omega, \alpha, \beta)}_{\substack{\text{GAUSSIAN} \\ \text{LINEAR ON } \omega}} \underbrace{P(\omega | \underline{x}, \underline{y}, \alpha, \beta)}_{\text{GAUSSIAN}} d\omega$$

GAUSSIAN

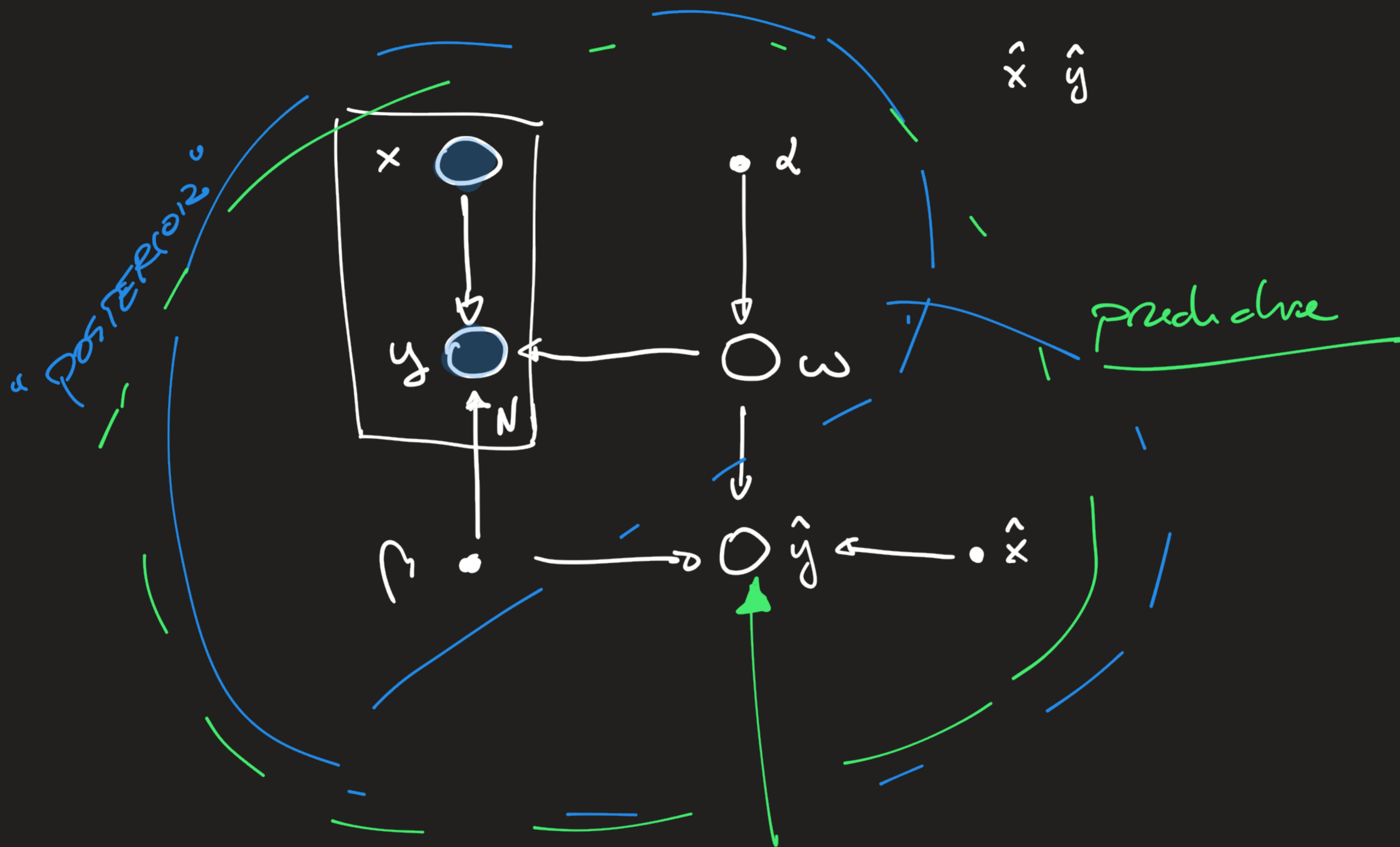
GAUSSIAN

$$= \mathcal{N}(y / w_N^T \phi(x), \sigma_N^2(x))$$

$$\sigma_N^2(x) = \frac{1}{\beta} + \phi(x)^T S_N \phi(x)$$

$$\sigma_N^2(x) \geq \sigma_{N+1}^2(x), \quad \sigma_N^2(x) \rightarrow \frac{1}{\beta}, \quad N \rightarrow \infty$$

PGM of BLR!



MARGINAL LIKELIHOOD

$$P(\underline{y} | \underline{x}, \alpha, \mathcal{M}) = \int P(\underline{y} | \underline{x}, \omega, \alpha, \mathcal{M}) P(\omega | \alpha) d\omega$$

MARGINAL LIKELIHOOD

as the posterior $P(\omega | \underline{y}, \underline{x}, \alpha, \mathcal{M}) = \frac{P(\underline{y} | \underline{x}, \omega, \alpha, \mathcal{M}) P(\omega | \alpha)}{P(\underline{y} | \underline{x}, \alpha, \mathcal{M})}$

The Bayesian way for α and $\beta \Rightarrow P(\alpha, \mathcal{M})$ HYPERPRIOR

$$P(\alpha, \mathcal{M} | \underline{y}, \underline{x}) \propto P(\underline{y} | \underline{x}, \alpha, \mathcal{M}) P(\alpha, \mathcal{M})$$

(a) UNINFORMATIVE PRIOR
!!
 $P(\alpha, \mathcal{M})$ IS CONSTANT

(b) POSTERIOR IS SHARPLY PEAKED AROUND MAP \Rightarrow

$$P(\alpha, \mathcal{M} | \underline{y}, \underline{x}) \approx \delta_{\text{MAP}(\alpha, \mathcal{M})}$$

(a) + (b) \Rightarrow we can FIX α, \mathcal{M} TO M.L. \Rightarrow WE CAN MAXIMIZE THE MARGINAL LIKELIHOOD.

$$\log p(\underline{y} | \underline{x}, \alpha, \beta) = \frac{M}{2} \log \alpha + \frac{N}{2} \log \beta - E(m_N) - \frac{1}{2} \log |S_N^{-1}| - \frac{N}{2} \log 2\pi$$

$$E(m_N) = \frac{\beta}{2} \|y - \Phi m_N\|^2 + \frac{\alpha}{2} m_N^T m_N$$

$$\left[m_N = \beta S_N \Phi^T y \ ; \ S_N^{-1} = \alpha I + \beta \Phi^T \Phi \right]$$

FIX-POINT METHOD

($\nabla \log p(\underline{y} | \underline{x}, \alpha, \beta) = 0$, and derive F.P. equations)

$$\begin{cases} \alpha = g_\alpha(\alpha, \beta) \\ \beta = g_\beta(\alpha, \beta) \end{cases}$$

ALGORITHM: Fix α_0, β_0 , then compute

$$\alpha_{n+1} = g_\alpha(\alpha_n, \beta_n)$$

$$\beta_{n+1} = g_\beta(\alpha_n, \beta_n)$$

iterate until convergence

$$\|\alpha_{n+1} - \alpha_n\| + \|\beta_{n+1} - \beta_n\| < \epsilon$$

$$\log p(\underline{y} | \underline{x}, \alpha, \beta) = \frac{M}{2} \log \alpha + \frac{N}{2} \log \beta - E(m_N) - \frac{1}{2} \log |S_N^{-1}| - \frac{N}{2} \log 2\pi$$

$$E(m_N) = \frac{\beta}{2} \|\underline{y} - \Phi \cdot m_N\|^2 + \frac{\alpha}{2} m_N^T m_N$$

$$[S_N^{-1} = \alpha I + \beta \Phi^T \Phi]$$

$\beta \Phi^T \Phi$ and compute eigenvalues $\lambda_i > 0$
do not depend on α

$$|S_N^{-1}| = \prod_{i=0}^{M-1} (\alpha + \lambda_i)$$

$$\frac{\partial \log |S_N^{-1}|}{\partial \alpha} = \frac{\partial}{\partial \alpha} \sum \log (\alpha + \lambda_i) = \sum_{i=0}^{M-1} \frac{1}{\alpha + \lambda_i}$$

$$\frac{\partial \lambda_i}{\partial \beta} = \frac{\lambda_i}{\beta}$$

$$\Rightarrow \alpha = \left[\frac{\gamma}{m_N^T m_N} \right] = g_\alpha(\alpha, \beta), \text{ where } \gamma = \sum_{i=0}^{M-1} \frac{\lambda_i}{\alpha + \lambda_i}$$

$$\frac{1}{\beta} = \frac{1}{N - \gamma} \sum_{n=1}^N (y_n - m_N^T \phi(x_n))^2 =: \frac{1}{g_\beta(\alpha, \beta)}$$

$$\gamma = \sum_{i=0}^{M-1} \frac{\lambda_i}{\alpha + \lambda_i}$$

$\lambda_i =$ eigenvalues of $\Phi^T \Phi$

λ_i gives info on w_{ML}^i

SMALL $\lambda_i \Rightarrow$ LARGE UNCERTAINTY on w_{ML}

LARGE $\lambda_i \Rightarrow$ SMALL UNCERTAINTY

EFFECTIVE # OF PARAMETERS

$$\left(\begin{array}{l} \lambda_i \ll \alpha \Rightarrow \frac{\lambda_i}{\lambda_i + \alpha} \approx 0 \\ \lambda_i \gg \alpha \Rightarrow \frac{\lambda_i}{\lambda_i + \alpha} \approx 1 \end{array} \right)$$

for $N \gg M$ $\gamma \approx M$

(MARGINAL LIKELIHOOD EVIDENCE)

BAYESIAN MODEL COMPARISON

M_1, M_2 ARE DIFFERENT MODELS

WHICH ONE IS THE BEST TO EXPLAIN DATA $\mathcal{D} = \{ \underline{x}, \underline{y} \}$

$P(M_j)$ PRIOR ON THE MODELS

$P(M_j | \mathcal{D}) =$ POSTERIOR

$$P(\mathcal{D} | M_j) \cdot P(M_j)$$

$$\frac{P(\mathcal{D} | M_j) \cdot P(M_j)}{\sum_j P(\mathcal{D} | M_j) \cdot P(M_j)}$$

MARGINAL LIKELIHOOD

1) MODEL AVERAGING.

$$\underbrace{P(y|x, D)}_{\text{prediction}} = \sum_{\theta} \underbrace{P(y|x, D, M_{\theta})}_{\text{likelihood}} \cdot P(M_{\theta}|D)$$

2) CHOOSE THE BEST MODEL

$$\frac{P(D|M_1)}{P(D|M_2)} \equiv \text{BAYES FACTOR}$$

choose M_{θ} with largest Bayes Factor

because $\int P(D|M_2) \log \frac{P(D|M_1)}{P(D|M_2)} dD > 0$ if M_1 is the True model

EQUIVALENT KERNEL

$$p(y | x, \underline{x}, \underline{y}, \alpha, \beta) = \mathcal{N}(y | \underbrace{m_N^T \phi(x)}, \sigma^2(x))$$

$$y(x, m_N) = m_N^T \phi(x) = \beta \phi(x)^T S_N \begin{matrix} \Phi^T \underline{y} \\ \vdots \end{matrix} = \sum_{n=1}^N \beta \phi^T(x) S_N \phi(x_n) y_n$$
$$= \sum_{n=1}^N k(x, x_n) y_n$$

$$k(x, x') = \beta \phi^T(x) S_N \phi(x') \quad \text{or}$$

equivalent kernel

$$\text{cov}[y(x), y(x')] = \phi^T(x) S_N \phi(x') = \underline{\beta^{-1} k(x, x')}$$

$$k(x, x') = \Psi^T(x) \Psi(x') \quad \Psi(x) = \beta^{1/2} S_N^{1/2} \phi(x) \quad \text{or}$$