Linear and Generalized Linear Models for Supervised Learning

CS772A: Probabilistic Machine Learning Piyush Rai

Plan today

- Probabilistic Linear Models for supervised learning (learning p(y|x))
- Focus will be the "discriminative" setting (assume a form and learn p(y|x) directly)
 - Probabilistic Linear Regression
 - Logistic and Softmax Classification
 - Generalized Linear Models for supervised learning



A discriminative model **Probabilistic Linear Regression** for regression problems

- Assume training data $\{x_n, y_n\}_{n=1}^N$, with features $x_n \in \mathbb{R}^D$ and responses $y_n \in \mathbb{R}$ Unknown to be estimated
- Assume y_n generated by a noisy linear model with wts $w = [w_1, \dots, w_n] \in \mathbb{R}^n$

$$y_n = w^{\mathsf{T}} x_n + \epsilon_n^{\mathsf{Gaussian noise drawn}}$$
 from $\mathcal{N}(\epsilon_n | 0, \beta^{-1})$

• Notation alert: β is the <u>precision</u> of Gaussian noise (and β^{-1} the <u>variance</u>)



Each weight assumed real-valued



Probabilistic Linear Regression

• For all the training data, we can write the above model in matrix-vector notation



- Linear Gaussian model and w is the Gaussian r.v. with $p(w|\lambda) = \mathcal{N}(w|0, \lambda^{-1}\mathbf{I}_D)$
- A simple "plate diagram" for this model would look like this (hyperparameters not shown in the diagram)
 White nodes denote unknown guantities are nodes denote





On compact notations..

• When writing the likelihood (assuming y_n 's are i.i.d. given \boldsymbol{w} and \boldsymbol{x}_n)

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^{N} \mathcal{N}(y_n | \mathbf{w}^{\mathsf{T}} \mathbf{x}_n, \beta^{-1})$$

$$= \mathcal{N}(\boldsymbol{y}|\boldsymbol{X}\boldsymbol{w},\beta^{-1}\boldsymbol{I}_N)$$

- Thus a product of N univariate Gaussians here (not always) is equivalent to an N-dim Gaussian over the vector $\mathbf{y} = [y_1, y_2, \dots, y_N]$
- We will prefer to use this equivalence at other places too whenever we have multiple i.i.d. random variables, each having a univariate Gaussian distribution

The Posterior



- The posterior over $oldsymbol{w}$ (for now, assume hyperparams $oldsymbol{eta}$ and $oldsymbol{\lambda}$ to be known)

$$p(w|y, X, \beta, \lambda) = \frac{p(w|\lambda)p(y|w, X, \beta)}{p(y|X, \beta, \lambda)} \propto p(w|\lambda)p(y|w, X, \beta)$$
Must be a Gaussian
due to conjugacy
Must be a Gaussian

$$abla p(w|y, \mathbf{X}, eta, \lambda) \propto \mathcal{N}(w|\mathbf{0}, \lambda^{-1}\mathbf{I}_D) imes \mathcal{N}(y|\mathbf{X}w, eta^{-1}\mathbf{I}_N)$$

Using the "completing the squares" trick (or linear Gaussian model results)

$$p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \beta, \lambda) = \mathcal{N}(\boldsymbol{\mu}_{N}, \boldsymbol{\Sigma}_{N})$$
where $\boldsymbol{\Sigma}_{N} = (\beta \sum_{n=1}^{N} x_{n} x_{n}^{\top} + \lambda \mathbf{I}_{D})^{-1} = (\beta \mathbf{X}^{\top} \mathbf{X} + \lambda \mathbf{I}_{D})^{-1}$
(posterior's covariance matrix)

$$\mu_{N} = \boldsymbol{\Sigma}_{N} \begin{bmatrix} \beta \sum_{n=1}^{N} y_{n} x_{n} \end{bmatrix} = \boldsymbol{\Sigma}_{N} \begin{bmatrix} \beta \mathbf{X}^{\top} \mathbf{y} \end{bmatrix} = (\mathbf{X}^{\top} \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I}_{D})^{-1} \mathbf{X}^{\top} \mathbf{y}$$
(posterior's mean)

$$\mu_{N} = \boldsymbol{\Sigma}_{N} \begin{bmatrix} \beta \sum_{n=1}^{N} y_{n} x_{n} \end{bmatrix} = \boldsymbol{\Sigma}_{N} \begin{bmatrix} \beta \mathbf{X}^{\top} \mathbf{y} \end{bmatrix} = (\mathbf{X}^{\top} \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I}_{D})^{-1} \mathbf{X}^{\top} \mathbf{y}$$
(posterior's mean)

The Posterior: A Visualization

- Assume a lin. reg. problem with true $w = [w_0, w_1], w_0 = -0.3, w_1 = 0.5$
- Assume data generated by a linear regression model $y = w_0 + w_1 x + "noise"$
 - Note: It's actually 1-D regression (w_0 is just a bias term), or 2-D reg. with feature [1, x]
- Figures below show the "data space" and posterior of \boldsymbol{w} for different number of observations (note: with no observations, the posterior = prior)





Posterior Predictive Distribution

• To get the prediction y_* for a new input x_* , we can compute its PPD

$$p(y_*|x_*, \mathbf{X}, \mathbf{y}, \beta, \lambda) = \int p(y_*|x_*, \mathbf{w}, \beta) p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) d\mathbf{w} < \sum_{\substack{\text{posterior distribution so only}\\ \mathcal{N}(y_*|\mathbf{w}^{\mathsf{T}}x_*, \beta^{-1})} \mathcal{N}(w|\mu_N, \Sigma_N)$$

• The above is the marginalization of \boldsymbol{w} from $\mathcal{N}(y_*|\boldsymbol{w}^{\mathsf{T}}\boldsymbol{x}_*, \beta^{-1})$. Using LGM results

$$\mathcal{D}(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}, \beta, \lambda) = \mathcal{N}(\boldsymbol{\mu}_N^\top \mathbf{x}_*, \beta^{-1} + \mathbf{x}_*^\top \mathbf{\Sigma}_N \mathbf{x}_*)$$
 Can also derive it by writing $y_* = \mathbf{w}^\top \mathbf{x}_* + \epsilon$
where $\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N)$ and $\epsilon \sim \mathcal{N}(0, \beta^{-1})$

- So we have a predictive mean $\mu_N^T x_*$ as well as an input-specific predictive variance
- In contrast, MLE and MAP make "plug-in" predictions (using the point estimate of \boldsymbol{w})

$$p(y_*|x_*, w_{MLE}) = \mathcal{N}(w_{MLE}^\top x_*, \beta^{-1}) - \text{MLE prediction}$$

Since PPD also takes into account the uncertainty in w , $p(y_*|x_*, w_{MAP}) = \mathcal{N}(w_{MAP}^\top x_*, \beta^{-1}) - \text{MAP prediction}$
Since PPD also takes into the uncertainty in w , the predictive variance is larger

• Unlike MLE/MAP, variance of y_* also depends on the input \boldsymbol{x}_* (this, as we will see later, will be very useful in sequential decision-making problems such as active learning?

 $= w^{\mathsf{T}} x_* + \epsilon$

Posterior Predictive Distribution: An Illustration

Black dots are training examples



- Width of the shaded region at any x denotes the predictive uncertainty at that x (+/- one std-dev)
- Regions with more training examples have smaller predictive variance

CS772A: PML

10

Nonlinear Regression



Can extend the linear regression model to handle nonlinear regression problems

• One way is to replace the feature vectors \boldsymbol{x} by a nonlinear mapping $\boldsymbol{\phi}(\boldsymbol{x})$

$$p(y|\boldsymbol{x}, \boldsymbol{w}) = \mathcal{N}(\boldsymbol{w}^{ op} \phi(\boldsymbol{x}), \beta^{-1})$$

Can be pre-defined (e.g., replace a scalar x by polynomial mapping $[1, x, x^2]$) or extracted by a pretrained deep neural net

- Alternatively, a kernel function can be used to <u>implicitly</u> define the nonlinear mapping
- More on nonlinear regression when we discuss Gaussian Processes

Estimating Hyperparameters via MLE-II

• The probabilistic linear req. model we saw had two hyperparams (β, λ)



For regression with Gaussian likelihood

and Gaussian prior on \boldsymbol{w} , the marginal

likelihood has an exact expression

17

methods like EM, variational

CS772A: PML

inference, MCMC, etc later

And then compute $p(\boldsymbol{w}|\boldsymbol{X},\boldsymbol{y},\hat{\boldsymbol{\beta}},\hat{\boldsymbol{\lambda}})$ treating $\hat{\beta}, \hat{\lambda}$ as given

Prob. Linear Regression: Some Other Variations

13

• Can use other likelihoods $p(y_n | \boldsymbol{x}_n, \boldsymbol{w})$ and/or prior distribution $p(\boldsymbol{w})$



Logistic Regression

There are other ways too that can convert the score into a probability, such as a CDF: $p(y = 1 | x, w) = \mu = \Phi(w^T x)$ where Φ is the CDF of $\mathcal{N}(0,1)$. This model is known as "Probit Regression".



- A discriminative model for binary classification $(y \in \{0,1\})$
- A linear model with parameters $w \in \mathbb{R}^D$ computes a score $w^\top x$ for input x
- ullet A sigmoid function maps this real-valued score into probability of label being $oldsymbol{1}$

Also used as a nonlinear "activation function" in deep neural networks

$$p(y = 1 | \boldsymbol{x}, \boldsymbol{w}) = \mu = \sigma(\boldsymbol{w}^{\mathsf{T}} \boldsymbol{x})$$



- Thus conditional distribution of label $y \in \{0,1\}$ given x is the following Bernoulli

Likelihood

$$p(y|\mathbf{x}, \mathbf{w}) = \text{Bernoulli}[y|\mu] = \mu^{y}(1-\mu)^{1-y} = \left[\frac{\exp(\mathbf{w}^{\mathsf{T}}\mathbf{x})}{1+\exp(\mathbf{w}^{\mathsf{T}}\mathbf{x})}\right]^{y} \left[\frac{1}{1+\exp(\mathbf{w}^{\mathsf{T}}\mathbf{x})}\right]^{1-y}$$

- NLL is the binary cross-entropy loss: $-[y_n \log \mu_n + (1 y_n) \log (1 \mu_n)]$
- NLL is convex in w. Can also use a prior $p(w|\lambda) = \mathcal{N}(w|0, \lambda^{-1}I)$ if interested in MAP or full posterior on w

Logistic Regression: MAP and Posterior

- The posterior will be $p(w|X, y) = \frac{p(w)p(y|X, w)}{p(y|X)} = \frac{p(w)\prod_{n=1}^{N} p(y_n|w, x_n)}{\int p(w)\prod_{n=1}^{N} p(y_n|w, x_n) dw}$ Bernoulli
- MAP estimation is easy. $-\log p(w|X, y)$ is convex for LR. Unique minima
 - Can use first or second order optimization with gradient and Hessian being

$$g = -\sum_{n=1}^{N} (y_n - \mu_n) \mathbf{x}_n + \lambda \mathbf{I} \mathbf{w} = \mathbf{X}^{\top} (\boldsymbol{\mu} - \mathbf{y}) + \lambda \mathbf{w} \quad (a \ D \times 1 \text{ vector})$$
$$\mathbf{H} = \sum_{n=1}^{N} \mu_n (1 - \mu_n) \mathbf{x}_n \mathbf{x}_n^{\top} + \lambda \mathbf{I} = \mathbf{X}^{\top} \mathbf{S} \mathbf{X} + \lambda \mathbf{I} \quad (a \ D \times D \text{ matrix})$$
$$\mu_n = \sigma(\mathbf{w}^{\top} \mathbf{x}_n)$$

- Full posterior is intractable because of non-conjugacy
 - A popular option is to use the Laplace's approximation (other methods like MCMC and variational inference can also be used; will see them later)
 CS772A: PML

Laplace's (or Gaussian) Approximation

Consider a posterior distribution that is intractable to compute

$$p(\theta|D) = \frac{p(D, \theta)}{p(D)} = \frac{p(D|\theta)p(\theta)}{p(D)}$$

Laplace approximation approximates the above using a Gaussian distribution





Properties of Laplace's Approximation

- Straightforward if posterior's derivatives (first/second) can be computed easily
- Expensive if parameter θ is very high dimensional $< \frac{1}{2}$ simpler models (e.g., logistic reg with a very large number of features
 - Reason: We need to compute and invert Hessian of size $D \times D$ (D is the # of params)

E.g., a deep neural network, or even in

CS772A: PML



- Used only when θ is a real-valued vector (because of Gaussian approximation)
- Note: Even if we have a <u>non-probabilistic</u> model (loss function + regularization), we can obtain an approx "posterior" for that model using the Laplace's approximation
 - Optima of the regularized loss function will be Gaussian's mean
 - Inverse of the second derivative of the regularized loss function will be covariance matrix

*Mixtures of Laplace Approximations for Improved Post-Hoc Uncertainty in Deep Learning (Eschenhagen et al, 2021)

LR: Posterior Predictive Distribution

The posterior predictive distribution can be computed as

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int p(y_* = 1 | \mathbf{w}, \mathbf{x}_*) p(\mathbf{w} | \mathbf{X}, \mathbf{y}) d\mathbf{w}$$

Integral not tractable and must be approximated sigmoid Gaussian (if using Laplace approx.)

- Monte-Carlo approximation of this integral is one possible way
 - Draw M samples w_1, w_2, \dots, w_M , from the approx. of posterior
 - Approximate the PPD as follows

$$p(y_* = 1 | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) \approx \frac{1}{M} \sum_{m=1}^{M} p(y_* = 1 | \mathbf{w}_m, \mathbf{x}_*) = \frac{1}{M} \sum_{m=1}^{M} \sigma(\mathbf{w}_m^{\mathsf{T}} \mathbf{x}_n)$$

In contrast, when using MLE/MAP solution \widehat{w}_{opt} , the plug-in pred. distribution

$$p(y_* = 1 | \boldsymbol{x}_*, \boldsymbol{X}, \boldsymbol{y}) = \int p(y_* = 1 | \boldsymbol{w}, \boldsymbol{x}_*) p(\boldsymbol{w} | \boldsymbol{X}, \boldsymbol{y}) d\boldsymbol{w}$$

$$\approx p(y_* = 1 | \widehat{\boldsymbol{w}}_{opt}, \boldsymbol{x}_*) = \sigma(\widehat{\boldsymbol{w}}_{opt}^{\mathsf{T}} \boldsymbol{x}_n)$$

LR: Plug-in Prediction vs Bayesian Averaging

- Plug-in prediction uses a single w (point est) to make prediction
- \blacksquare PPD does an averaging using all possible w 's from the posterior



Multiclass Logistic (a.k.a. Softmax) Regression

- Also called multinoulli/multinomial regression: Basically, LR for K > 2 classes
- In this case, $y_n \in \{1, 2, ..., K\}$ and label probabilities are defined as



- K weight vecs w_1, w_2, \dots, w_K (one per class), each D-dim, and $W = [w_1, w_2, \dots, w_K]$
- Each likelihood $p(y_n | x_n, W)$ is a multinoulli distribution. Therefore total likelihood

$$p(\mathbf{y}|\mathbf{X}, \mathbf{W}) = \prod_{n=1}^{N} \prod_{\ell=1}^{K} \mu_{n\ell}^{y_{n\ell}} - \frac{\text{Notation: } y_{n\ell} = 1 \text{ if true class of }}{x_n \text{ is } \ell \text{ and } y_{n\ell'} = 0 \forall \ell' \neq \ell}$$

Can do MLE/MAP/fully Bayesian estimation for W similar to LR model

21

Generalized Linear Models

- (Probabilistic) Linear Regression: when response y is real-valued $p(y|x,w) = \mathcal{N}(y|w^{\mathsf{T}}x,\beta^{-1})$
- Logistic Regression: when response y is binary (0/1)

 $p(y|\mathbf{x}, \mathbf{w}) = \text{Bernoulli}[y|\sigma(\mathbf{w}^{\mathsf{T}}\mathbf{x})] = \left[\frac{\exp(\mathbf{w}^{\mathsf{T}}\mathbf{x})}{1 + \exp(\mathbf{w}^{\mathsf{T}}\mathbf{x})}\right]^{y} \left[\frac{1}{1 + \exp(\mathbf{w}^{\mathsf{T}}\mathbf{x})}\right]^{1-y}$

- Both are examples of a Generalized Linear Model (GLM)
 - The model depends on the inputs \boldsymbol{x} via a linear model $\boldsymbol{w}^{\mathsf{T}}\boldsymbol{x}$
- GLM is defined using an exponential family distribution $p(y|\mathbf{x}, \mathbf{w}) = \text{ExpFam}[y|f(\mathbf{w}^{\mathsf{T}}\mathbf{x})]^{\mathsf{L}}$

MLE/MAP of *w* is easy for GLMs (due to convex objective, thanks to exp-family). Posterior usually requires approximations if likelihood and prior are not conjugate pairs (Laplace approximation or other methods used)

- ExpFam can be any suitable distribution depending on the nature of outputs, e.g.,
 Gaussian for reals, Bernoulli for binary, Poisson for Count, gamma for positive reals
- ExpFam distributions are more generally useful in other contexts as well