

Week 10a Bayesian logistic regression

Recap:

Task: given data $D = \{x, y\}$, fit probabilistic model

Discriminative model, y given x $P(y | x, w)$

Generative model, modeling x and y directly $P(y, x | w)$

Then specify likelihood = $P(D | w) = P(x_{\text{train}}, y_{\text{train}} | w)$

Discriminative model: product over each datapoints m , then product over sigmoid, $z = [-1, 1] = 2y - 1$ (or other functions)

Generative model: product over each datapoints m , then product rule

$$\begin{aligned} \text{Likelihood} &= P(D | w) \\ p(D | w) &= \prod_{m=1}^n p(y^{(m)}, x^{(m)} | w) = \prod_{m=1}^n p(y^{(m)} | x^{(m)}, w) p(x^{(m)} | w) \\ p(D | w) &= \prod_{m=1}^n p(y^{(m)} | x^{(m)}, w) \\ &\stackrel{\text{DISC.}}{=} N(y^{(m)}; f(x^{(m)}, w), \sigma_y^2) \\ &\stackrel{\text{BINARY}}{=} \text{Bernoulli}(y^{(m)}; \sigma(f(x^{(m)}, w))) \end{aligned}$$

For discriminative logistic regression: why product over $P(y^{(m)} | w, x^{(m)})$
minimize negative log likelihood, find single w^*

Likelihood is a function of parameters, not a distribution over the parameters

likelihood is product over the probabilities according to the model of observing the datapoints

Assume data is identically distributed and independently drawn from this distribution (IID) (so get the normal distribution symbol)

$$\begin{aligned} \text{Likelihood} &= P(D | w) \\ p(D | w) &= \prod_n p(x^{(n)} | w) \cdot p(y^{(n)} | w, x^{(n)}) \\ &\quad \text{i "known" here} \\ &\quad \sigma(w^T x^{(n)} (2y^{(n)} - 1)) \end{aligned}$$

Bayesian approach: true set of weights are unknown, described by prob dist

predictions $P(y_{\text{test}} | x_{\text{test}}, D)$ are made considering all possible settings (sigmoids), weighted by how plausible they are given the training data (posterior $p(w | D)$)

$$\begin{aligned} p(y | x, D) &= \int p(y, w | x, D) dw = \\ &= \int p(y | x, w, D) p(w | x, D) dw \end{aligned}$$

$$p(w | D) \underset{\text{Posterior}}{=} \frac{p(w, D)}{p(D)} = \frac{p(D | w) p(w)}{\int p(D | w) p(w) dw}$$

$$P(D) = \int P(D | w) p(w) dw.$$

$$\begin{aligned}
 & \text{Predictions} \\
 p(y=1 | \underline{x}, D) &= \int p(y=1, \underline{w} | \underline{x}, D) d\underline{w} \\
 &= \underbrace{\int p(y=1 | \underline{w}, \underline{x}, D)}_{\sigma(\underline{w}^\top \underline{x})} \cdot \underbrace{p(\underline{w} | \underline{x}, D)}_{\text{posterior}} d\underline{w}
 \end{aligned}$$

Interpretation: average the predictive distributions $P(y | x, w)$ for different parameters weighted by how plausible those parameters are, $p(w | D)$

How to compute this integral?

(= expectation of sigmoid function $P(y_{\text{test}} | x_{\text{test}}, w)$ over posterior)

if prior * likelihood yields a parametric form that we know of (normal), can calculate integral in closed form

or approximate it by Gaussian (simpler form), optimize a “simpler surrogate”

or match moments (first moment = mean & second moment = variance)

or MAP (= normal logistic regression)

MAP (maximum probability inference):

Get point estimate

= just pick mode of the distribution (pdf curve highest)

= argmax of posterior $P(w|D)$, (= get most probable weight in posterior)

= argmax of $P(w, D)$, since $P(D)$ does not depend on parameters w

then chain rule, take log

=> Under certain assumptions (prior is simple independent gaussian, zero mean & scalar variance), prior = regularization term

=> MAP = L2 regularized MLE (logistic regression)

MLE $\underline{w}^* = \underset{\underline{w}}{\operatorname{argmax}} p(D \underline{w})$ $= \underset{\underline{w}}{\operatorname{argmin}} -\log p(D \underline{w}) + \lambda \ \underline{w}\ ^2$	MAP estimation “pick ‘mode of the distribution’” $b^* = \underset{b}{\operatorname{argmax}} p(b)$ $\underline{w}^{\text{MAP}} = \underset{\underline{w}}{\operatorname{argmax}} p(\underline{w} D) =$ $= \underset{\underline{w}}{\operatorname{argmax}} p(\underline{w}, D) =$ $= \underset{\underline{w}}{\operatorname{argmax}} p(D \underline{w}) p(\underline{w})$ $= \underset{\underline{w}}{\operatorname{argmax}} \underbrace{\log p(D \underline{w})}_{\frac{1}{2\sigma_w^2} \underline{w}^\top \underline{w}}$ + $\underbrace{\log p(\underline{w})}_{\frac{1}{2\sigma_w^2} \underline{w}^\top \underline{w}}$
--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmax}} [\log p(\mathbf{w} | \mathcal{D})] = \underset{\mathbf{w}}{\operatorname{argmax}} \left[\log P(\mathcal{D} | \mathbf{w}) - \frac{1}{2\sigma_w^2} \mathbf{w}^\top \mathbf{w} \right]$$

Advantages of bayesian:

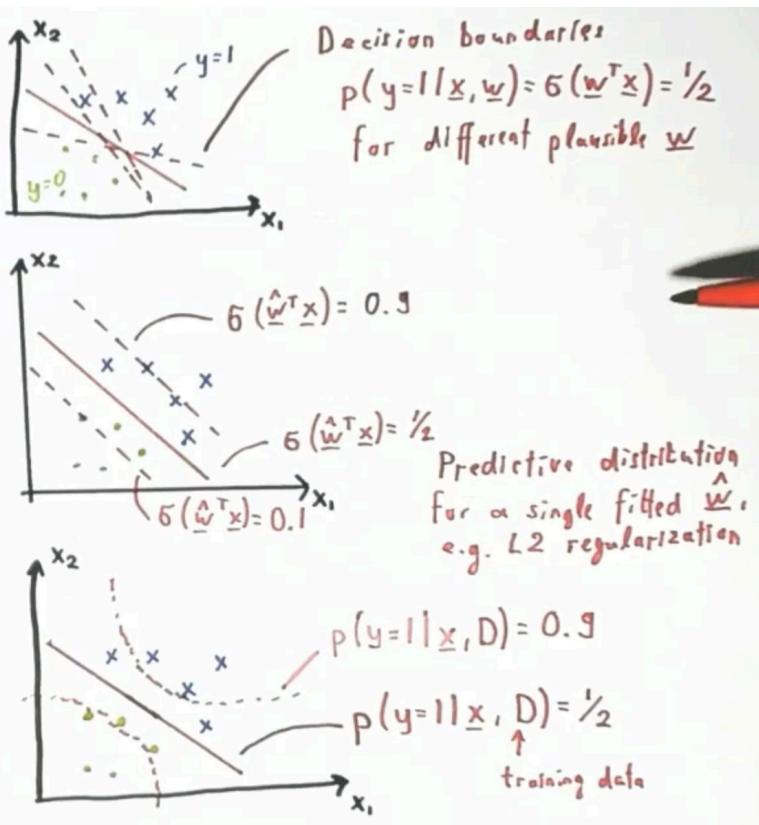
model uncertainty

look at marginal likelihood, make better use of available data (to choose model, also don't need validation set)

posterior represent our belief about the model

Represent uncertainty (in first figure)

Last figure: when further away from the training inputs, Contours curved, predictions less certain



Sketch logistic regression posterior:

posterior $p(w | D)$ is not Gaussian, but can approximate it with a Gaussian

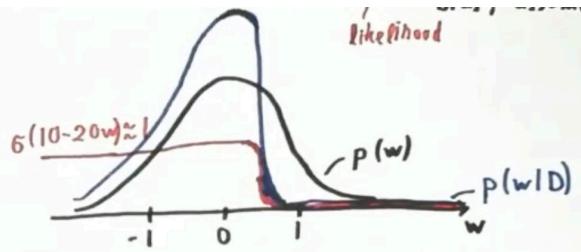
e.g. known prior, known 1 datapoint ($y = 1$ at $x = -20$, bias=10)

since weight isn't a large positive value, slide off positive region & renormalize (area under curve sums to 1)

$$p(w) \propto \mathcal{N}(w; 0, 1)$$

$$p(w | D) \propto \mathcal{N}(w; 0, 1) \sigma(10 - 20w)$$

Normal curve: prior; Red: sliced off, before normalize; Black: after normalized



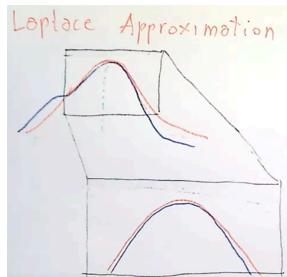
Since posterior = sigmoid multiplied,

every time multiply the posterior by a sigmoidal likelihood, softly carve away half of the weight space in some direction

Bayesian central limit theorem: after many datapoints, distribution over plausible weights looks gaussian

Week 10b Laplace approximation

approximate a distribution $P(w | D)$ (not exactly gaussian) with a Gaussian incremental improvement of the MAP approximation to Bayesian inference



matches mode (of Gaussian approximation = mode of posterior w^*)

matches the **curvature** of the log probability density at location i,j (Hessian)
=> goal: find mode & second derivative

Approximate posterior by Boltzmann distribution

positive, so e; integrate to one, so divide by normalization constant (partition function) Z

Energy E

$$Z = P(D); e^{-E(w)} = p(w, D)$$

$$P(w|D) = \frac{e^{-E(w)}}{Z} \underset{\substack{\text{partition function} \\ \text{norm. const.}}}{=} \frac{p(w, D)}{P(D)}$$

Energy E: complex function we try to approximate ($-\log P(w, D)$)

most probable weight w^*

= maximize posterior (MAP)

proportional to energy, drop other terms (Z and exp)

= minimum of energy (since exp is an increasing function):

(minimum = turning point, first derivative = 0, second derivative = curvature)

w^* = L2 regularization or MAP fit

$$E(w) = -\log p(w, D), \quad w^* = \arg \min_w E(w).$$

Take log of posterior (put in MAP)

By property of conditional probability, $P(D) = \text{some constant w.r.t } w$:

since log of distribution is closer to quadratic function (good for optimizer)

$$\log p(w | D) = \log p(w, D) + \text{const wrt. } w$$

Energy being approximated by

its value observed in local points $E(w^*)$

+ gradient, evaluated at local point

+ quadratic term (Taylor approximation)

$$\begin{aligned} E(w) &\approx E(w^*) + \cancel{(w - w^*)^T \frac{\partial E}{\partial w}(w^*)} + \\ &e^{-\left(\frac{1}{2} (w - w^*)^T H (w - w^*)\right)} \\ H_{ij} &= \left. \frac{\partial^2 E}{\partial w_i \partial w_j} \right|_{w^*} \\ \boxed{P(w | D) &\approx \mathcal{N}(w; w^*, H^{-1})} \end{aligned}$$

Hessian, second derivative at the optimum, $w = w^*$ (by MAP):

= curvature of the function

= how sharply the distribution is peaked in different directions

= how spread out the gaussian would be

$$H_{ij} = \left. \frac{\partial^2 E(w)}{\partial w_i \partial w_j} \right|_{w=w^*}$$

$E(w)$ approx. by Gaussian, with $N(\mu, \Sigma)$,
 energy = $-\log P(w, D)$ is:

$$E_N(w) = \frac{1}{2}(w - \mu)^\top \Sigma^{-1}(w - \mu),$$

So $w^* = \mu$ and $H = \Sigma^{-1}$, covariance is $\Sigma = H^{-1}$

Laplace approximation to the posterior distribution (by matching minimum & curvature):

When prior & likelihood are gaussian, L approx is exact

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

Approximate marginal likelihood $P(D | M)$ (= partition function Z)

Expand the Gaussian PDF:

$$p(w | \mathcal{D}) = \frac{p(w, \mathcal{D})}{P(\mathcal{D})} \approx N(w; w^*, H^{-1}) = \frac{|H|^{1/2}}{(2\pi)^{D/2}} \exp\left(-\frac{1}{2}(w - w^*)^\top H(w - w^*)\right)$$

In order to get normalizing constant $P(D)$, does not depends on w , could plug in any values for w ,

if plug in w^* ($w=w^*$), exp term disappears, very convenient

get $P(D)$:

get w^* by SGD on negative log likelihood with regularization term; = MAP
 DIFFERENTIATE: D in $P(D)$: training data; D in $D/2$: # of parameters w

$$\frac{p(w^*, \mathcal{D})}{P(\mathcal{D})} \approx \frac{|H|^{1/2}}{(2\pi)^{D/2}}, \quad P(\mathcal{D}) \approx \boxed{\frac{p(w^*, \mathcal{D})(2\pi)^{D/2}}{|H|^{1/2}}}.$$

Approximate $p(D)$ for different models, choose the model with the highest marginal likelihood $P(D | M)$

since it is approximation, can go wrong if approx. by gaussian is poor fit

Prediction

Given test input & observations, want to take our belief about different ws into account

but distribution in integral (eq bellow) is not gaussian anymore (different from Gaussian linear reg)

$$P(y | \mathbf{x}, \mathcal{D}) = \int P(y | \mathbf{x}, \mathbf{w}) p(\mathbf{w} | \mathcal{D}) d\mathbf{w}$$

$P(y | \mathbf{x}, \mathbf{w})$ = sigmoid, since logistic regression;
use a normal distribution to weight the sigmoid

$$\begin{aligned} P(y=1 | \mathbf{x}, \mathcal{D}) &\approx \int \sigma(\mathbf{w}^\top \mathbf{x}) \mathcal{N}(\mathbf{w}; \mathbf{w}^*, H^{-1}) d\mathbf{w} \\ &= \mathbb{E}_{\mathcal{N}(\mathbf{w}; \mathbf{w}^*, H^{-1})} [\sigma(\mathbf{w}^\top \mathbf{x})]. \end{aligned}$$

Can not compute high dim integral,
Goal: convert to 1D, transform sigma to gaussian (to get closed form solution)

since $\mathbf{w}^\top \mathbf{x}$ is scalar, convert to 1D integral
scalar $a = \mathbf{w}^\top \mathbf{x}$; $p(a)$ = gaussian;
 a is linear trans of the distribution of w by X
take the expectation over this scalar quantity

$$= \underbrace{\mathbb{E}_{p(a)} [\sigma(a)]}_{\mathcal{N}(a; \mathbf{w}^* \mathbf{x}, \mathbf{x}^\top H^{-1} \mathbf{x})}$$

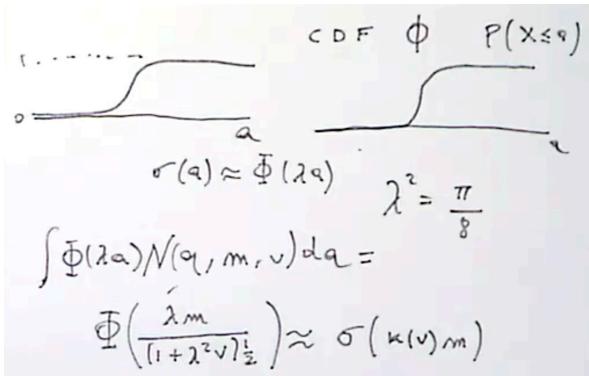
prediction = 1D integral, could compute numerically

$$\begin{aligned} P(y=1 | \mathbf{x}, \mathcal{D}) &\approx \mathbb{E}_{\mathcal{N}(a; \mathbf{w}^* \mathbf{x}, \mathbf{x}^\top H^{-1} \mathbf{x})} [\sigma(a)] \\ &= \int \sigma(a) \mathcal{N}(a; \mathbf{w}^* \mathbf{x}, \mathbf{x}^\top H^{-1} \mathbf{x}) da. \end{aligned}$$

Probit approximation: approx. sigmoid by gaussian CDF

Since we know how to compute and integrate out CDF of gaussian * gaussian
lambda: adjust the slopes

term inside sigmoid: kappa; $m = \mathbf{w}^* \mathbf{x}$



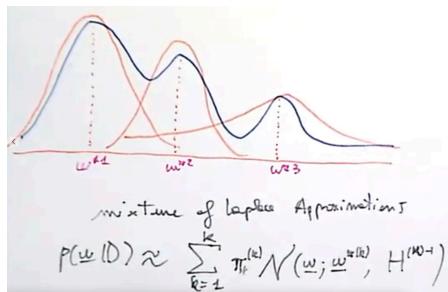
Kappa of covariance $\kappa(v)$:

predictions use the most probable or MAP weights w^*

BUT $w^* \mathbf{x}$ term will be scaled down by kappa when uncertain about activation
 \Rightarrow predictions will be less confident far from the data

$$P(y=1 | \mathbf{x}, \mathcal{D}) \approx \sigma(\kappa w^{*\top} \mathbf{x}), \quad \kappa = \frac{1}{\sqrt{1 + \frac{\pi}{8} \mathbf{x}^\top H^{-1} \mathbf{x}}}.$$

Mixture (sum) of Laplace approx. by coefficient pi



Is the Laplace approximation reasonable?

Laplace approximation failure

when distribution is not Gaussian, or multi-modal (posterior for neural net)

