

Probabilistic Models

Likelihood

$$p(\mathcal{D} | \underline{w}) = \prod_{n=1}^N p(\underline{x}^{(n)} | \underline{w}) \underbrace{p(y^{(n)} | \underline{x}^{(n)}, \underline{w})}$$

Often "1" in supervised models.
We know $\underline{x}^{(n)}$

Data $\{ \underline{x}^{(n)}, y^{(n)} \}$

Examples: $N(y^{(n)}; f(\underline{x}^{(n)}; \underline{w}), \sigma_y^2)$

For each can minimize negative log likelihood (+ regularizer)

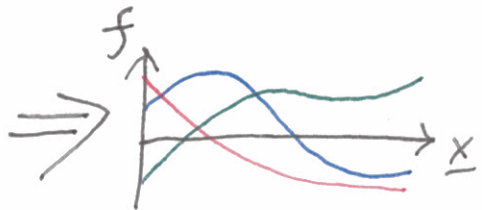
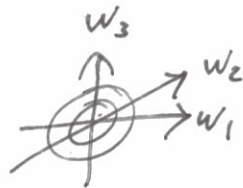
$N(y^{(n)}; f(\underline{x}^{(n)}; \underline{w}), \sigma_y^2)$

Bernoulli ($y^{(n)}; \sigma(f(\underline{x}^{(n)}; \underline{w}))$)

Robust versions

...

Prior $p(\underline{w})$



Prior over functions

Bayesian Inference

✓ $\Phi, \sigma_y^2, \sigma_w^2, \dots$

Posterior

Model assumptions. Can mention everywhere. Often don't.

$$P(\underline{w} | D, M) = \frac{P(D | \underline{w}, M) P(\underline{w} | M)}{P(D | M)} \propto P(D | \underline{w}, M) P(\underline{w} | M)$$

$P(D | M)$

Marginal Likelihood (dropping "M")

$$= \int P(D | \underline{w}) p(\underline{w}) d\underline{w}$$

Not Bayesian

$$\operatorname{argmax}_{\underline{w}} P(\underline{w} | D) = \text{"MAP" params}$$

$$= \operatorname{argmax}_{\underline{w}} \left[\log P(D | \underline{w}) + \underbrace{\log p(\underline{w})}_{L2 \text{ regularizer}} \right]$$

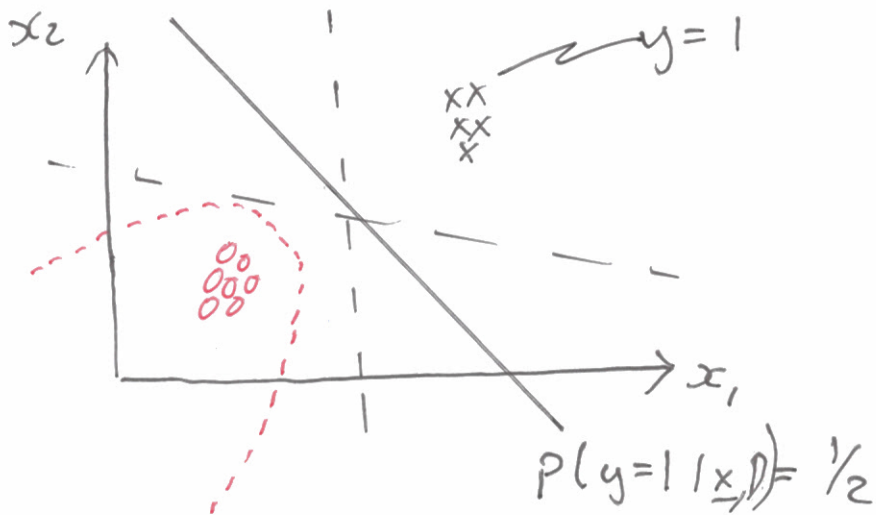
L2 regularizer

Logistic Regression

$$P(y=1 | \underline{x}, \underline{w}) = \sigma(\underline{w}^T \underline{x}) = \frac{1}{1 + e^{-\underline{w}^T \underline{x}}}$$

$$P(y | X, \underline{w}) = \prod_{n=1}^N \sigma(\underline{w}^T \underline{x}^{(n)} z^{(n)})$$

↑
= $P(D | \underline{w})$ (ish) $2y-1$



$$P(y=1 | \underline{x}, D) = 0.1$$

Make predictions

$$p(y^* | \underline{x}^*, D) = \int p(y^* | \underline{w}, \underline{x}^*, D) d\underline{w}$$

test output test input

$$= \int \underbrace{p(y^* | \underline{w}, \underline{x}^*)}_{\text{Prediction for weights } \underline{w}} \underbrace{p(\underline{w} | D)}_{\text{Posterior over weights}} d\underline{w}$$

Monte Carlo

$$p(y^* = 1 | \underline{x}, D) = \mathbb{E}_{p(\underline{w} | D)} [p(y^* = 1 | \underline{w}, \underline{x}^*)]$$

$$= \mathbb{E}_{p(\underline{w} | D)} [\sigma(\underline{w}^T \underline{x}^*)]$$

$$\approx \frac{1}{S} \sum_{s=1}^S \sigma(\underline{w}^{(s)T} \underline{x}^*)]$$

$$\underline{w}^{(s)} \sim p(\underline{w} | D).$$

How? $\underline{w}^{(s)} \sim p(\underline{w} | D)$

Approximately with "MCMC"

(not this course)

Importance Sampling

$$\int g(x) p(x) dx = \int \left[g(x) \frac{p(x)}{q(x)} \right] q(x) dx$$

$$\mathbb{E}_p [g(x)] = \mathbb{E}_q [\dots] \quad \begin{array}{l} \text{[If } q(x) \neq 0 \\ \text{if } p(x) \neq 0] \end{array}$$

For logistic regression

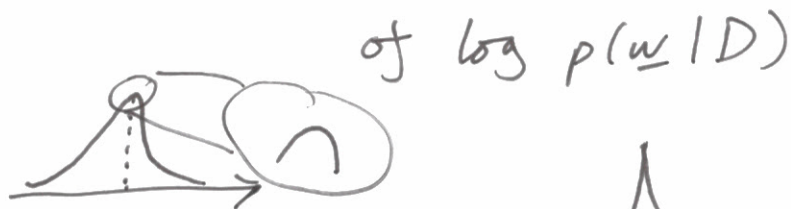
$$\underline{w}^{(s)} \sim q(\underline{w}) = \text{prior } p(\underline{w})$$

$$p(y^* | \underline{x}^*, D) \propto \frac{\frac{1}{S} \sum_{s=1}^S p(y^* | \underline{x}^*, \underline{w}^{(s)}) p(D | \underline{w}^{(s)})}{\sum_{s=1}^S p(D | \underline{w}^{(s)})}$$

Laplace Approximation

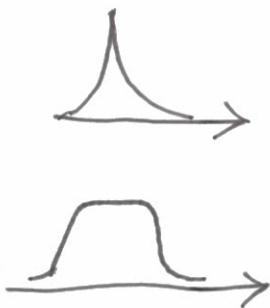
Fits a Gaussian to $p(\underline{w} | D)$

- Matching mode
- Matching curvature (2nd derivative)



$$\underline{w}^* = \underset{\underline{w}}{\operatorname{argmax}} p(\underline{w} | D)$$

MAP Parameters



$$= \underset{\underline{w}}{\operatorname{argmin}} -\log p(\underline{w} | D)$$

Numerical
methods

"Energy"

$$E(\underline{w}) = -\log p(\underline{w}, D)$$

① $w^* \equiv \underset{\underline{w}}{\operatorname{argmin}} E(\underline{w})$

② Find curvature:

$$1D \quad H = \left. \frac{\partial^2 E}{\partial w^2} \right|_{w=w^*}$$

In general

Hessian, $H_{ij} = \frac{\partial^2 E}{\partial w_i \partial w_j}$

Laplace Approximation

$$p(\underline{w} | D) \approx N(\underline{w}; \underline{w}^*, H^{-1})$$